

# Guiding human gaze with convolutional neural networks

Leon A. Gatys    Matthias Kümmerer    Thomas S. A. Wallis    Matthias Bethge  
University of Tübingen

{leon.gatys, matthias.kuemmerer, tom.wallis, matthias.bethge}@bethgelab.org

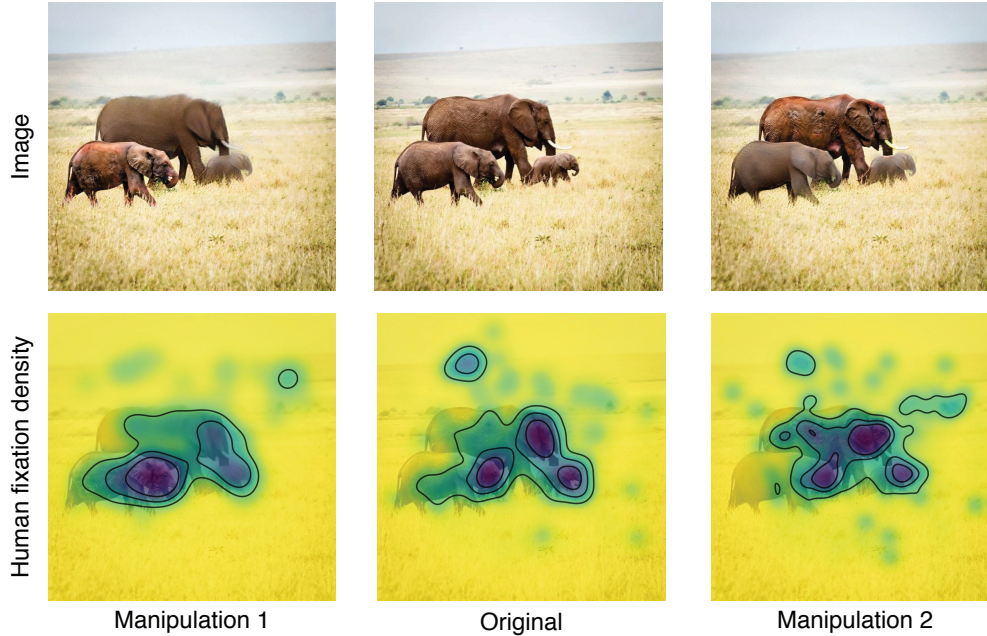


Figure 1: Manipulating images to change where people look. Images in the first and third column were generated by our model by transforming the original image shown in the second column. The presented images are shown in the first row and the measured human fixation densities are shown in the second row. The first manipulation aimed to make people fixate on the elephant in the foreground. The second manipulation aimed to make people fixate on the large elephant in the background. Compared to the original image, the total probability to fixate on the target elephant increased by 0.17 (86%), in the first image and 0.1 (28%) in the second image.

## Abstract

The eye fixation patterns of human observers are a fundamental indicator of the aspects of an image to which humans attend. Thus, manipulating fixation patterns to guide human attention is an exciting challenge in digital image processing. Here, we present a new model for manipulating images to change the distribution of human fixations in a controlled fashion. We use the state-of-the-art model for fixation prediction to train a convolutional neural network to transform images so that they satisfy a given fixation distribution. For network training, we carefully design a loss function to achieve a perceptual effect while preserving naturalness of the transformed images. Finally, we evaluate the success of our model by measuring human fixations for

a set of manipulated images. On our test images we can in-/decrease the probability to fixate on selected objects on average by 43/22% but show that the effectiveness of the model depends on the semantic content of the manipulated images.<sup>1</sup>

## 1. Introduction

Humans typically guide their visual attention selectively to different parts of an image and the spatial distribution of visual attention over an image strongly shapes perception. Since human vision is foveated, the most important measurable correlate of visual attention is the position at which the fovea is fixated.

<sup>1</sup>Supplement at: [bethgelab.org/media/uploads/gazeguide/Supplement.zip](http://bethgelab.org/media/uploads/gazeguide/Supplement.zip)

Recently, models based on features from convolutional neural networks (CNNs) trained on object recognition have lead to major advances in predicting human fixation locations, also called saliency prediction [22, 23, 20, 14]. Furthermore, the same feature spaces allow to generate and manipulate images with respect to important perceptual properties such as objects [33], text [34], image texture [8] or artistic style [9]. These image manipulations on perceptual variables are achieved by optimising perceptual loss functions defined in neural representations of CNNs trained on object recognition.

In this work, we aim to combine the recent success in predicting human fixation patterns with the advances in CNN-based image manipulation. We explore to what extent we can use the information captured by a CNN trained on fixation prediction to inform image manipulations that change the distribution of human fixations (also called *saliency map* [21, 24]) of images in a controlled way.

This problem is challenging for two main reasons. First, the saliency map of an image contains far less information than the image itself. Thus, there exists a myriad of images that would satisfy a given target saliency map, many of which are quite unnatural. Second, we find that the state-of-the-art model for human fixation prediction [23] is prone to adversarial examples, similar to other CNN-based prediction models [43, 5]. That means, one can apply changes to an image that are imperceptible for humans but make the image satisfy an arbitrary target saliency map for the prediction model. We address these challenges by carefully designing a loss function to preserve identity and naturalness of the transformed images. Moreover, we reduce the flexibility of the admissible image transformation by training a fixed CNN architecture to manipulate the saliency map of a large dataset of images such that it cannot overfit on adversarial noise for specific examples

Finally, we conduct a behavioural experiment to evaluate the success of our method. We construct a set of manipulated images and measure human fixation patterns in response to these images. Analysing the results, we find that the effectiveness of our method depends on the semantic content of the image but in most cases we successfully change the fixation patterns in the desired fashion.

## 2. Related work

Several previous studies aim to manipulate images with respect to saliency (for review see [30]). However, the nature of the image manipulations as well as the measure of image saliency varies. Image saliency is typically measured by some hand-crafted model (e.g. [16]) and the image manipulations aim to directly change the features that are used for saliency prediction (e.g. colour [31, 35], frequency bands [42] and luminance-contrast [7, 47] or a mixture of them [48, 12]). There are often some additional constraints

to preserve naturalness of the output image (e.g. total limits on the manipulation [48, 42] or only to use colours from within the image [32] or of similar objects from a database [35]). In contrast to our work, none of the previous studies aimed to learn a general saliency-manipulating image transformation directly from a data-driven state-of-the-art model for saliency prediction.

Similar to us, a number of studies have used pre-trained CNN features to synthesise or manipulate images with regard to perceptual properties (for review see [10]). Successful examples include attribute-based image synthesis [40, 33, 34] and manipulation [46], texture synthesis and style transfer [8, 9]. Technically most closely related to our work are studies that train a CNN to transform one image into another image. Example tasks include image generation from segmentations [15, 39, 3], style transfer [17, 44, 26, 50] or superresolution [17, 25, 38]. Note that in contrast to many other image to image translation tasks (e.g. [15]), there exists no ground truth data in pixel space for our saliency manipulation task.

## 3. Saliency manipulation method

The basic problem we are addressing can be formulated as follows: We want to find a transformation  $\mathcal{T}$  that maps an image  $I(x, y)$  and a target fixation distribution  $p_t(x, y)$  to a new image  $I_t(x, y)$  such that, when observing  $I_t(x, y)$ , human fixation patterns satisfy the specified target distribution  $p_t(x, y)$  (domain labelled ‘Human perception’ in Fig. 2).

