

[Open in Colab](#)[https://colab.research.google.com/github/shreyasrajesh0308/ECE188DeepLearning/blob/main/Adv\\_](https://colab.research.google.com/github/shreyasrajesh0308/ECE188DeepLearning/blob/main/Adv_)

## Perform an Adversarial attack.

For the second part of the project we consider a trained model (MobileNet) which is trained on the imagenet dataset.

We use an evasion attack called [FGSM](https://neptune.ai/blog/adversarial-attacks-on-neural-networks-exploring-the-fast-gradient-sign-method#:~:text=The%20Fast%20Gradient%20Sign%20Method%20(FGSM)%20combines%20a%20gradient%20sign%20method%20with%20a%20perturbation%20to%20create%20an%20adversarial%20image) ([https://neptune.ai/blog/adversarial-attacks-on-neural-networks-exploring-the-fast-gradient-sign-method#:~:text=The%20Fast%20Gradient%20Sign%20Method%20\(FGSM\)%20combines%20a%20gradient%20sign%20method%20with%20a%20perturbation%20to%20create%20an%20adversarial%20image](https://neptune.ai/blog/adversarial-attacks-on-neural-networks-exploring-the-fast-gradient-sign-method#:~:text=The%20Fast%20Gradient%20Sign%20Method%20(FGSM)%20combines%20a%20gradient%20sign%20method%20with%20a%20perturbation%20to%20create%20an%20adversarial%20image)) to fool the neural network into making incorrect predictions.

## Import Packages.

Import the necessary packages we continue to use Tensorflow and Keras

```
In [1]:  
1 import tensorflow as tf  
2 import matplotlib as mpl  
3 import matplotlib.pyplot as plt  
4 from keras.preprocessing import image  
5  
6 mpl.rcParams['figure.figsize'] = (8, 8)  
7 mpl.rcParams['axes.grid'] = False
```

## Load the Pretrained model.

We use the [MobileNetV2](https://arxiv.org/abs/1801.04381) (<https://arxiv.org/abs/1801.04381>) model trained on the [Imagenet](https://www.image-net.org/) (<https://www.image-net.org/>) dataset.

```
In [2]:  
1 pretrained_model = tf.keras.applications.MobileNetV2(include_top=True,  
2                                                       weights='imagenet')  
3 pretrained_model.trainable = False  
4  
5 # ImageNet Labels  
6 decode_predictions = tf.keras.applications.mobilenet_v2.decode_predictions
```

## Helper Function for Data Processing

Following functions can be used for data processing. Dont worry about these, just use them.

```
In [3]:  
1 # Helper function to preprocess the image so that it can be inputted in MobileNet  
2 def preprocess(image):  
3     image = tf.cast(image, tf.float32)  
4     image = tf.image.resize(image, (224, 224))  
5     image = tf.keras.applications.mobilenet_v2.preprocess_input(image)  
6     image = image[None, ...]  
7     return image  
8  
9 # Helper function to extract labels from probability vector  
10 def get_imagenet_label(probs):  
11     return decode_predictions(probs, top=1)[0][0]
```

## Load an Image.

Load any image, we consider an image of a Panda Bear.

```
In [4]:  
1 image_raw = tf.io.read_file('panda.jpg')  
2 image = tf.image.decode_image(image_raw)  
3  
4 image = preprocess(image)  
5 image_probs = pretrained_model.predict(image)
```

1/1 [=====] - 1s 556ms/step

In [5]:

```
1 plt.figure()
2 plt.imshow(image[0] * 0.5 + 0.5) # To change [-1, 1] to [0,1]
3 _, image_class, class_confidence = get_imagenet_label(image_probs)
4 plt.title('{} : {:.2f}% Confidence'.format(image_class, class_confidence*100
5 plt.show()
```

Downloading data from [https://storage.googleapis.com/download.tensorflow.org/data/imagenet\\_class\\_index.json](https://storage.googleapis.com/download.tensorflow.org/data/imagenet_class_index.json) ([https://storage.googleapis.com/download.tensorflow.org/data/imagenet\\_class\\_index.json](https://storage.googleapis.com/download.tensorflow.org/data/imagenet_class_index.json))

