

# Network Visualization-TensorFlow

May 23, 2018

## 1 Network Visualization (TensorFlow)

In this notebook we will explore the use of *image gradients* for generating new images.

When training a model, we define a loss function which measures our current unhappiness with the model's performance; we then use backpropagation to compute the gradient of the loss with respect to the model parameters, and perform gradient descent on the model parameters to minimize the loss.

Here we will do something slightly different. We will start from a convolutional neural network model which has been pretrained to perform image classification on the ImageNet dataset. We will use this model to define a loss function which quantifies our current unhappiness with our image, then use backpropagation to compute the gradient of this loss with respect to the pixels of the image. We will then keep the model fixed, and perform gradient descent *on the image* to synthesize a new image which minimizes the loss.

In this notebook we will explore three techniques for image generation:

1. **Saliency Maps:** Saliency maps are a quick way to tell which part of the image influenced the classification decision made by the network.
2. **Fooling Images:** We can perturb an input image so that it appears the same to humans, but will be misclassified by the pretrained network.
3. **Class Visualization:** We can synthesize an image to maximize the classification score of a particular class; this can give us some sense of what the network is looking for when it classifies images of that class.

This notebook uses **TensorFlow**; we have provided another notebook which explores the same concepts in PyTorch. You only need to complete one of these two notebooks.

```
In [1]: # As usual, a bit of setup
import time, os, json
import numpy as np
import matplotlib.pyplot as plt
import tensorflow as tf

from cs231n.classifiers.squeezenet import SqueezeNet
from cs231n.data_utils import load_tiny_imagenet
from cs231n.image_utils import preprocess_image, deprocess_image
from cs231n.image_utils import SQUEEZENET_MEAN, SQUEEZENET_STD

%matplotlib inline
```

```

plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

def get_session():
    """Create a session that dynamically allocates memory."""
    # See: https://www.tensorflow.org/tutorials/using\_gpu#allowing\_gpu\_memory\_growth
    config = tf.ConfigProto()
    config.gpu_options.allow_growth = True
    session = tf.Session(config=config)
    return session

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

/Users/ianscottknight/anaconda/envs/cs231n/lib/python3.6/site-packages/h5py/__init__.py:36: FutureWarning:
from ._conv import register_converters as _register_converters

```

## 2 Pretrained Model

For all of our image generation experiments, we will start with a convolutional neural network which was pretrained to perform image classification on ImageNet. We can use any model here, but for the purposes of this assignment we will use SqueezeNet [1], which achieves accuracies comparable to AlexNet but with a significantly reduced parameter count and computational complexity.

Using SqueezeNet rather than AlexNet or VGG or ResNet means that we can easily perform all image generation experiments on CPU.

We have ported the PyTorch SqueezeNet model to TensorFlow; see: `cs231n/classifiers/squeezenet.py` for the model architecture.

To use SqueezeNet, you will need to first **download the weights** by descending into the `cs231n/datasets` directory and running `get_squeezenet_tf.sh`. Note that if you ran `get_assignment3_data.sh` then SqueezeNet will already be downloaded.

Once you've downloaded the Squeezenet model, we can load it into a new TensorFlow session:

[1] Iandola et al, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5MB model size", arXiv 2016

```

In [2]: tf.reset_default_graph()
        sess = get_session()

SAVE_PATH = 'cs231n/datasets/squeezenet.ckpt'
if not os.path.exists(SAVE_PATH + ".index"):
    raise ValueError("You need to download SqueezeNet!")
model = SqueezeNet(save_path=SAVE_PATH, sess=sess)

INFO:tensorflow:Restoring parameters from cs231n/datasets/squeezenet.ckpt

```

## 2.1 Load some ImageNet images

We have provided a few example images from the validation set of the ImageNet ILSVRC 2012 Classification dataset. To download these images, descend into `cs231n/datasets/` and run `get_imagenet_val.sh`.

Since they come from the validation set, our pretrained model did not see these images during training.