However, since collecting human fixation patterns is expensive, we cannot directly guide our search for an appropriate image transformation by human behavioural data. Instead, we need to use a model that predicts human fixations and can be easily evaluated on new images.

### 3.1. Predicting human fixation patterns

To model the density of human fixations, we use the most recent DeepGaze model [23], a saliency prediction model that achieves state-of-the-art performance in fixation prediction as evaluated on the MIT300 saliency benchmark [2]. The model takes an image as input and outputs a probability distribution over pixels indicating the fixation probability at each image location:

$$DG(I(x, y)) = p(x, y) \quad (1)$$

To predict fixation densities, the model uses the feature spaces of the VGG-19 Network [41] that was trained on the ImageNet object recognition challenge [37]. In particular, it uses features from layers conv5\_1, relu5\_1, relu5\_2, conv5\_3, relu5\_4 giving a three-dimensional tensor with 2560 ( $5 \times 512$ ) channels. It computes a point-wise non-linear combination of the VGG features using a four-layer

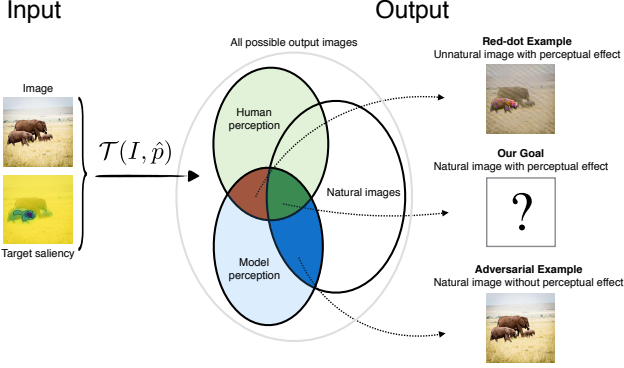


Figure 2: Conceptual setup: An input image and a target saliency map are transformed into an output image. Output images can have the desired saliency effect for humans (light green area, ‘Human perception’) and for the saliency model (light blue area, ‘Model perception’) and can be natural or artificial. The transformation is guided by model perception but our goal is to produce natural images that affect human perception in the desired way (dark green area). Artificial images can easily achieve desired saliency effects for humans (red area). Adversarial examples can easily achieve desired saliency effects for the model while maintaining naturalness (dark blue area).

readout network with  $1 \times 1$  rectified convolutions in each layer to produce a single output channel. This channel is up-sampled by a factor 8 and blurred with a gaussian kernel to regularise the predictions. Finally, a center-bias is added to the channel to model the prior distribution over fixations. The result,  $S(x, y)$ , is converted into a probability distribution over the image using a soft-max function over spatial positions:

$$p(x, y) = \frac{\exp(S(x, y))}{\sum_{x, y} \exp(S(x, y))} \quad (2)$$

Here, we omit the center-bias of DeepGaze when using it for saliency manipulation, since we want to inform our image transformations only with image dependent saliency information.

### 3.2. Saliency loss

To manipulate the saliency map of an image  $I$  to match a target saliency map  $p_t$ , we transform it by  $\mathcal{T}$  to generate a new image  $\hat{I}$ :  $\hat{I} = \mathcal{T}(I, p)$ . Next, we compute the saliency map  $\hat{p}$  of the transformed image:  $\hat{p} = DG(\hat{I})$ . To measure the success of the image transformation, we compute the KL-Divergence from the saliency map of the transformed image to the target saliency map:

$$\mathcal{L}_{sal} = \sum_{x, y} p_t(x, y) \log \left( \frac{p_t(x, y)}{\hat{p}(x, y)} \right) \quad (3)$$

Importantly, we are limited how well we can manipulate where people look by the agreement of our saliency model

with human fixations. In fact, we can only directly search for images that affect the perception of DeepGaze (domain labelled ‘Model perception’ in Fig. 2) but not for images that affect human perception. Still, we want to find images that not only affect the model perception but also human perception (intersection between ‘Model perception’ and ‘Human perception’ domains in Fig. 2).

### 3.3. Preserve naturalness

Manipulating the saliency map of images in any way is not necessarily useful. For example, guiding the observer’s attention by placing a bright red dot in the image would hardly be considered an interesting image manipulation. Similarly, there exist many transformations of the input image that strongly distort the image to match the target saliency map. These images can produce a perceptual effect for both the saliency model and humans, but lie outside the domain of natural images (red area of intersection between ‘Model perception’ and ‘Human perception’ domains Fig. 2) and thus the transformations that produce them are of limited use. Here we employ several measures to ensure that we are searching for image transformations that stay in the domain of natural images and somewhat close to the input image in particular.

#### 3.3.1 Feature loss

We aim to preserve the overall structure and content of the input image but leave flexibility for identity-preserving image transformations. To that end, we penalise the distance of the transformed image to the input image in a feature space provided by a deep layer of the VGG network. Say  $\mathbf{F}_\ell(I)$  is the feature representation of an image  $I$  in layer  $\ell$  of the VGG-network. Each column of  $\mathbf{F}_\ell(I)$  is a vectorised feature map and thus  $\mathbf{F}_\ell \in \mathbb{R}^{M_\ell(I) \times N_\ell}$  where  $N_\ell$  is the number of feature maps in layer  $\ell$  and  $M_\ell(I) = H_\ell(I) \times W_\ell(I)$  is the product of height and width of each feature map. We measure the mean-squared error between the feature representation of the input image  $I$  and the transformed image  $\hat{I}$ :

$$\mathcal{L}_{feat} = \frac{1}{N_\ell M_\ell(I)} \sum_{ij} \left( \mathbf{F}_\ell(\hat{I}) - \mathbf{F}_\ell(I) \right)_{ij}^2 \quad (4)$$

This is the same as the well-known content loss from Neural Style Transfer [9], only that here we use layer relu5\_2 instead of layer relu4\_2. This choice worked well in our experiments but we did not exhaustively compare between different choices of feature representations.

#### 3.3.2 Texture loss

We also want to preserve the overall appearance of the input image and its low-level structure without enforcing its exact

reconstruction. For that purpose, we employ a texture loss [8] that measures the difference between the texture of the transformed image and the texture of the output image:

$$\mathcal{L}_{tex} = \sum_{\ell} w_{\ell} E_{\ell} \quad (5)$$

$$E_{\ell} = \frac{1}{N_{\ell}^2} \sum_{ij} \left( \mathbf{G}_{\ell}(\hat{I}) - \mathbf{G}_{\ell}(I) \right)_{ij}^2 \quad (6)$$

where  $\mathbf{G}_{\ell}(I) = \frac{1}{M_{\ell}(I)} \mathbf{F}_{\ell}(I)^T \mathbf{F}_{\ell}(I)$  is the Gram Matrix of the feature maps in layer  $\ell$  in response to image  $I$ . We include Gram Matrices from layers relu1\_1, relu2\_1, relu3\_1, relu4\_1, relu5\_1 with equal weights to model the texture of the input image. For both the texture and the feature loss, we used the same normalised VGG-network [8] that is also used by DeepGaze [23].

