

c. Q3의 결과를 작성한 코드와 함께 출력하세요.(20점)

## 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]:

```
1 # As usual, a bit of setup
2 import sys
3 sys.path.append('../')
4 import time, os, json
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import tensorflow as tf
8
9 from cs231n.classifiers.squeezenet import SqueezeNet
10 from cs231n.data_utils import load_tiny_imagenet
11 from cs231n.image_utils import preprocess_image, deprocess_image
12 from cs231n.image_utils import SQUEEZENET_MEAN, SQUEEZENET_STD
13
14 %matplotlib inline
15 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
16 plt.rcParams['image.interpolation'] = 'nearest'
17 plt.rcParams['image.cmap'] = 'gray'
18
19 # for auto-reloading external modules
20 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
21 %load_ext autoreload
22 %autoreload 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]:

```
1 SAVE_PATH = 'cs231n/datasets/squeezenet.ckpt'
2
3 if not os.path.exists(SAVE_PATH + ".index"):
4     raise ValueError("You need to download SqueezeNet!")
5
6 model = SqueezeNet()
7 status = model.load_weights(SAVE_PATH)
8
9 model.trainable = False
```

## 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]:

```

1 from cs231n.data_utils import load_imagenet_val
2 X_raw, y, class_names = load_imagenet_val(num=5)
3
4 plt.figure(figsize=(12, 6))
5 for i in range(5):
6     plt.subplot(1, 5, i + 1)
7     plt.imshow(X_raw[i])
8     plt.title(class_names[y[i]])
9     plt.axis('off')
10 plt.gcf().tight_layout()

```



## 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]:

```

1 X = np.array([preprocess_image(img) for img in X_raw])

```

## 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.

Open the file `cs231n/classifiers/squeezenet.py` and read the code to make sure you understand how to use the model. You will have to use `tf.GradientTape()` ([https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/GradientTape](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/GradientTape)) to compute gradients with respect to the pixels of the image. In particular, it will be very useful to read this [section](https://www.tensorflow.org/alpha/tutorials/eager/automatic_differentiation#gradient_tapes) ([https://www.tensorflow.org/alpha/tutorials/eager/automatic\\_differentiation#gradient\\_tapes](https://www.tensorflow.org/alpha/tutorials/eager/automatic_differentiation#gradient_tapes)) for better understanding.

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

## Hint: Tensorflow gather\_nd method

Recall in Assignment 1 you needed to select one element from each row of a matrix; if `s` is a numpy array of shape `(N, C)` and `y` is a numpy array of shape `(N, )` containing integers  $0 \leq y[i] < C$ , then `s[np.arange(N), y]` is a numpy array of shape `(N, )` which selects one element from each element in `s` using the indices in `y`.

In Tensorflow you can perform the same operation using the `gather_nd()` method. If `s` is a Tensor of shape `(N, C)` and `y` is a Tensor of shape `(N, )` containing longs in the range  $0 \leq y[i] < C$ , then

```
tf.gather_nd(s, tf.stack((tf.range(N), y), axis=1))
```

will be a Tensor of shape `(N, )` containing one entry from each row of `s`, selected according to the indices in `y`.

You can also read the documentation for the [gather\\_nd method](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/gather_nd) ([https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/gather\\_nd](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/gather_nd)).

In [5]:

```

1  def compute_saliency_maps(X, y, model):
2      """
3      Compute a class saliency map using the model for images X and labels y.
4
5      Input:
6      - X: Input images, numpy array of shape (N, H, W, 3)
7      - y: Labels for X, numpy of shape (N,)
8      - model: A SqueezeNet model that will be used to compute the saliency map.
9
10     Returns:
11     - saliency: A numpy array of shape (N, H, W) giving the saliency maps for the
12     input images.
13     """
14     saliency = None
15     # Compute the score of the correct class for each example.
16     # This gives a Tensor with shape [N], the number of examples.
17     #
18     # Note: this is equivalent to scores[np.arange(N), y] we used in NumPy
19     # for computing vectorized losses.
20
21     #####
22     # TODO: Produce the saliency maps over a batch of images.
23     #
24     # 1) Define a gradient tape object and watch input Image variable
25     # 2) Compute the "loss" for the batch of given input images.
26     #     - get scores output by the model for the given batch of input images
27     #     - use tf.gather_nd or tf.gather to get correct scores
28     # 3) Use the gradient() method of the gradient tape object to compute the
29     #     gradient of the loss with respect to the image
30     # 4) Finally, process the returned gradient to compute the saliency map.
31     #####
32     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
33     N, _, _, _ = X.shape
34     X = tf.convert_to_tensor(X)
35     with tf.GradientTape() as t:
36         t.watch(X)
37         scores = model.call(X)
38         correct_scores = tf.gather_nd(scores, tf.stack((tf.range(N), y), axis=1))
39         loss = tf.nn.softmax_cross_entropy_with_logits(y, correct_scores)
40     grad = t.gradient(loss, X)
41     saliency = tf.keras.backend.max(grad, axis = 3)
42
43     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
44     #####
45     #                                     END OF YOUR CODE
46     #####
47     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 [6]:

```

1 def show_saliency_maps(X, y, mask):
2     mask = np.asarray(mask)
3     Xm = X[mask]
4     ym = y[mask]
5
6     saliency = compute_saliency_maps(Xm, ym, model)
7
8     for i in range(mask.size):
9         plt.subplot(2, mask.size, i + 1)
10        plt.imshow(deprocess_image(Xm[i]))
11        plt.axis('off')
12        plt.title(class_names[ym[i]])
13        plt.subplot(2, mask.size, mask.size + i + 1)
14        plt.title(mask[i])
15        plt.imshow(saliency[i], cmap=plt.cm.hot)
16        plt.axis('off')
17        plt.gcf().set_size_inches(10, 4)
18    plt.show()
19
20 mask = np.arange(5)
21 show_saliency_maps(X, y, mask)

```



## 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?

**Your Answer:** saliency map은 이미지 segmentation과 비슷한 형태를 보여준다. 그래서 image에 대한 gradient를 구하는 방식은 Simonyan et al. 2014가 보여주었다 이미지 재건설에 saliency map을 사용하여 가능하다는 것을 실험적으로 증명하였다.

## 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 [7]:

```

1  def make_fooling_image(X, target_y, model):
2      """
3      Generate a fooling image that is close to X, but that the model classifies
4      as target_y.
5
6      Inputs:
7      - X: Input image, a numpy array of shape (1, 224, 224, 3)
8      - target_y: An integer in the range [0, 1000)
9      - model: Pretrained SqueezeNet model
10
11      Returns:
12      - X_fooling: An image that is close to X, but that is classified as target_y
13      by the model.
14      """
15
16      # Make a copy of the input that we will modify
17      X_fooling = X.copy()
18
19      # Step size for the update
20      learning_rate = 1
21
22      #####
23      # TODO: Generate a fooling image X_fooling that the model will classify as
24      # the class target_y. Use gradient *ascent* on the target class score, using
25      # the model.scores Tensor to get the class scores for the model.image. #
26      # When computing an update step, first normalize the gradient:
27      #   dX = learning_rate * g / ||g||_2
28      #
29      # You should write a training loop, where in each iteration, you make an
30      # update to the input image X_fooling (don't modify X). The loop should
31      # stop when the predicted class for the input is the same as target_y.
32      #
33      # HINT: Use tf.GradientTape() to keep track of your gradients and
34      # use tape.gradient to get the actual gradient with respect to X_fooling.
35      #
36      # HINT 2: For most examples, you should be able to generate a fooling image
37      # in fewer than 100 iterations of gradient ascent. You can print your
38      # progress over iterations to check your algorithm.
39      #####
40      # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
41
42      X = tf.convert_to_tensor(X)
43      X_fooling = tf.convert_to_tensor(X_fooling)
44      target_score = 0
45      i = 0
46      while i != 100:
47          with tf.GradientTape() as t:
48              t.watch(X_fooling)
49              scores = model.call(X_fooling)
50              target_score = scores[0, target_y]
51              grad = t.gradient(target_score, X_fooling)
52              dX = learning_rate * tf.math.l2_normalize(grad)
53              X_fooling += dX
54              i += 1
55
56      # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
57      #####
58      #                                     END OF YOUR CODE
59      #####