Run the following cell to visualize some of these images, along with their ground-truth labels.

```
In [3]: from cs231n.data_utils import load_imagenet_val
        X_raw, y, class_names = load_imagenet_val(num=5)
```

```
plt.figure(figsize=(12, 6))
for i in range(5):
    plt.subplot(1, 5, i + 1)
    plt.imshow(X_raw[i])
    plt.title(class_names[y[i]])
    plt.axis('off')
plt.gcf().tight_layout()
```



## 2.2 Preprocess images

The input to the pretrained model is expected to be normalized, so we first preprocess the images by subtracting the pixelwise mean and dividing by the pixelwise standard deviation.

```
In [4]: X = np.array([preprocess_image(img) for img in X_raw])
```

## 3 Saliency Maps

Using this pretrained model, we will compute class saliency maps as described in Section 3.1 of [2].

A **saliency map** tells us the degree to which each pixel in the image affects the classification score for that image. To compute it, we compute the gradient of the unnormalized score corresponding to the correct class (which is a scalar) with respect to the pixels of the image. If the image has shape  $(H, W, 3)$  then this gradient will also have shape  $(H, W, 3)$ ; for each pixel in the image, this gradient tells us the amount by which the classification score will change if the pixel changes by a small amount. To compute the saliency map, we take the absolute value of this gradient, then take the maximum value over the 3 input channels; the final saliency map thus has shape  $(H, W)$  and all entries are nonnegative.

You will need to use the `model.scores` Tensor containing the scores for each input, and will need to feed in values for the `model.image` and `model.labels` placeholder when evaluating the gradient. Open the file `cs231n/classifiers/squeezenet.py` and read the documentation to make sure you understand how to use the model. For example usage, you can see the `loss` attribute.

[2] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

```
In [109]: def compute_saliency_maps(X, y, model):
          """
          Compute a class saliency map using the model for images X and labels y.

          Input:
          - X: Input images, numpy array of shape (N, H, W, 3)
          - y: Labels for X, numpy of shape (N,)
          - model: A SqueezeNet model that will be used to compute the saliency map.

          Returns:
          - saliency: A numpy array of shape (N, H, W) giving the saliency maps for the
            input images.
          """
          saliency = None
          # Compute the score of the correct class for each example.
          # This gives a Tensor with shape [N], the number of examples.
          #
          # Note: this is equivalent to scores[np.arange(N), y] we used in NumPy
          # for computing vectorized losses.
          correct_scores = tf.gather_nd(model.scores,
                                       tf.stack((tf.range(X.shape[0]), model.labels), axis=0))
          #####
          # TODO: Produce the saliency maps over a batch of images.
          #
          # 1) Compute the loss using the correct scores tensor provided for you.
          #    (We'll combine losses across a batch by summing)
          # 2) Use tf.gradients to compute the gradient of the loss with respect
          #    to the image (accessible via model.image).
          # 3) Compute the actual value of the gradient by a call to sess.run().
          #    You will need to feed in values for the placeholders model.image and
          #    model.labels.
          # 4) Finally, process the returned gradient to compute the saliency map.
          #####
          d_X = tf.gradients(correct_scores, model.image)
          maximum = tf.reduce_max(tf.abs(d_X), axis=4)

          saliency = sess.run(tf.squeeze(maximum), feed_dict={model.image: X, model.labels: y})
          #####
          #
          #                               END OF YOUR CODE
          #
          #####
```

```
return saliency
```