### 3.3.3 Adversarial loss

Finally, we employ a patch-based conditional adversarial loss [15, 39, 26]. An adversarial loss [11] is an adaptive loss term aiming to correct for systematic differences between the input and the transformed images. In the adversarial loss, a discriminator learns to discriminate between image patches of transformed and input images and is jointly optimised with the image transformation. If the image transformation produces images whose patches are systematically different from the input images, the discriminator can learn this difference and inform the image transformation to correct for it.

For the discriminator we used the implementation of the patch-based adversarial loss by [15] with receptive field size of 70 px. We used the LSGAN [29] objective for improved training stability and thus, the following term is added to the loss function for the image transformation:

$$\mathcal{L}_{adv} = \sum_{x,y} (1 - D(\hat{I})(x, y))^2 \quad (7)$$

where  $D(I)(x, y)$  denotes the single channel output of the discriminator in response to image  $I$ . At the same time, the discriminator CNN is optimised to distinguish between input and transformed images:

$$\mathcal{L}_D = \frac{1}{2} \sum_{x,y} \left( D(\hat{I})(x, y)^2 + (1 - D(I)(x, y))^2 \right) \quad (8)$$

This patch-based conditional adversarial loss can also be understood as a texture loss. However, in our previously described texture loss the loss function is fixed and based on the pre-trained VGG features, whereas the adversarial loss function is adaptive and trained from scratch.

### 3.3.4 Avoid adversarial transformations

The total loss function we aim to minimise with respect to the image transformation is:

$$\mathcal{L}_{total} = \lambda_s \mathcal{L}_{sal} + \lambda_f \mathcal{L}_{feat} + \lambda_t \mathcal{L}_{tex} + \lambda_a \mathcal{L}_{adv} \quad (9)$$

For a given input image and target saliency map, the most flexible approach to minimise this loss function is to directly optimise the pixels of the input image, as was previously done for CNN-based image generation [40, 28, 8, 49]. Unfortunately though, we find that optimising the image directly leads to adversarial examples [43]: transformed images that are indistinguishable from the input images for humans, but match the required target saliency map predicted by DeepGaze. Thus, they remain natural images and generate a perceptual effect for the model but do not generate a perceptual effect for humans (dark blue area of intersection between ‘Model perception’ and ‘Natural images’ domains Fig. 2)

Adversarial examples can exist for any model whose prediction is based on different image information than human perception. Since our saliency prediction model is not perfect, the existence of adversarial examples is not surprising. Because the saliency loss is the only term that encourages the transformed image to be different from the input image, an image transformation that generates adversarial examples for DeepGaze will minimise our total loss function (Eq. 9). Thus, to avoid adversarial image transformations, we need to constrain the class of image transformations over which we optimise.

### 3.4. Image Transformation

We define the class of admissible image transformations by parameterising  $\mathcal{T}$  as a CNN with parameters  $\theta$ . Furthermore, we require that the same transformation simultaneously minimises our loss function for a large set of images (Fig. 3(a)). Hence, searching for the optimal image transformation  $\mathcal{T}_{opt}(I, p_t, \theta_{opt})$  means training a CNN to minimise the loss function (Eq. 9) for many input images  $I$  and target saliency maps  $p_t$ :

$$\theta_{opt} = \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{I, p_t} [\mathcal{L}_{total}(I, p_t, \theta)] \quad (10)$$

During training, input images and target saliencies are passed to the transformer network, which produces output images. The output image is passed to DeepGaze. The VGG features of DeepGaze are used to compute the feature and texture losses. The final output of DeepGaze is used to compute the saliency loss (Fig. 3(a)). At the same time, the input and the transformed image are passed to the discriminator network. The discriminator computes the adversarial loss for the transformer network (Eq. 7) and optimises its own discrimination target (Eq. 8) (Fig. 3(a)).



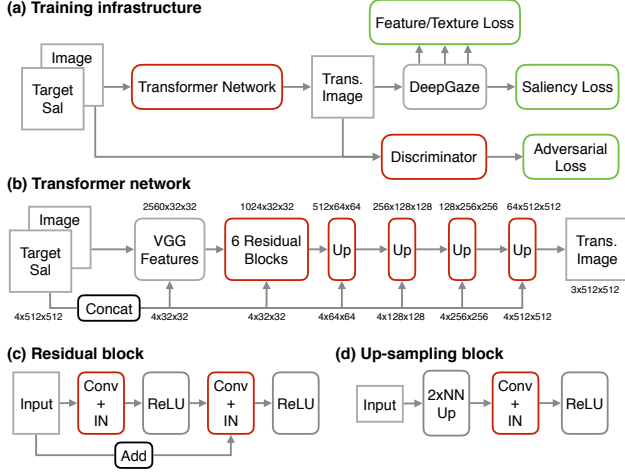


Figure 3: **(a)** Training infrastructure. Modules with trainable parameters are red, contributions to the loss function are green. Grey modules are not optimised. **(b)** Network architecture. VGG features are computed from the input image and in subsequent stages transformed and up-sampled to generate the transformed image. The input image and target saliency map is concatenated with the feature maps at several stages of the transformer network. **(c)** Residual block: Two convolutional layers with InstanceNorm (IN) and ReLU non-linearities. The input is added before the second ReLU. **(d)** Up-sampling block: Nearest-neighbour up-sampling followed by a convolutional layer with IN and ReLU.

### 3.4.1 Transformer network

Since existing CNN architectures for image-to-image mapping [17, 15] did not perform well on our task (Fig. 4), we developed a new CNN architecture (Fig. 3(b)). First, the input image is transformed into the 2560 VGG feature channels used by DeepGaze to predict the saliency map. This feature space extracts the rich image information with large receptive fields that DeepGaze uses for saliency prediction. Its spatial dimensions are 16 times smaller than that of the input image. At that resolution, the image is transformed using 6 blocks of residual layers [17, 13] that have 1024 feature channels. Each residual block consists of two stages of convolutional layer with kernel size  $3 \times 3$  (Conv), InstanceNorm (IN) [45] and rectifying-linear unit (ReLU). The input to the residual block is added before the last ReLU (Fig. 3(c)). Next, we have 4 up-sampling stages that increase each spatial dimension by a factor of 2 and decrease the number of channels by a factor of 2 (Fig. 3(b)). Each up-sampling stage consists of nearest neighbour up-sampling (NN Up), Conv, IN and ReLU (Fig. 3(d)). At each stage of the processing hierarchy, we want the transformer network to have access to the input image and target saliency map to inform the transformation and allow the preservation of low-level information. Therefore, we down-sample the 4 input channels (RGB image and target saliency map) to the respective size and concatenate them with the feature output

of the network at each processing stage (Fig. 3(b)). Thus, the first residual block receives 2564 (2560+4), the first up-sampling block 1028 (1024+4) and the second up-sampling block 516 (512+4) input channels and so on. In that way we have a powerful transformer architecture built on the rich object-based VGG features that can also preserve all low-level image information. In the following we call this architecture ‘feature-guided transform’ (FGTransform).

### 3.4.2 Network training

We use the MSCOCO training images [27] to train our network on saliency manipulation. The images are spatially resized to 512x512 pixels, and pre-processed for the VGG-network (transformed to BGR, subtraction of channel mean, and scaled by 255) [41].

We have to generate target saliency maps for each training sample. We hypothesise that there is a class of ‘natural saliency maps’ that arise from the set of natural images. We tried to construct natural target saliency maps by changing the original saliency map of the input image in either of two ways.