```

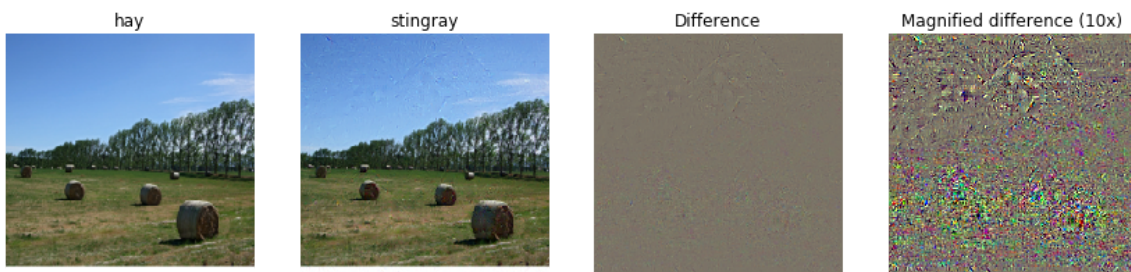


```
60     return X_fooling
```

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 [8]:

```
1  idx = 0
2  Xi = X[idx][None]
3  target_y = 6
4  X_fooling = make_fooling_image(Xi, target_y, model)
5
6  # Make sure that X_fooling is classified as y_target
7  scores = model(X_fooling)
8  assert tf.math.argmax(scores[0]).numpy() == target_y, 'The network is not fooled'
9
10 # Show original image, fooling image, and difference
11 orig_img = deprocess_image(Xi[0])
12 fool_img = deprocess_image(X_fooling[0])
13 plt.figure(figsize=(12, 6))
14
15 # Rescale
16 plt.subplot(1, 4, 1)
17 plt.imshow(orig_img)
18 plt.axis('off')
19 plt.title(class_names[y[idx]])
20 plt.subplot(1, 4, 2)
21 plt.imshow(fool_img)
22 plt.title(class_names[target_y])
23 plt.axis('off')
24 plt.subplot(1, 4, 3)
25 plt.title('Difference')
26 plt.imshow(deprocess_image((Xi-X_fooling)[0]))
27 plt.axis('off')
28 plt.subplot(1, 4, 4)
29 plt.title('Magnified difference (10x)')
30 plt.imshow(deprocess_image(10 * (Xi-X_fooling)[0]))
31 plt.axis('off')
32 plt.gcf().tight_layout()
```



## 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 [9]:

```
1 from scipy.ndimage.filters import gaussian_filter1d
2 def blur_image(X, sigma=1):
3     X = gaussian_filter1d(X, sigma, axis=1)
4     X = gaussian_filter1d(X, sigma, axis=2)
5     return X
```

In [10]:

```
1 def jitter(X, ox, oy):
2     """
3     Helper function to randomly jitter an image.
4
5     Inputs
6     - X: Tensor of shape (N, H, W, C)
7     - ox, oy: Integers giving number of pixels to jitter along W and H axes
8
9     Returns: A new Tensor of shape (N, H, W, C)
10    """
11    if ox != 0:
12        left = X[:, :, :-ox]
13        right = X[:, :, -ox:]
14        X = tf.concat([right, left], axis=2)
15    if oy != 0:
16        top = X[:, :-oy]
17        bottom = X[:, -oy:]
18        X = tf.concat([bottom, top], axis=1)
19    return X
```