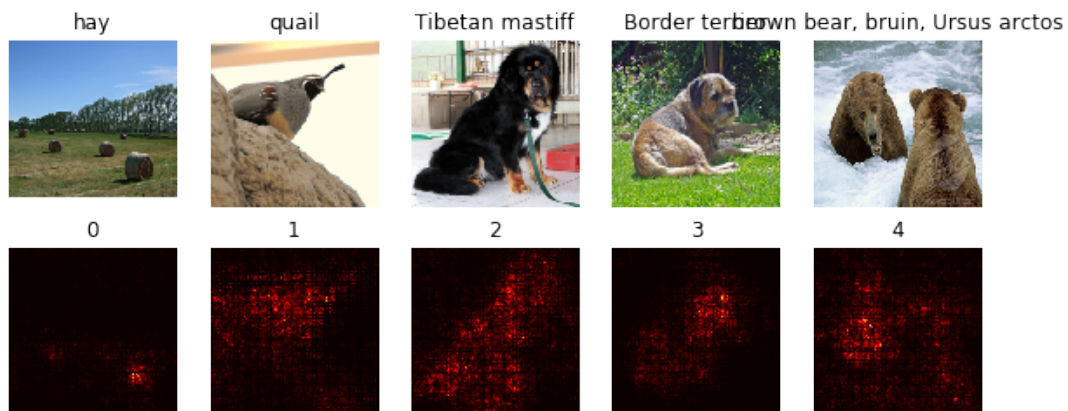
Once you have completed the implementation in the cell above, run the following to visualize some class saliency maps on our example images from the ImageNet validation set:

```
In [110]: def show_saliency_maps(X, y, mask):
            mask = np.asarray(mask)
            Xm = X[mask]
            ym = y[mask]

            saliency = compute_saliency_maps(Xm, ym, model)

            for i in range(mask.size):
                plt.subplot(2, mask.size, i + 1)
                plt.imshow(deprocess_image(Xm[i]))
                plt.axis('off')
                plt.title(class_names[ym[i]])
                plt.subplot(2, mask.size, mask.size + i + 1)
                plt.title(mask[i])
                plt.imshow(saliency[i], cmap=plt.cm.hot)
                plt.axis('off')
                plt.gcf().set_size_inches(10, 4)
            plt.show()

            mask = np.arange(5)
            show_saliency_maps(X, y, mask)
```



## 4 INLINE QUESTION

A friend of yours suggests that in order to find an image that maximizes the correct score, we can perform gradient ascent on the input image, but instead of the gradient we can actually use the saliency map in each step to update the image. Is this assertion true? Why or why not?

## 5 Fooling Images

We can also use image gradients to generate "fooling images" as discussed in [3]. Given an image and a target class, we can perform gradient **ascent** over the image to maximize the target class, stopping when the network classifies the image as the target class. Implement the following function to generate fooling images.

[3] Szegedy et al, "Intriguing properties of neural networks", ICLR 2014

```
In [111]: def make_fooling_image(X, target_y, model):
    """
    Generate a fooling image that is close to X, but that the model classifies
    as target_y.

    Inputs:
    - X: Input image, a numpy array of shape (1, 224, 224, 3)
    - target_y: An integer in the range [0, 1000)
    - model: Pretrained SqueezeNet model

    Returns:
    - X_fooling: An image that is close to X, but that is classified as target_y
    by the model.
    """

    # Make a copy of the input that we will modify
    X_fooling = X.copy()

    # Step size for the update
    learning_rate = 1

    #####
    # TODO: Generate a fooling image X_fooling that the model will classify as #
    # the class target_y. Use gradient *ascent* on the target class score, using #
    # the model.scores Tensor to get the class scores for the model.image.    #
    # When computing an update step, first normalize the gradient:              #
    #   dX = learning_rate * g / ||g||_2                                       #
    #                                                                            #
    # You should write a training loop, where in each iteration, you make an   #
    # update to the input image X_fooling (don't modify X). The loop should    #
    # stop when the predicted class for the input is the same as target_y.      #
    #                                                                            #
    # HINT: It's good practice to define your TensorFlow graph operations       #
    # outside the loop, and then just make sess.run() calls in each iteration. #
    #                                                                            #
    # HINT 2: For most examples, you should be able to generate a fooling image #
    # in fewer than 100 iterations of gradient ascent. You can print your      #
    # progress over iterations to check your algorithm.                        #
    #####
    scores = model.scores
```

```

gradient = tf.squeeze(tf.gradients(scores[:, target_y], model.image))
norm = tf.norm(gradient)

d_X = learning_rate * gradient / norm

while True:
    step, out = sess.run([d_X, scores], feed_dict={model.image: X_fooling})
    X_fooling += step
    if out[0].argmax() == target_y:
        break

#####
#                                     END OF YOUR CODE                                     #
#####
return X_fooling