35363/35363 [=====] - 0s 1us/step



## Create the Adversarial Image.

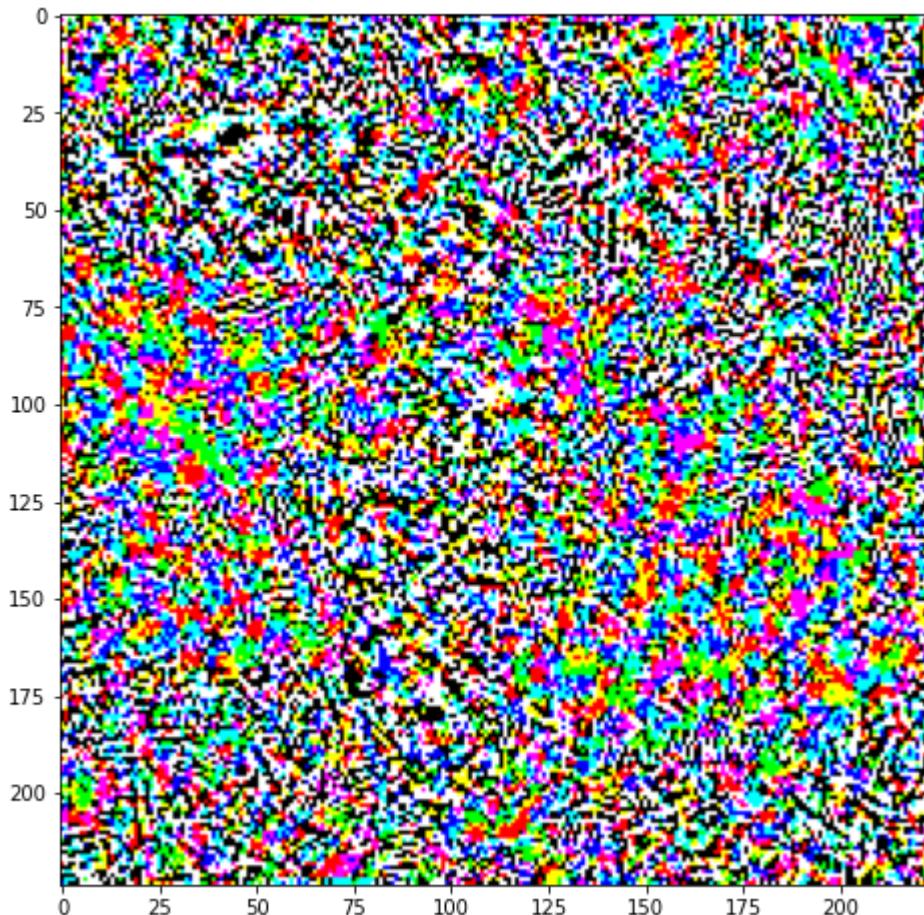
We use the FGSM method to create an adversarial image. Be sure to read about FGSM to understand how the attack works.

In [6]:

```
1 loss_object = tf.keras.losses.CategoricalCrossentropy()
2
3 def create_adversarial_pattern(input_image, input_label):
4     with tf.GradientTape() as tape:
5         tape.watch(input_image)
6         prediction = pretrained_model(input_image)
7         loss = loss_object(input_label, prediction)
8
9     # Get the gradients of the loss w.r.t to the input image.
10    gradient = tape.gradient(loss, input_image)
11    # Get the sign of the gradients to create the perturbation
12    signed_grad = tf.sign(gradient)
13    return signed_grad
```

In [7]:

```
1 # Get the input label of the image.
2 giant_panda_index = 388
3 label = tf.one_hot(giant_panda_index, image_probs.shape[-1])
4 label = tf.reshape(label, (1, image_probs.shape[-1]))
5
6 perturbations = create_adversarial_pattern(image, label)
7 plt.imshow(perturbations[0] * 0.5 + 0.5); # To change [-1, 1] to [0,1]
```



In [8]:

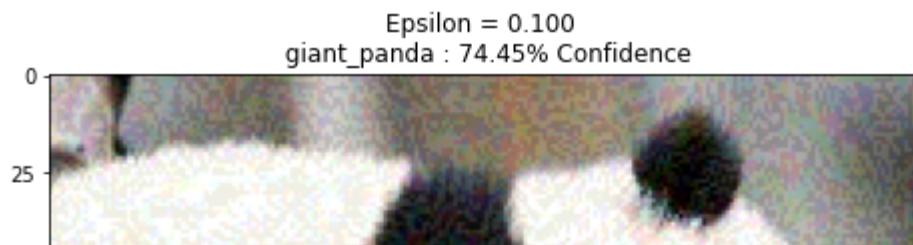
```
1 def display_images(image, description):
2     _, label, confidence = get_imagenet_label(pretrained_model.predict(image))
3     plt.figure()
4     plt.imshow(image[0]*0.5+0.5)
5     plt.title('{} \n {} : {:.2f}% Confidence'.format(description,
6                                                       label, confidence*100))
7     plt.show()
```

In [11]:

```

1  epsilons = [0, 0.01, 0.1, 0.15, 0.3]
2  descriptions = [('Epsilon = {:.3f}'.format(eps) if eps else 'Input')
3                  for eps in epsilons]
4  print(descriptions)
5  for i, eps in enumerate(epsilons):
6      adv_x = image + eps*perturbations
7      adv_x = tf.clip_by_value(adv_x, -1, 1)
8      display_images(adv_x, descriptions[i])

```



## Task2: Perform an Analysis to understand the potency of the attack.

Your task here is to understand how small a change could change the class output and this is measured by the epsilon value needed to change the class.

Your task is as follows:

- Pick 10 images each from different classes in imagenet.
- Perform a perturbation analysis on each of these images.
- In the analysis you are required to report the smallest epsilon value for which you notice a class change.
- Make a table for each of the images considered with the minimum epsilon value for the FGSM attack.

Write the Code for the above below. You can add the table also below.

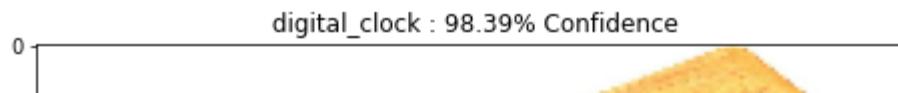
10 Imagenet classes chosen

- Goldfinch 11
- Boa Constrictor 61
- Scorpion 71
- Pizza 963

- Buckle 464
- Digital clock 530
- Dutch oven 544
- Scabbard 777
- Wooden Spoon 910
- Geyser 974

In [26]:

```
1 # Get Images
2 chosen_images = []
3 chosen_images_probs = []
4 chosen_indexes = [11, 61, 71, 963, 464, 530, 544, 777, 910, 974]
5 for i in range(0, 10):
6     image_raw = tf.io.read_file('image'+ str(i) + '.jpg')
7     image = tf.image.decode_image(image_raw)
8     image = preprocess(image)
9     chosen_images.append(image)
10    image_probs = pretrained_model.predict(image)
11    chosen_images_probs.append(image_probs)
12
13    display_images(image, "")
```



In [27]:

```
1 chosen_perturbations = []
2
3 for i in range(0, 10):
4     label = tf.one_hot(chosen_indexes[i], chosen_images_probs[i].shape[-1])
5     label = tf.reshape(label, (1, chosen_images_probs[i].shape[-1]))
6
7     perturbations = create_adversarial_pattern(chosen_images[i], label)
8     # plt.imshow(perturbations[0] * 0.5 + 0.5); # To change [-1, 1] to [0,1]
9     chosen_perturbations.append(perturbations)
```

In [28]:

```
1 import numpy as np
2 epsilons = np.arange(0, 1, 0.01)
3 descriptions = [('Epsilon = {:.3f}'.format(eps) if eps else 'Input')
4                  for eps in epsilons]
```

In [29]:

```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[0] + eps*chosen_perturbations[0]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```

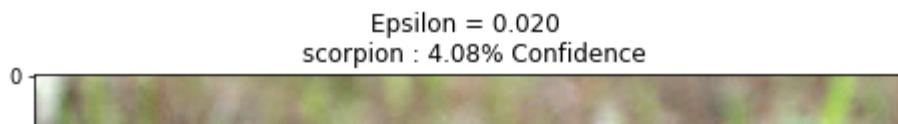
In [30]:

```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[1] + eps*chosen_perturbations[1]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```



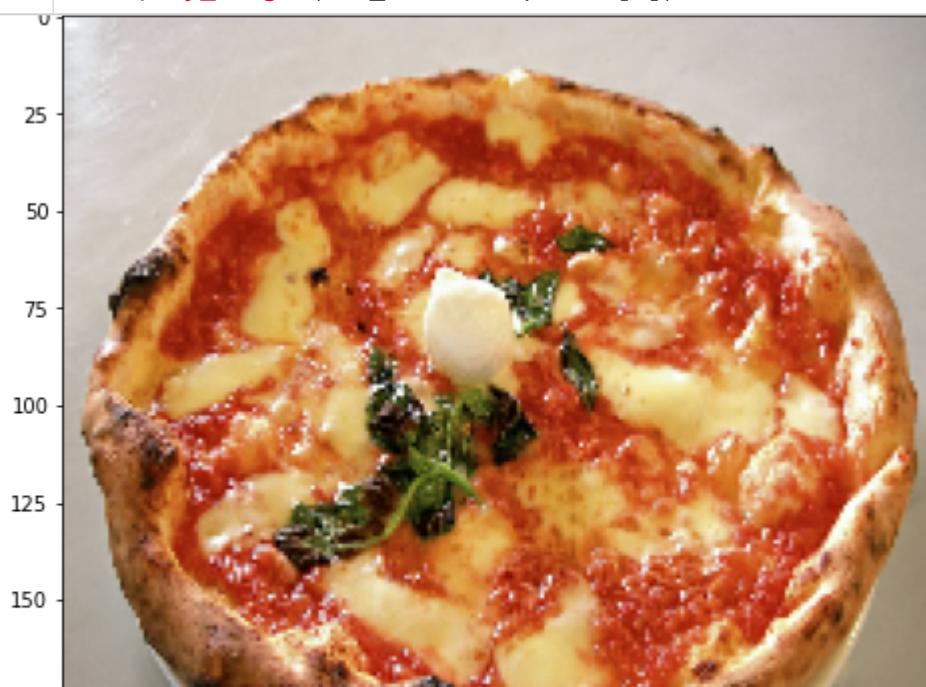
In [31]:

```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[2] + eps*chosen_perturbations[2]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```



In [32]:

```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[3] + eps*chosen_perturbations[3]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```



In [33]:

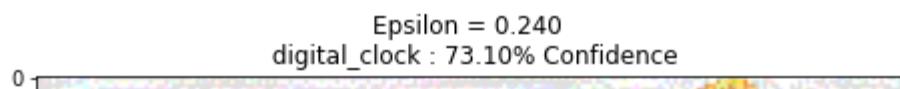
```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[4] + eps*chosen_perturbations[4]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```

1/1 [=====] - 0s 24ms/step



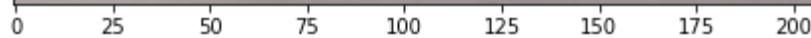
In [34]:

```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[5] + eps*chosen_perturbations[5]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```



In [35]:

```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[6] + eps*chosen_perturbations[6]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```



1/1 [=====] - 0s 24ms/step

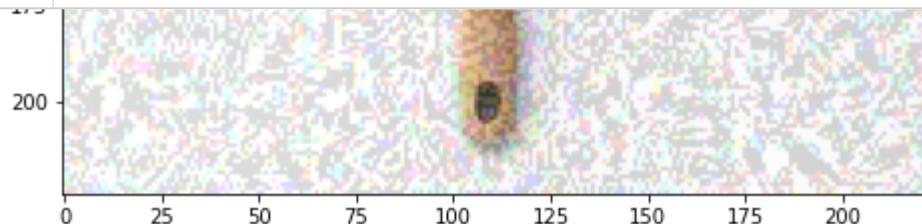
In [36]:

```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[7] + eps*chosen_perturbations[7]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```



In [37]:

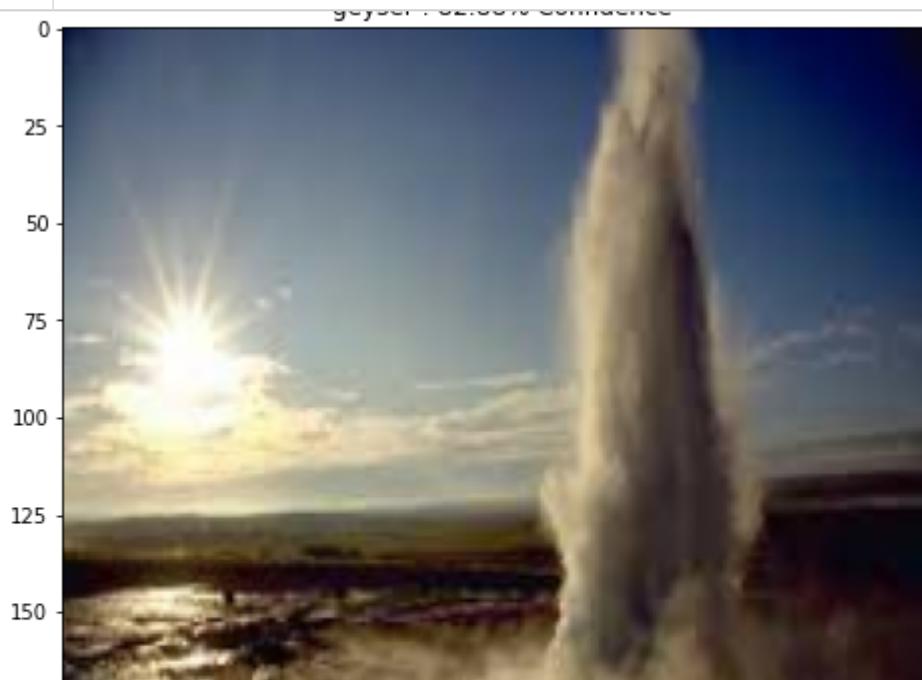
```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[8] + eps*chosen_perturbations[8]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```



1/1 [=====] - 0s 24ms/step

In [38]:

```
1 for i, eps in enumerate(epsilons):
2     adv_x = chosen_images[9] + eps*chosen_perturbations[9]
3     adv_x = tf.clip_by_value(adv_x, -1, 1)
4     display_images(adv_x, descriptions[i])
```



1	Image Class	Minimum Epsilon for FGSM
2	-----	-----
3	Goldfinch	0.17
4	Boa Constrictor	0.01
5	Scorpion	0.03
6	Pizza	0.01

7	Buckle   0.01
8	Digital clock   0.44
9	Dutch oven   0.01
10	Scabbard   0.01
11	Wooden Spoon   0.28
12	Geyser   0.01

As one can see, the FGSM attack is very effective against MobileNet, with only a small amount of epsilon required to change an image's class most of the time. In 6/10 of the images, only an epsilon of 0.01 was required. In 3/10, a somewhat larger epsilon was required to change the class, perhaps because of the clearer backgrounds/more defining features of the object classes present, meaning the images are more far away from any decision boundaries, meaning even in the direction of greatest loss, it requires more epsilon to actually cross a decision boundary. Thus, these perturbed images are slightly more identifiable.

## Task3: Compare the robustness of the considered model with other models.

Your task here is to compare how this model (MobileNetV2) compares with other popular object detection models.

Your task is as follows:

- Consider 5 different models (you can consider various RESNET architectures, any models you find interesting).
- Load the pre-trained weights of the model (trained on imagenet).
- Perform Task2 on all the considered models.
- Add all the results in the table. Hence the final table you have 6 columns for each model and epsilon values for each of the 10 images for all 6 models.

What do you observe? Why do you think this is the case?

Write the Code for the above below. You can also add the table and answer to the question below.

Models used:

ResNet101V2

ResNet50V2

ResNet152V2

NasNetMobile

DenseNet201

I had issues with VGGNet and so was not able to use it

In [59]:

```

1 # Helper function to preprocess the image so that it can be inputted in Mobi
2 def preprocess_v(image, model_ref):
3     image = tf.cast(image, tf.float32)
4     image = tf.image.resize(image, (224, 224))
5     image = model_ref.preprocess_input(image)
6     image = image[None, ...]
7     return image
8
9 # Helper function to extract labels from probability vector
10 def get_imagenet_label_v(probs, model_ref):
11     # ImageNet labels
12     return model_ref.decode_predictions(probs, top=1)[0][0]

```

In [42]:

```

1 models = [
2     tf.keras.applications.ResNet101V2(include_top=True,weights='imagenet'),
3     tf.keras.applications.ResNet50V2(include_top=True,weights='imagenet'),
4     tf.keras.applications.ResNet152V2(include_top=True,weights='imagenet'),
5     tf.keras.applications.VGG19(include_top=True,weights='imagenet'),
6     tf.keras.applications.DenseNet201(include_top=True,weights='imagenet')
7 ]

```

Downloading data from [https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet101v2\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet101v2_weights_tf_dim_ordering_tf_kernels.h5) ([https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet101v2\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet101v2_weights_tf_dim_ordering_tf_kernels.h5))

179518384/179518384 [=====] - 272s 2us/step

Downloading data from [https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet50v2\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet50v2_weights_tf_dim_ordering_tf_kernels.h5) ([https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet50v2\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet50v2_weights_tf_dim_ordering_tf_kernels.h5))

102869336/102869336 [=====] - 50s 0us/step

Downloading data from [https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet152v2\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet152v2_weights_tf_dim_ordering_tf_kernels.h5) ([https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet152v2\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/resnet/resnet152v2_weights_tf_dim_ordering_tf_kernels.h5))

242745792/242745792 [=====] - 108s 0us/step

Downloading data from [https://storage.googleapis.com/tensorflow/keras-applications/vgg19/vgg19\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/vgg19/vgg19_weights_tf_dim_ordering_tf_kernels.h5) ([https://storage.googleapis.com/tensorflow/keras-applications/vgg19/vgg19\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/vgg19/vgg19_weights_tf_dim_ordering_tf_kernels.h5))

574710816/574710816 [=====] - 568s 1us/step

Downloading data from [https://storage.googleapis.com/tensorflow/keras-applications/densenet/densenet201\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/densenet/densenet201_weights_tf_dim_ordering_tf_kernels.h5) ([https://storage.googleapis.com/tensorflow/keras-applications/densenet/densenet201\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels.h5](https://storage.googleapis.com/tensorflow/keras-applications/densenet/densenet201_weights_tf_dim_ordering_tf_kernels.h5))

82524592/82524592 [=====] - 46s 1us/step

In [116]:

```

1 models[3] = tf.keras.applications.NASNetMobile(include_top=True,weights='ima'

```

Downloading data from <https://storage.googleapis.com/tensorflow/keras-applications/nasnet/NASNet-mobile.h5> (<https://storage.googleapis.com/tensorflow/keras-applications/nasnet/NASNet-mobile.h5>)

24227760/24227760 [=====] - 12s 1us/step

In [117]:

```
1 model_refs = [
2     tf.keras.applications.resnet_v2,
3     tf.keras.applications.resnet_v2,
4     tf.keras.applications.resnet_v2,
5     tf.keras.applications.nasnet,
6     tf.keras.applications.densenet
7 ]
```

In [118]:

```
1 for i in range(0, 5):
2     models[i].trainable = False
```

In [119]:

```
1 loss_object = tf.keras.losses.CategoricalCrossentropy()
2
3 def create_adversarial_pattern(input_image, input_label, model):
4     with tf.GradientTape() as tape:
5         tape.watch(input_image)
6         prediction = model(input_image)
7         loss = loss_object(input_label, prediction)
8
9     # Get the gradients of the loss w.r.t to the input image.
10    gradient = tape.gradient(loss, input_image)
11    # Get the sign of the gradients to create the perturbation
12    signed_grad = tf.sign(gradient)
13    return signed_grad
```

In [120]:

```
1 # print(tf.keras.applications.resnet101v2.decode_predictions)
```

In [121]:

```
1 model_images = []
2 model_images_probs = []
3 chosen_indexes = [11, 61, 71, 963, 464, 530, 544, 777, 910, 974]
4 for j in range(0, 5):
5     model = models[j]
6     model_ref = model_refs[j]
7     images = []
8     probs = []
9     for i in range(0, 10):
10         image_raw = tf.io.read_file('image'+ str(i) + '.jpg')
11         image = tf.image.decode_image(image_raw)
12         image = preprocess_v(image, model_ref)
13         images.append(image)
14         image_probs = model.predict(image)
15         probs.append(image_probs)
16
17 #         display_images(image, "")
18 model_images.append(images)
19 model_images_probs.append(probs)
20
```