In the first manipulation, we add a constant to a local region of the un-normalised saliency map:

$$S_t(x, y) = S(x, y) + k_{sh}M(x, y) \quad (11)$$

Here,  $k_{sh}$  is the constant shift to the saliency map and  $M(x, y)$  denotes a mask that defines the local region which will be changed. Intuitively, this manipulation corresponds to increasing (for  $k_{sh} > 0$ ) or decreasing (for  $k_{sh} < 0$ ) the saliency of local parts of the image. We sample the masks  $M(x, y)$  from the object segmentation labels of the COCO dataset. Each mask is blurred with a gaussian kernel to keep the target saliency maps smooth. Thus, during training the network learns to in-/decrease the saliency of annotated objects and people in the training set (Fig. 4(a), third column).

For the second manipulation, we globally scale the saliency map with a constant factor:

$$S_t(x, y) = k_{sc} \times S(x, y) \quad (12)$$

Here  $k_{sc}$  denotes the constant scaling factor that is applied to the saliency map. Intuitively, this manipulation changes the consistency of the fixation patterns. It either increases the clustering of the fixations on few, very salient regions (for  $k_{sc} > 1$ ) or spreads the fixation patterns more uniformly over the image (for  $0 < k_{sc} < 1$ ). The factor  $k_{sc}$  can also be thought of a temperature parameter that in-/decreases the entropy of the fixation distribution (Fig. 4(a), fourth column).

In both cases, the target saliency distribution  $p_t$  is obtained by normalising  $S_t$  using the soft-max function (Eq. 2). During training, we generate new target saliency maps

on-the-fly for each training sample by either locally shifting all annotated objects in the image by a constant  $k_{sh} \in [-4, 4]$  or globally scaling by a factor  $k_{sc} \in [0.5, 2]$ . We always used the log saliency map during training, so the input to the transformer network is  $\{I, \log(p_t)\}$ .

Further details of the training procedure can be found in the Appendix.

## 4. Results

### 4.1. Visual inspection

We compare the training results between our FGTransform and the previously published ‘Resnet-9’ and ‘U-Net’ architectures [17, 15] (Fig. 4(b),(c)). For the ‘Resnet-9’ and ‘U-Net’ architectures, we observed problems with unnatural distortions in the transformed images (Fig. 4(b), head of the salient dog in second and third column). Interestingly, these distortions appear to resemble scribbled text, which is one of the main features driving human fixations [23]. Thus, it makes sense to put text at locations that should have high saliency, but this does not preserve the naturalness of the image well. Ideally though, the transformation can increase the saliency of the features in the input image rather than always putting text scribbles in the corresponding location. That is why we designed our own feature-guided architecture hoping that a better-suited network can learn a more image-dependent transformation. We did not encounter similar problems with artificial distortions with our architecture (Fig. 4(b), first column) indicating that it can learn a more image-dependent transformation. Model comparisons for all stimuli used in the behavioural study in section 5 can be found in the Supplement.

### 4.2. Quantitative evaluation

To compare the models quantitatively, we computed the training and test error over the training iterations for each saliency manipulation. We sampled 1000 training images from the COCO training set and 1000 test images from the COCO validation set that were not used during training of the model. As during training, we randomly sampled either  $k_{sh}$  from  $[-4, 4]$  or  $k_{sc}$  from  $[0.5, 2]$  to generate the target saliency for each image. We find that the training and test error are on a similar level and even after 150k training iterations there is no sign of overfitting (Fig. 4(c)). Also, when inspecting the transformed images we were unable to tell the difference between training and test images. Furthermore, our network architecture consistently leads to smaller loss values than the compared ‘Resnet-9’ and ‘U-Net’ architectures (Fig. 4(c)). This is true for the saliency loss as well as the regularisation losses meaning that our model provides a better solution to the problem and not only a better trade-off between naturalness and saliency manipulation.

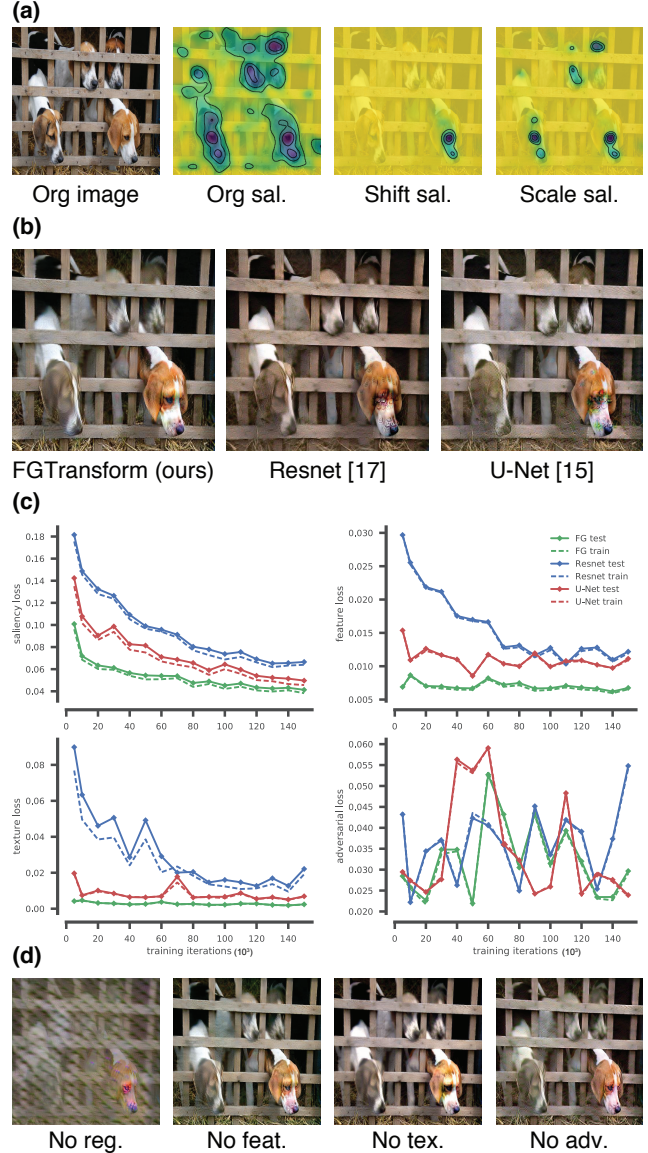


Figure 4: (a) Target saliency maps are either created by in-/decreasing the saliency of particular objects (third column) or globally scaling the saliency map (fourth column). (b) Existing image-to-image networks generate artificial distortions in our task. (c) Our feature-guided architecture minimises all parts of the loss function better than existing architectures. Results shown for training on local saliency shifts. Results for global saliency scaling look similar and can be found in the Appendix. (d) Ablation studies. Without regularisation output images are very distorted (first column). Leaving out other parts of the loss function has subtle but noticeable effects. Images best viewed with zoom.

### 4.3. Ablation studies

Finally, we re-trained our FGTransform while setting different parts of the loss function equal to zero. When only training on the saliency loss ( $\lambda_f, \lambda_t, \lambda_a = 0$ ), the

transformed images are strongly distorted (Fig. 4(d), first column). When leaving out only one of the regularisation losses (either of  $\lambda_f, \lambda_t, \lambda_a$  equal to 0), the differences are more subtle. We see that in each case, the images are slightly more distorted than for the full loss function 4(d)). Looking at many examples, we found that the addition of every part of the loss function increased the perceptual quality of the transformed images. Nevertheless, the perceptually optimal trade-off between the different regularisation terms is hard to determine since our loss function is only a rough quantitative measure of perceptual quality. Ablation results for all stimuli used in the behavioural study in section 5 can be found in the Supplement.