```

Run the following to generate a fooling image. You should ideally see at first glance no major difference between the original and fooling images, and the network should now make an incorrect prediction on the fooling one. However you should see a bit of random noise if you look at the 10x magnified difference between the original and fooling images. Feel free to change the idx variable to explore other images.

```

In [112]: idx = 0
          Xi = X[idx][None]
          target_y = 6
          X_fooling = make_fooling_image(Xi, target_y, model)

          # Make sure that X_fooling is classified as y_target
          scores = sess.run(model.scores, {model.image: X_fooling})
          assert scores[0].argmax() == target_y, 'The network is not fooled!'

          # Show original image, fooling image, and difference
          orig_img = deprocess_image(Xi[0])
          fool_img = deprocess_image(X_fooling[0])
          # Rescale
          plt.subplot(1, 4, 1)
          plt.imshow(orig_img)
          plt.axis('off')
          plt.title(class_names[y[idx]])
          plt.subplot(1, 4, 2)
          plt.imshow(fool_img)
          plt.title(class_names[target_y])
          plt.axis('off')
          plt.subplot(1, 4, 3)
          plt.title('Difference')
          plt.imshow(deprocess_image((Xi-X_fooling)[0]))
          plt.axis('off')
          plt.subplot(1, 4, 4)

```

```
plt.title('Magnified difference (10x)')
plt.imshow(deprocess_image(10 * (Xi-X_fooling)[0]))
plt.axis('off')
plt.gcf().tight_layout()
```



## 6 Class visualization

By starting with a random noise image and performing gradient ascent on a target class, we can generate an image that the network will recognize as the target class. This idea was first presented in [2]; [3] extended this idea by suggesting several regularization techniques that can improve the quality of the generated image.

Concretely, let  $I$  be an image and let  $y$  be a target class. Let  $s_y(I)$  be the score that a convolutional network assigns to the image  $I$  for class  $y$ ; note that these are raw unnormalized scores, not class probabilities. We wish to generate an image  $I^*$  that achieves a high score for the class  $y$  by solving the problem

$$I^* = \arg \max_I (s_y(I) - R(I))$$

where  $R$  is a (possibly implicit) regularizer (note the sign of  $R(I)$  in the argmax: we want to minimize this regularization term). We can solve this optimization problem using gradient ascent, computing gradients with respect to the generated image. We will use (explicit) L2 regularization of the form

$$R(I) = \lambda \|I\|_2^2$$

**and** implicit regularization as suggested by [3] by periodically blurring the generated image. We can solve this problem using gradient ascent on the generated image.

In the cell below, complete the implementation of the `create_class_visualization` function.