In [11]:

```

1  def create_class_visualization(target_y, model, **kwargs):
2      """
3      Generate an image to maximize the score of target_y under a pretrained model
4
5      Inputs:
6      - target_y: Integer in the range [0, 1000) giving the index of the class
7      - model: A pretrained CNN that will be used to generate the image
8
9      Keyword arguments:
10     - l2_reg: Strength of L2 regularization on the image
11     - learning_rate: How big of a step to take
12     - num_iterations: How many iterations to use
13     - blur_every: How often to blur the image as an implicit regularizer
14     - max_jitter: How much to jitter the image as an implicit regularizer
15     - show_every: How often to show the intermediate result
16     """
17     l2_reg = kwargs.pop('l2_reg', 1e-3)
18     learning_rate = kwargs.pop('learning_rate', 25)
19     num_iterations = kwargs.pop('num_iterations', 100)
20     blur_every = kwargs.pop('blur_every', 10)
21     max_jitter = kwargs.pop('max_jitter', 16)
22     show_every = kwargs.pop('show_every', 25)
23
24     # We use a single image of random noise as a starting point
25     X = 255 * np.random.rand(224, 224, 3)
26     X = preprocess_image(X)[None]
27
28     loss = None # scalar loss
29     grad = None # gradient of loss with respect to model.image, same size as model.image
30
31     X = tf.Variable(X)
32     for t in range(num_iterations):
33         # Randomly jitter the image a bit; this gives slightly nicer results
34         ox, oy = np.random.randint(0, max_jitter, 2)
35         X = jitter(X, ox, oy)
36
37         #####
38         # TODO: Compute the value of the gradient of the score for
39         # class target_y with respect to the pixels of the image, and make a
40         # gradient step on the image using the learning rate. You should use
41         # the tf.GradientTape() and tape.gradient to compute gradients.
42         #
43         # Be very careful about the signs of elements in your code.
44         #####
45         # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
46
47         with tf.GradientTape() as tape:
48             tape.watch(X)
49             scores = model.call(X)
50             target_score = scores[0, target_y]
51             grad = tape.gradient(target_score, X)
52             dX = learning_rate * grad
53             X -= l2_reg * ((X - tf.reduce_mean(X)) ** 2)
54             X += dX
55
56         # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
57         #####
58         #
59         # END OF YOUR CODE
60         #####

```

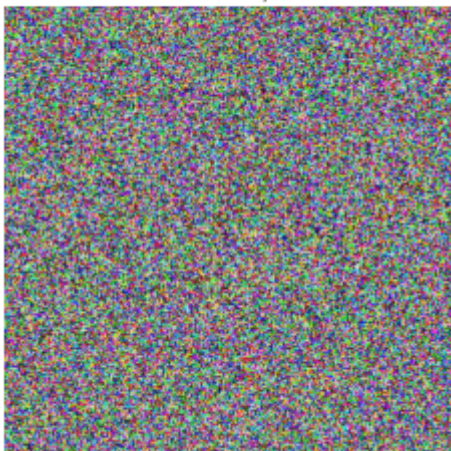
```
60
61     # Undo the jitter
62     X = jitter(X, -ox, -oy)
63     # As a regularizer, clip and periodically blur
64
65     X = tf.clip_by_value(X, -SQUEEZENET_MEAN/SQUEEZENET_STD, (1.0 - SQUEEZENET_MEAN/SQUEEZENET_STD))
66     if t % blur_every == 0:
67         X = blur_image(X, sigma=0.5)
68
69     # Periodically show the image
70     if t == 0 or (t + 1) % show_every == 0 or t == num_iterations - 1:
71         plt.imshow(deprocess_image(X[0]))
72         class_name = class_names[target_y]
73         plt.title('%s\nIteration %d / %d' % (class_name, t + 1, num_iterations))
74         plt.gcf().set_size_inches(4, 4)
75         plt.axis('off')
76         plt.show()
77     return X
```