## 5. Behavioural study

We have developed an image transformation that can manipulate arbitrary input images to change the saliency prediction of DeepGaze while preserving naturalness. However, the existence of adversarial examples illustrates that model prediction and human perception can be quite different for images optimised with respect to the model. Thus, to evaluate our work it is vital to measure human fixation patterns in response to images generated with our model.

### 5.1. Stimulus generation

We picked 24 images from the COCO validation set that were not used during training. For each image, we created two modified versions leading to a total of 72 images. We designed the modifications to generate opposite changes in human behaviour compared to the original image. For 21 images, we aimed to in-/decrease the saliency of different objects in the two versions using a mixture of the local shifting and global scaling of the saliency map (e.g. Fig. 5(a),(b)). For three images, we purely scaled the saliency map to in-/decrease its entropy (e.g. Fig. 5(c)). Since our image transformation is fast (117 ms on a GTX 1080 GPU), we can manipulate images with respect to the parameters  $k_{sh}$  and  $k_{sc}$  in an online, interactive fashion.

When applying both local shifting and global scaling, we first transformed the image by the network trained on local shifting and afterwards transformed the output of that manipulation by the network trained on global scaling of the saliency map. In difference to the training, we did not blur the local shift mask when computing the target saliency as this slightly improved the perceptual quality of the manipulations. Images of all stimuli, the target saliency map to generate them, their saliency map predicted by DeepGaze and their saliency map measured from human behaviour can be found in the Supplement.

DeepGaze is trained on images that are down-sampled by a factor 2 compared to the size of the images on which the fixation data was collected [22]. Therefore we needed to up-sample the generated images by the same factor be-

fore collecting human fixations. We used a state-of-the art network for superresolution [38] to up-sample the generated and original images of size  $512 \times 512$  to size  $1024 \times 1024$  as this gave better results than bicubic up-sampling.

### 5.2. Experimental setup

We measured fixation responses of 23 subjects to our 72 stimuli in a free-viewing task. Experimental details can be found in the Appendix.

To obtain an empirical fixation density from the raw fixation data, we fit a kernel density estimate together with a uniform component and a center-bias. We fit this estimate separately for every image and cross-validate over subjects.

To compute DeepGaze’s prediction for each stimulus, we use the center-bias computed from the empirical densities from all other stimuli.

### 5.3. Results

Each stimulus was generated to produce a specific behavioural effect. For 21 source images, we generated two manipulated versions whose desired effect was to in-/decrease the probability of looking at certain objects (e.g. Fig. 5(a)). After obtaining empirical fixation densities from measuring people’s fixation patterns, we use the COCO segmentation masks to compute the probability of people looking at the objects in the image:  $p_{obj} = \sum_{x,y} M(x,y)p(x,y)$ . In case of manipulations that aimed to increase the probability of looking at certain objects, we measured an average increase in the probability of looking at objects targeted by the manipulation of 0.09. Relative to the probability of looking at the corresponding object in the original image this result constitutes an increase of 43%. For manipulations that aimed to decrease the probability of looking at certain objects, we measured an average decrease in the probability of looking at objects targeted by the manipulation of 0.04 (22%).

For three source images, we generated two manipulated versions that aimed to change the entropy of the fixation density (Fig. 5(c)). We compute the entropy of the empirical fixation densities as  $H = -\sum_{x,y} p(x,y) \log(p(x,y))$ . The fixation densities of images manipulated to in-/decrease the entropy of their saliency map showed an average in-/decrease in entropy by 0.19/0.48 bits (1/3%). Although the effect size for this manipulation seems small, it can considerably change the perception of an image (Fig. 5(c)).

We can also compare the empirical fixation densities measured from human behaviour with the prediction produced by our saliency model (Fig. 5). We find that DeepGaze usually predicts much stronger effects for the manipulated images than we find from human behaviour. DeepGaze predicts an in-/decrease in fixation probability for manipulated objects of 0.43/0.15 (142/74%). For the entropy manipulation, DeepGaze predicts an average in-



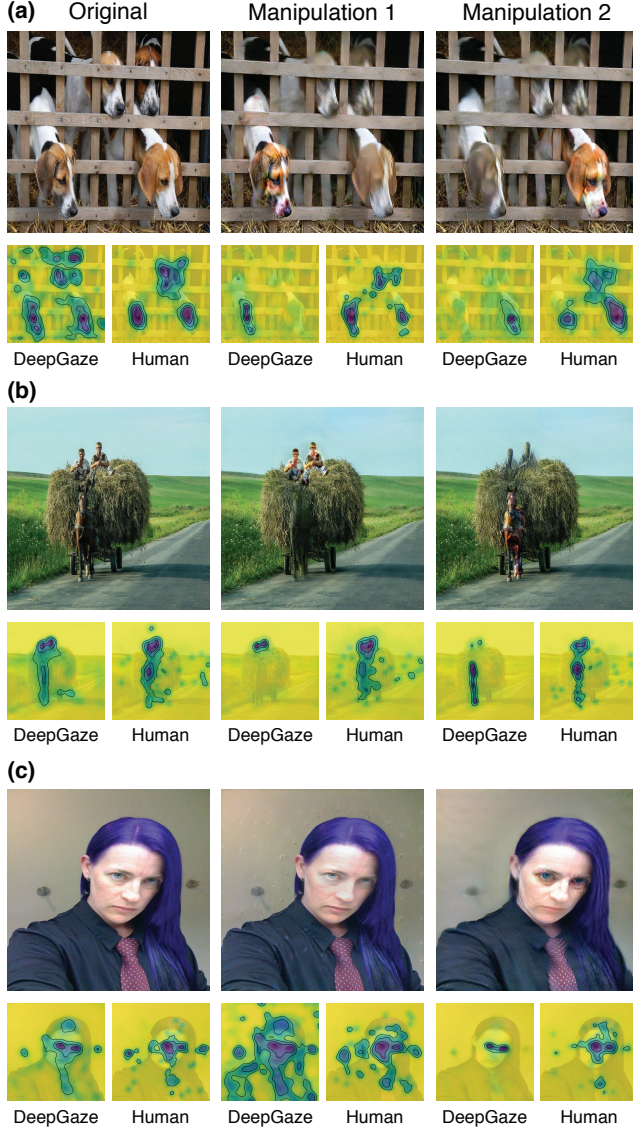


Figure 5: Example results of behavioural study. **(a)** The image is manipulated to highlight either the dog in the lower right or lower left of the image. Fixation probability on the targeted dog increases by 0.17 (77%) for the dog in the left and 0.17 (70%) for the dog in the right. Deep Gaze predicts stronger increases of 0.64 (285%) and 0.54 (274%) respectively. **(b)** The image is manipulated to highlight either two humans on the carriage or the horse pulling it. Fixation probability increases on average by 0.02 (13%) for the humans and 0.15 (66%) for the horse. Deep Gaze predicts stronger increases of 0.15 (105%) and 0.56 (262%) respectively. **(c)** The image is manipulated to in-/decrease the entropy of the fixation density. The empirical entropy in-/decreases by 0.76/0.18 bits. DeepGaze predicts a stronger in-/decrease of 1.08/1.04 bits. Images best viewed with zoom.