[2] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

[3] Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML 2015 Deep Learning Workshop

```
In [102]: from scipy.ndimage.filters import gaussian_filter1d
def blur_image(X, sigma=1):
    X = gaussian_filter1d(X, sigma, axis=1)
    X = gaussian_filter1d(X, sigma, axis=2)
    return X
```



```

In [124]: def create_class_visualization(target_y, model, **kwargs):
          """
          Generate an image to maximize the score of target_y under a pretrained model.

          Inputs:
          - target_y: Integer in the range [0, 1000) giving the index of the class
          - model: A pretrained CNN that will be used to generate the image

          Keyword arguments:
          - l2_reg: Strength of L2 regularization on the image
          - learning_rate: How big of a step to take
          - num_iterations: How many iterations to use
          - blur_every: How often to blur the image as an implicit regularizer
          - max_jitter: How much to jitter the image as an implicit regularizer
          - show_every: How often to show the intermediate result
          """
          l2_reg = kwargs.pop('l2_reg', 1e-3)
          learning_rate = kwargs.pop('learning_rate', 25)
          num_iterations = kwargs.pop('num_iterations', 100)
          blur_every = kwargs.pop('blur_every', 10)
          max_jitter = kwargs.pop('max_jitter', 16)
          show_every = kwargs.pop('show_every', 25)

          # We use a single image of random noise as a starting point
          X = 255 * np.random.rand(224, 224, 3)
          X = preprocess_image(X)[None]

          #####
          # TODO: Compute the loss and the gradient of the loss with respect to #
          # the input image, model.image. We compute these outside the loop so #
          # that we don't have to recompute the gradient graph at each iteration #
          # #
          # Note: loss and grad should be TensorFlow Tensors, not numpy arrays! #
          # #
          # The loss is the score for the target label, target_y. You should #
          # use model.scores to get the scores, and tf.gradients to compute #
          # gradients. Don't forget the (subtracted) L2 regularization term! #
          #####

          loss = model.scores[0, target_y] # scalar loss
          grad = tf.squeeze(tf.gradients(loss, model.image)) # gradient of loss with respect to model.image

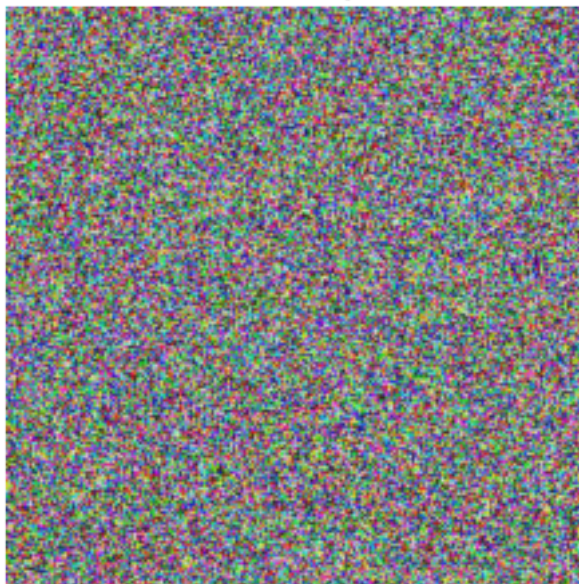
          #####
          #                                     END OF YOUR CODE                                     #
          #####

          for t in range(num_iterations):

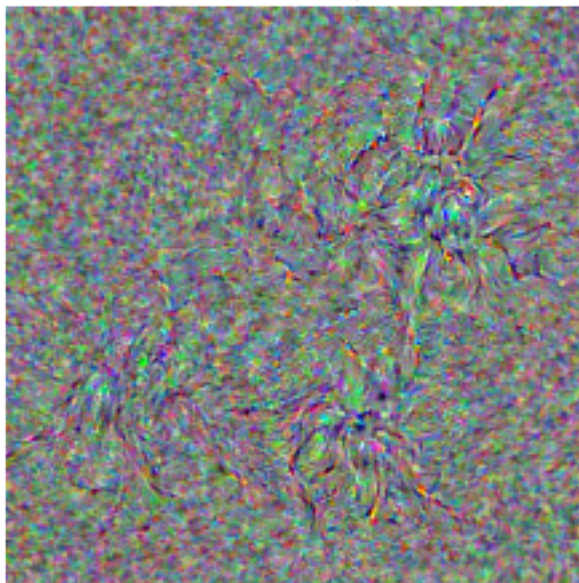
```