Once you have completed the implementation in the cell above, run the following cell to generate an image of Tarantula:

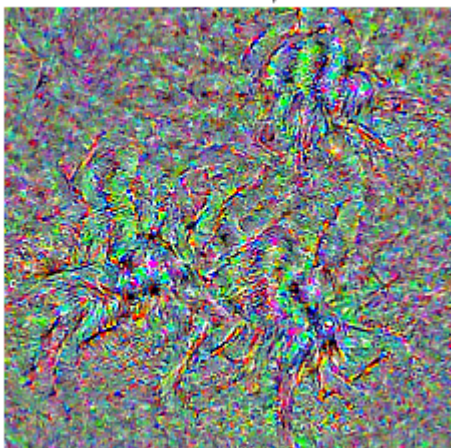
In [12]:

```
1 target_y = 76 # Tarantula
2 out = create_class_visualization(target_y, model)
```

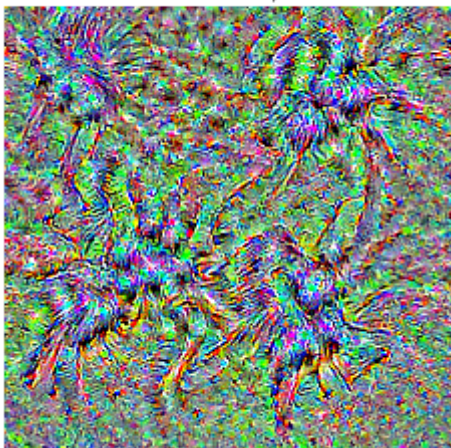
tarantula  
Iteration 1 / 100



tarantula  
Iteration 25 / 100

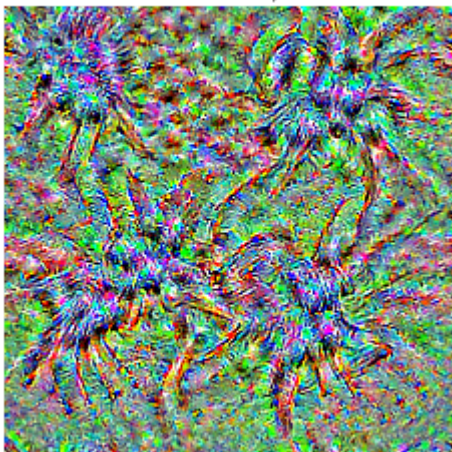


tarantula  
Iteration 50 / 100

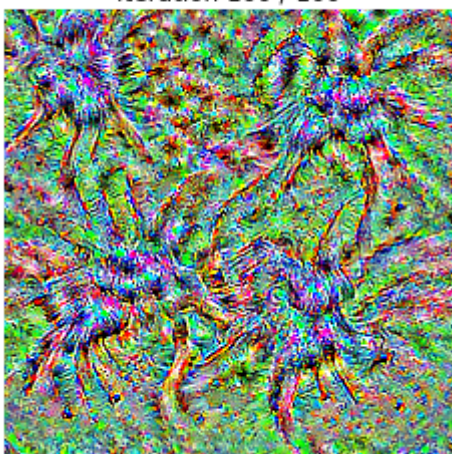




tarantula  
Iteration 75 / 100



tarantula  
Iteration 100 / 100



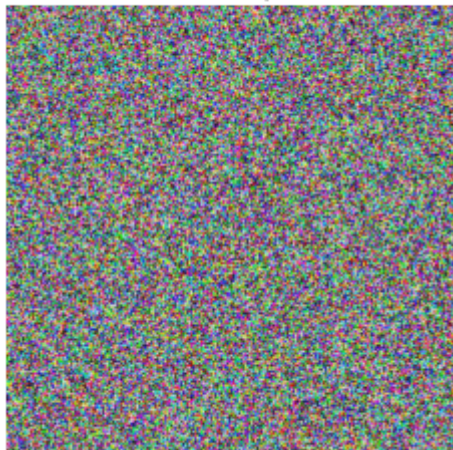
Try out your class visualization on other classes! You should also feel free to play with various hyperparameters to try and improve the quality of the generated image, but this is not required.