/decrease by 0.89/1.32 bits (5/8%). The discrepancy between model prediction and human behaviour shows that our image transformation still has an adversarial component

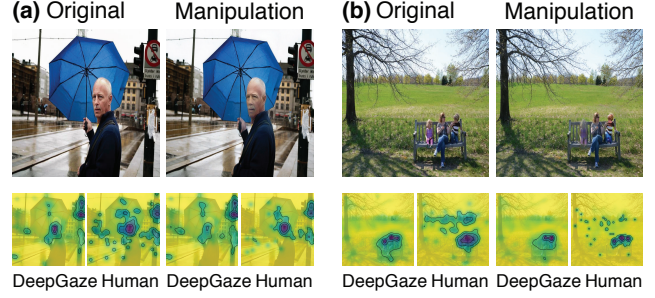


Figure 6: Failure examples. **(a)** The probability to fixate on the man increases by 0.07 (19%) although DeepGaze predicts a decrease by 0.32 (44%). **(b)** The probability to fixate on the girl increases by 0.07 (150%) although DeepGaze predicts a decrease by 0.1 (82%). Images best viewed with zoom.

to it: It is more effective in manipulating model perception compared to human perception.

Inspecting the results, we find that the adversarial component is most apparent when human fixations are driven by high-level semantic information. For example, humans continue to look at a face in the middle of the image, even though the image is manipulated such that DeepGaze assigns considerably less saliency to the face (Fig. 6(a)).

## 6. Discussion

The ability of machines to predict perceptual properties of images has greatly improved with the rise of CNNs. Still, it is an open question to what extent these models can inform image manipulation and synthesis with respect to perceptual properties. Prediction models discard image information that is not informative for their task leading to two major problems when using them for image generation: Ensuring that the output is sufficiently natural and the existence of adversarial examples.

In this work we tackle these problems in the specific case of saliency manipulation. We use the best model for fixation prediction to learn an image-to-image mapping that can manipulate images to change human fixation patterns in a controlled fashion. Nevertheless, we also find some limitations of our saliency manipulation model. Most prominently, the learned image transformation has difficulties to decrease the saliency of semantically important objects such as human faces. This can probably only be achieved with severe changes to the semantic content of the image (e.g. deleting a face in some way), which our model has not managed to learn. It is a compelling question for future work how to improve the design of the transformation network to enable such strong semantic image manipulations. In summary, we believe this work sets the stage for an exciting new path to edit images and contributes towards enabling the use of powerful prediction models for image manipulation.



## 7. Appendix

### 7.1. Training details

#### 7.1.1 Transformer network

We set the weights of the loss function to  $\lambda_s = 1$ ,  $\lambda_f = 1e-2$ ,  $\lambda_t = 2e-2$ ,  $\lambda_a = 1e-1$  after a preliminary exploration phase on smaller datasets. We trained separate network instance for each type of saliency manipulation, the local shifting and the global scaling. Every network was trained for 150k iterations with batch size 4 and learning rate  $1e-3$  using the Adam optimiser [18]. As a comparison to our FGTransform, we also trained instances of the ‘Resnet-9’ and ‘U-Net’ architectures from [17, 15] with two slight modifications: We replaced BatchNorm by InstanceNorm and we initialised them as an auto-encoder by first training them to reconstruct the input image, which slightly accelerated training in the beginning. We also experimented with the ‘context-aggregation network’ (CAN32) from [4] but did not get promising results in our task.

#### 7.1.2 Adversarial loss

The discriminator in the adversarial loss is a five-layer CNN with 64, 128, 256, 512 and 1 channels and LeakyReLU non-linearity. In the last layer, each unit has receptive field size of 70px and is trained to discriminate between transformed and real image patches. In difference to [15] we used InstanceNorm instead of BatchNorm layers in the discriminator. To implement the adversarial loss, we slightly modified the code<sup>2</sup> from [15, 50].

#### 7.1.3 Extended results

In Fig. 7 we show the training and test loss for the networks trained with target saliency maps from global scaling of the original saliency maps. The Ablation and Model Comparison results for all stimuli used in the behavioural experiment can be found in the ‘**ModelComparison**’ and ‘**AblationStudies**’ folders in the Supplement at: [bethgelab.org/media/uploads/gazeguide/Supplement.zip](http://bethgelab.org/media/uploads/gazeguide/Supplement.zip).

### 7.2. Behavioural study

23 participants were recruited from an internal mailing list. Each participant saw three blocks of 24 stimuli; each block contained one version of the 24 source images (one block with the original images and two blocks with manipulated images). The order of the blocks was counterbalanced over participants (latin square design) such that an approximately equal number of participants saw each condition in their first block. On each trial, participants were presented with a fixation target in the centre of the screen. After fixating, the images were displayed for 3 seconds and could be

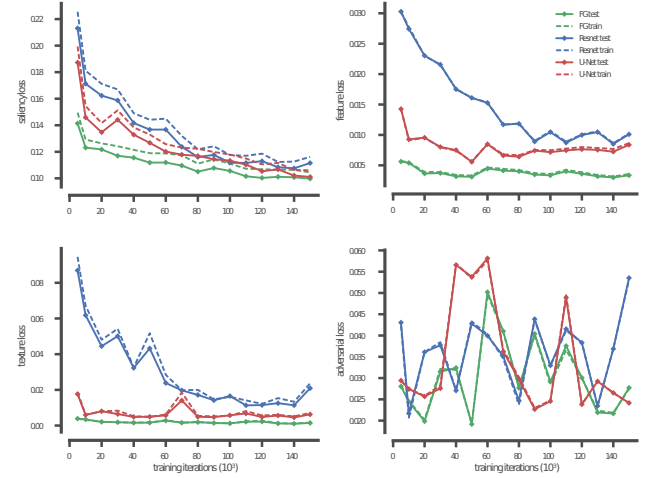


Figure 7: Quantitative evaluation of network trained with target saliency maps generated by global scaling of the original saliency map.

freely explored by the participants. The inter-trial interval was 2 seconds, in which a blank grey screen was presented.

Stimuli were displayed on a VIEWPixx 3D LCD (VIEWPixx Technologies Inc., Saint-Bruno-de-Montarville, Canada; spatial resolution 1920 1080 pixels, temporal resolution 120 Hz). Participants viewed the display from 60 cm in a darkened chamber. At this distance the images subtended approximately 25 degrees of angle at the retina. Stimuli were presented using the Psychtoolbox Library [1, 36, 19] version 3.0.12 under MATLAB (The Mathworks Inc., Natick MA, USA; R2015b). Participants’ gaze position was recorded monocularly (left eye) at 500 Hz with an Eyelink 1000 (SR Research, Ontario, Canada) video-based eyetracker in combination with the Eyelink toolbox [6] for MATLAB. Gaze traces were classified into fixations using the default settings of the SR Research processing software.

There are two potential problems in the data. First, every participant saw each image 3 times (once per block). Thus, memory effects potentially changed behaviour. We controlled for this by counterbalancing the order of the blocks over participants to minimise the influence of memory / familiarity on the average densities over conditions. We additionally analysed the data of only the first block for each participant and found a similar effect in human behaviour as for the full data (average in-/decrease to fixate the targeted objects: 0.08/0.03 (55/9%)).