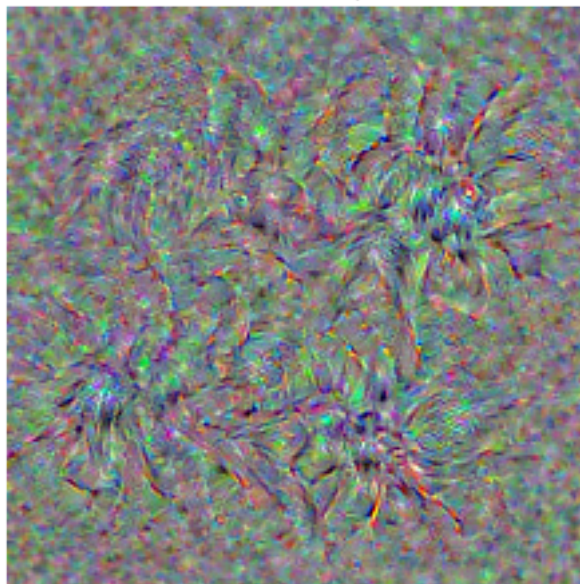
tarantula  
Iteration 1 / 100



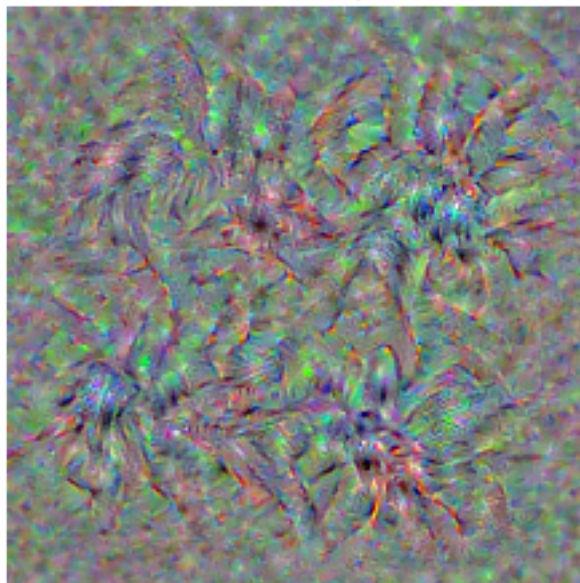
tarantula  
Iteration 25 / 100



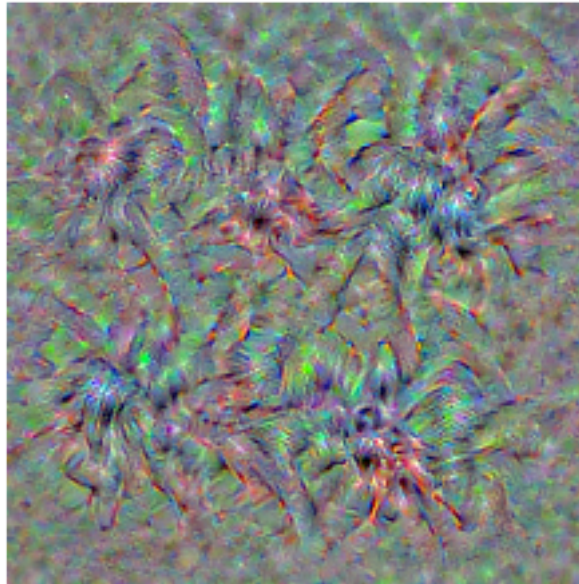
tarantula  
Iteration 50 / 100



tarantula  
Iteration 75 / 100



tarantula  
Iteration 100 / 100



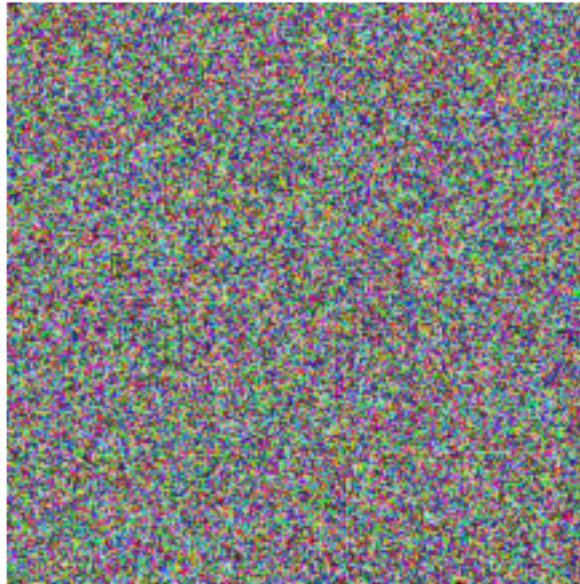
Try out your class visualization on other classes! You should also feel free to play with various hyperparameters to try and improve the quality of the generated image, but this is not required.

```
In [142]: target_y = np.random.randint(1000)
          #target_y = 78 # Tick
          #target_y = 187 # Yorkshire Terrier
          #target_y = 683 # Oboe
          #target_y = 366 # Gorilla
          #target_y = 604 # Hourglass
          target_y = 583 # Guillotine
          print(class_names[target_y])
          X = create_class_visualization(target_y, model)
          #print([[name, class_names[name]] for name in class_names])
```