In [14]:

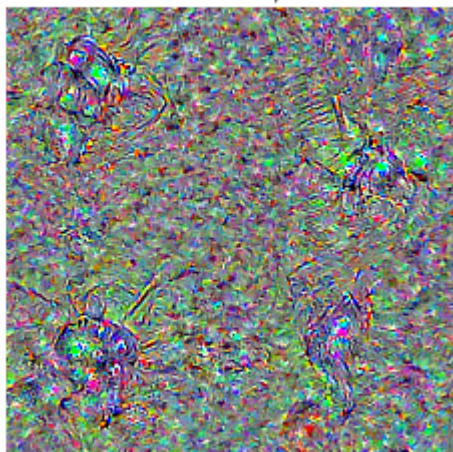
```
1 target_y = np.random.randint(1000)
2 # target_y = 78 # Tick
3 # target_y = 187 # Yorkshire Terrier
4 # target_y = 683 # Oboe
5 # target_y = 366 # Gorilla
6 # target_y = 604 # Hourglass
7 print(class_names[target_y])
8 X = create_class_visualization(target_y, model)
```

scuba diver

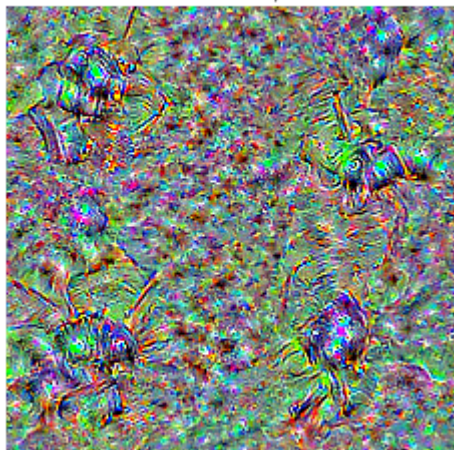
scuba diver  
Iteration 1 / 100



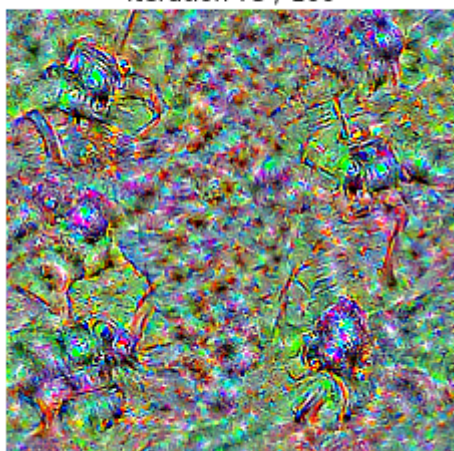
scuba diver  
Iteration 25 / 100



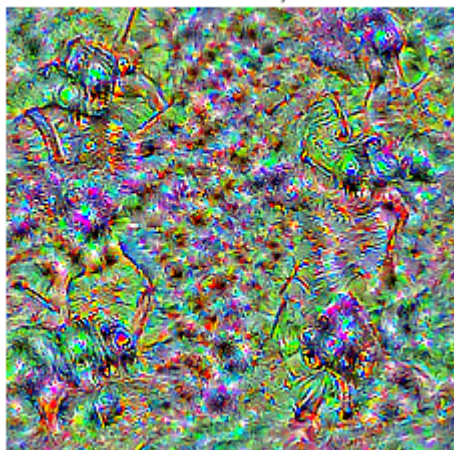
scuba diver  
Iteration 50 / 100



scuba diver  
Iteration 75 / 100



scuba diver  
Iteration 100 / 100



In [ ]:

1