Second, although all images except for one (the snowboarder) were unseen by all participants before the experiment, some participants were familiar with the general research question. Thus, fixations could be influenced by conscious behaviour. To check if this could potentially invalidate our results we also analysed only the first fixations

<sup>2</sup>[github.com/junyanz/pytorch-CycleGAN-and-pix2pix](https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix)

from all participants, under the assumption that the first fixation is difficult to control voluntarily. We found that for this subset of the data the effect size rather increased (average in-/decrease to fixate the targeted objects: 0.15/0.10 (3230/35%)). The reason for the very large average relative increase is that for some manipulated images, a significant amount of first fixations is guided to target objects that received no first fixation in the original image (e.g. see Fig. 8, fixations on the bicycle). This generates a huge relative increase in fixation probability for these images, which dominate the average value. The median relative in-/decrease in fixation probability for the first fixation only was 57/48%.

All data from the experiment is contained in the folder **‘BehaviouralExperiment’** in the Supplement at: [bethgelab.org/media/uploads/gazeguide/Supplement.zip](http://bethgelab.org/media/uploads/gazeguide/Supplement.zip):

- **‘BehaviouralExperiment/Stimuli’** contains the images shown in the experiment
- **‘BehaviouralExperiment/TargetSaliencyMaps’** contains the target saliency maps used to produce the stimuli. Note that the target saliency maps can look artificial, since we did not use blurring of the object masks as we did during training. Training directly with non-blurred object masks did not improve the results.
- **‘BehaviouralExperiment/DeepGazePrediction’** contains the saliency prediction by DeepGaze in response to the stimuli images. The difference between the target saliency and the saliency prediction is quantified by the test loss of the model (although the loss was again evaluated with blurred object masks).
- **‘BehaviouralExperiment/HumanFixations’** contains the human fixation data in response to the stimuli in the folder.
- **‘BehaviouralExperiment/DataOnlyFirstBlock’** contains the human fixations and DeepGaze predictions for only the data in the first block for each subject. Note that DeepGaze predictions are slightly different, because the center-bias is also adapted to contain only the data from the first block.
- **‘BehaviouralExperiment/DataOnlyFirstFixation’** contains the human fixations and DeepGaze predictions for only first fixation of all subjects and images. Note that DeepGaze predictions are slightly different, because the center-bias is also adapted to contain only the data from the first fixation.

## References

- [1] D. H. Brainard. The Psychophysics Toolbox. *Spatial Vision*, 10(4):433–436, 1997. 9

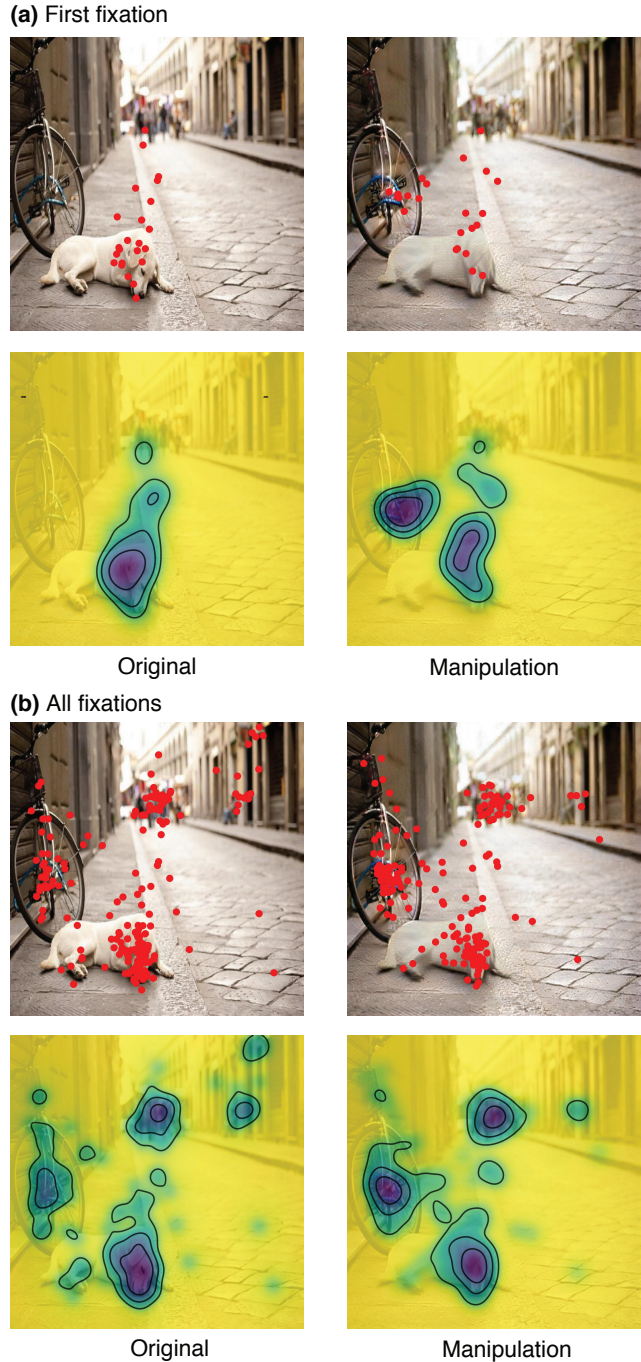


Figure 8: **(a)** Raw fixation data and estimated empirical density for only the first fixation for all subjects. There are no fixations on the bicycle on the left for the original but several for the manipulated image making the relative increase in fixation probability extreme (0.31/111648%). **(b)** Raw fixation data and estimated empirical density for all fixations from all subjects. Fixations are more distributed and the relative increase of fixation probability on the bicycle is therefore less extreme (0.12/70%).