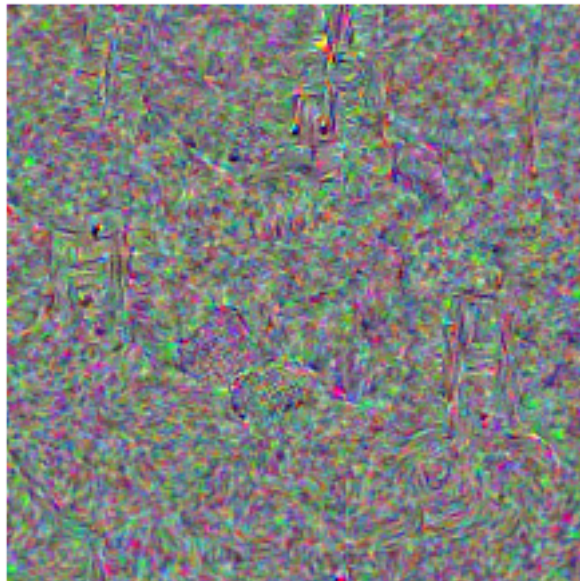
guillotine



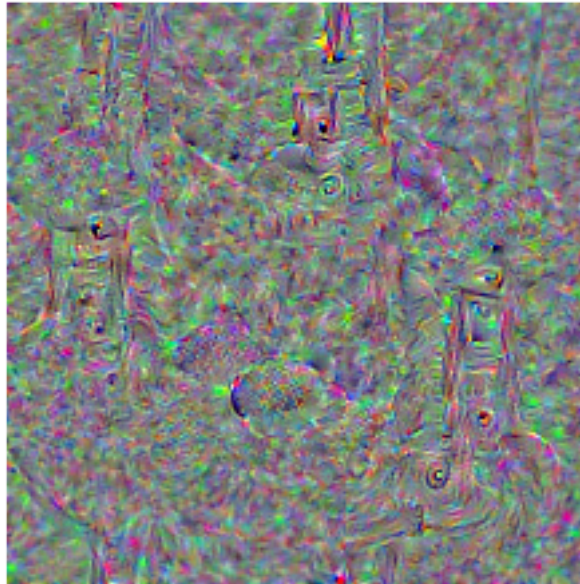
guillotine  
Iteration 1 / 100



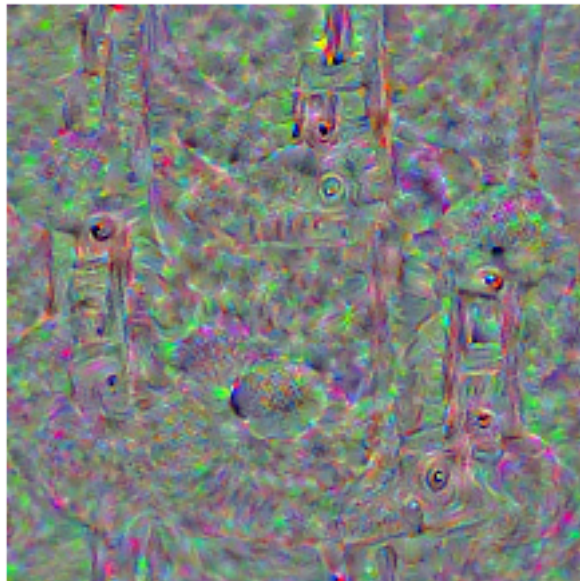
guillotine  
Iteration 25 / 100



guillotine  
Iteration 50 / 100



guillotine  
Iteration 75 / 100



guillotine  
Iteration 100 / 100