- [2] Z. Bylinskii, T. Judd, A. Borji, L. Itti, F. Durand, A. Oliva, and A. Torralba. Mit saliency benchmark. [2](#)
- [3] Q. Chen and V. Koltun. Photographic image synthesis with cascaded refinement networks. *arXiv preprint arXiv:1707.09405*, 2017. [2](#)
- [4] Q. Chen, J. Xu, and V. Koltun. Fast image processing with fully-convolutional networks. In *IEEE International Conference on Computer Vision*, 2017. [9](#)
- [5] M. Cisse, Y. Adi, N. Neverova, and J. Keshet. Houdini: Fooling Deep Structured Prediction Models. *arXiv:1707.05373 [cs, stat]*, July 2017. arXiv: 1707.05373. [2](#)
- [6] F. W. Cornelissen, E. M. Peters, and J. Palmer. The Eye-link Toolbox: Eye tracking with MATLAB and the Psychophysics Toolbox. *Behavior Research Methods, Instruments, & Computers*, 34(4):613–617, Nov. 2002. [9](#)
- [7] W. Einhäuser and P. König. Does luminance-contrast contribute to a saliency map for overt visual attention? *European Journal of Neuroscience*, 17(5):1089–1097, Mar. 2003. [2](#)
- [8] L. Gatys, A. S. Ecker, and M. Bethge. Texture Synthesis Using Convolutional Neural Networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 262–270. Curran Associates, Inc., 2015. [2](#), [4](#)
- [9] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2414–2423, 2016. [2](#), [3](#)
- [10] L. A. Gatys, A. S. Ecker, and M. Bethge. Texture and art with deep neural networks. *Current Opinion in Neurobiology*, 46:178–186, Oct. 2017. [2](#)
- [11] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative Adversarial Nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014. [4](#)
- [12] A. Hagiwara, A. Sugimoto, and K. Kawamoto. Saliency-based Image Editing for Guiding Visual Attention. In *Proceedings of the 1st International Workshop on Pervasive Eye Tracking & Mobile Eye-based Interaction*, PETMEI ’11, pages 43–48, New York, NY, USA, 2011. ACM. [2](#)
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. pages 770–778, 2016. [5](#)
- [14] X. Huang, C. Shen, X. Boix, and Q. Zhao. Salicon: Reducing the semantic gap in saliency prediction by adapting deep neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 262–270, 2015. [2](#)
- [15] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *arXiv preprint arXiv:1611.07004*, 2016. [2](#), [4](#), [5](#), [6](#), [9](#)
- [16] L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(11):1254–1259, Nov. 1998. [2](#)
- [17] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual Losses for Real-Time Style Transfer and Super-Resolution. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *Computer Vision – ECCV 2016*, number 9906 in Lecture Notes in Computer Science, pages 694–711. Springer International Publishing, Oct. 2016. DOI: 10.1007/978-3-319-46475-6\_43. [2](#), [5](#), [6](#), [9](#)
- [18] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]*, Dec. 2014. arXiv: 1412.6980. [9](#)
- [19] M. Kleiner, D. H. Brainard, and D. G. Pelli. What’s new in Psychtoolbox-3? *Perception*, 36(ECVP Abstract Supplement), 2007. [9](#)
- [20] S. S. Kruthiventi, K. Ayush, and R. V. Babu. Deepfix: A fully convolutional neural network for predicting human eye fixations. *IEEE Transactions on Image Processing*, 2017. [2](#)
- [21] M. Kuemmerer, T. Wallis, and M. Bethge. Information-theoretic model comparison unifies saliency metrics. *Proceedings of the National Academy of Science*, 112(52):16054–16059, 2015. [2](#)
- [22] M. Kümmeler, L. Theis, and M. Bethge. Deep Gaze I: Boosting Saliency Prediction with Feature Maps Trained on ImageNet. In *ICLR Workshop*, 2015. [2](#), [7](#)
- [23] M. Kümmeler, T. S. Wallis, L. A. Gatys, and M. Bethge. Understanding Low- and High-Level Contributions to Fixation Prediction. In *The IEEE International Conference on Computer Vision (ICCV)*, 2017. [2](#), [4](#), [6](#)
- [24] M. Kümmeler, T. S. A. Wallis, and M. Bethge. Saliency Benchmarking: Separating Models, Maps and Metrics. In *arxiv*. 2017. [2](#)
- [25] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network. *arXiv:1609.04802 [cs, stat]*, Sept. 2016. arXiv: 1609.04802. [2](#)
- [26] C. Li and M. Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *European Conference on Computer Vision*, pages 702–716. Springer, 2016. [2](#), [4](#)
- [27] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: Common Objects in Context. In *Computer Vision – ECCV 2014*, Lecture Notes in Computer Science, pages 740–755. Springer, Cham, Sept. 2014. [5](#)
- [28] A. Mahendran and A. Vedaldi. Understanding Deep Image Representations by Inverting Them. *arXiv:1412.0035 [cs]*, Nov. 2014. arXiv: 1412.0035. [4](#)
- [29] X. Mao, Q. Li, H. Xie, R. Y. K. Lau, Z. Wang, and S. P. Smolley. Least Squares Generative Adversarial Networks. *arXiv:1611.04076 [cs]*, Nov. 2016. arXiv: 1611.04076. [4](#)
- [30] V. Mateescu and I. Bajic. Visual Attention Retargeting. *IEEE Multimedia*, 23, May 2015. [2](#)
- [31] V. A. Mateescu and I. V. Bajić. Attention Retargeting by Color Manipulation in Images. In *Proceedings of the 1st International Workshop on Perception Inspired Video Processing*, PIVP ’14, pages 15–20, New York, NY, USA, 2014. ACM. [2](#)
- [32] R. Mechrez, E. Shechtman, and L. Zelnik-Manor. Saliency Driven Image Manipulation. *arXiv:1612.02184 [cs]*, Dec. 2016. arXiv: 1612.02184. [2](#)



- [33] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. *arXiv:1605.09304 [cs]*, May 2016. arXiv: 1605.09304. 2
- [34] A. Nguyen, J. Yosinski, Y. Bengio, A. Dosovitskiy, and J. Clune. Plug & Play Generative Networks: Conditional Iterative Generation of Images in Latent Space. *arXiv:1612.00005 [cs]*, Nov. 2016. arXiv: 1612.00005. 2
- [35] T. V. Nguyen, B. Ni, H. Liu, W. Xia, J. Luo, M. Kankanhalli, and S. Yan. Image Re-Attentionizing. *IEEE Transactions on Multimedia*, 15(8):1910–1919, Dec. 2013. 2
- [36] D. G. Pelli. The VideoToolbox software for visual psychophysics: Transforming numbers into movies. *Spatial Vision*, 10(4):437–442, 1997. 9
- [37] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *arXiv:1409.0575 [cs]*, Sept. 2014. arXiv: 1409.0575. 2
- [38] M. S. M. Sajjadi, B. Schölkopf, and M. Hirsch. EnhanceNet: Single Image Super-Resolution through Automated Texture Synthesis. *arXiv:1612.07919 [cs]*, Dec. 2016. arXiv: 1612.07919. 2, 7
- [39] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. *arXiv preprint arXiv:1612.07828*, 2016. 2, 4
- [40] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. *arXiv:1312.6034 [cs]*, Dec. 2013. 2, 4
- [41] K. Simonyan and A. Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv:1409.1556 [cs]*, Sept. 2014. arXiv: 1409.1556. 2, 5
- [42] S. L. Su, F. Durand, and M. Agrawala. De-emphasis of Distracting Image Regions Using Texture Power Maps. In *Proceedings of the 2Nd Symposium on Applied Perception in Graphics and Visualization*, APGV '05, pages 164–164, New York, NY, USA, 2005. ACM. 2
- [43] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *arXiv:1312.6199 [cs]*, Dec. 2013. 2, 4
- [44] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky. Texture Networks: Feed-forward Synthesis of Textures and Stylized Images. *arXiv:1603.03417 [cs]*, Mar. 2016. arXiv: 1603.03417. 2
- [45] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Improved Texture Networks: Maximizing Quality and Diversity in Feed-forward Stylization and Texture Synthesis. *arXiv preprint arXiv:1701.02096*, 2017. 5
- [46] P. Upchurch, J. Gardner, K. Bala, R. Pless, N. Snavely, and K. Weinberger. Deep Feature Interpolation for Image Content Changes. *arXiv:1611.05507 [cs]*, Nov. 2016. arXiv: 1611.05507. 2
- [47] E. Vig, M. Dorr, and E. Barth. Learned saliency transformations for gaze guidance. volume 7865, page 78650W. International Society for Optics and Photonics, Feb. 2011. 2
- [48] L. K. Wong and K. L. Low. Saliency retargeting: An approach to enhance image aesthetics. In *2011 IEEE Workshop on Applications of Computer Vision (WACV)*, pages 73–80, Jan. 2011. 2
- [49] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, and H. Lipson. Understanding Neural Networks Through Deep Visualization. *arXiv:1506.06579 [cs]*, June 2015. arXiv: 1506.06579. 4
- [50] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint arXiv:1703.10593*, 2017. 2, 9