```
1/1 [=====] - 0s 84ms/step
1/1 [=====] - 0s 80ms/step
1/1 [=====] - 0s 91ms/step
1/1 [=====] - 0s 83ms/step
1/1 [=====] - 0s 77ms/step
1/1 [=====] - 0s 80ms/step
1/1 [=====] - 0s 80ms/step
1/1 [=====] - 0s 81ms/step
1/1 [=====] - 0s 90ms/step
1/1 [=====] - 0s 81ms/step
1/1 [=====] - 0s 49ms/step
1/1 [=====] - 0s 51ms/step
1/1 [=====] - 0s 53ms/step
1/1 [=====] - 0s 55ms/step
1/1 [=====] - 0s 58ms/step
1/1 [=====] - 0s 57ms/step
1/1 [=====] - 0s 57ms/step
1/1 [=====] - 0s 52ms/step
1/1 [=====] - 0s 53ms/step
1/1 [=====] - 0s 49ms/step
1/1 [=====] - 0s 118ms/step
1/1 [=====] - 0s 116ms/step
1/1 [=====] - 0s 121ms/step
1/1 [=====] - 0s 118ms/step
1/1 [=====] - 0s 122ms/step
1/1 [=====] - 0s 115ms/step
1/1 [=====] - 0s 117ms/step
1/1 [=====] - 0s 114ms/step
1/1 [=====] - 0s 115ms/step
1/1 [=====] - 0s 119ms/step
1/1 [=====] - 2s 2s/step
1/1 [=====] - 0s 29ms/step
```

```
1/1 [=====] - 0s 28ms/step  
1/1 [=====] - 0s 28ms/step  
1/1 [=====] - 0s 29ms/step  
1/1 [=====] - 0s 29ms/step  
1/1 [=====] - 0s 93ms/step  
1/1 [=====] - 0s 88ms/step  
1/1 [=====] - 0s 95ms/step  
1/1 [=====] - 0s 87ms/step  
1/1 [=====] - 0s 88ms/step  
1/1 [=====] - 0s 91ms/step  
1/1 [=====] - 0s 91ms/step  
1/1 [=====] - 0s 91ms/step  
1/1 [=====] - 0s 88ms/step  
1/1 [=====] - 0s 100ms/step
```

In [122]:

```
1 model_perturbations = []  
2 for j in range(0, 5):  
3     perturbs = []  
4     for i in range(0, 10):  
5         label = tf.one_hot(chosen_indexes[i], model_images_probs[j][i].shape[  
6             1])  
7         label = tf.reshape(label, (1, model_images_probs[j][i].shape[-1]))  
8  
9         perturbations = create_adversarial_pattern(model_images[j][i], label)  
10        # plt.imshow(perturbations[0] * 0.5 + 0.5); # To change [-1, 1] to  
11        perturbs.append(perturbations)  
12    model_perturbations.append(perturbs)
```

In [126]:

```
1 im_labels = [
2     "goldfinch",
3     "boa_constrictor",
4     "scorpion",
5     "pizza",
6     "buckle",
7     "digital_clock",
8     "Dutch_oven",
9     "scabbard",
10    "wooden_spoon",
11    "geyser"
12 ]
13
14 def get_lowest_epsilon_and_display(image, pert, actual_label, model, model_r
15     eps = 0
16     adv_x = image + eps*pert
17     adv_x = tf.clip_by_value(adv_x, -1, 1)
18     _, label, confidence = get_imagenet_label_v(model.predict(adv_x), mode
19     while label == actual_label:
20         eps = eps + 0.01
21         adv_x = image + eps*pert
22         adv_x = tf.clip_by_value(adv_x, -1, 1)
23         _, label, confidence = get_imagenet_label_v(model.predict(adv_x), mode
24     #     print(eps)
25     #     print(label)
26     #     print(actual_label)
27     plt.figure()
28     plt.imshow(adv_x[0]*0.5+0.5)
29     plt.title('Epsilon = {:.3f}\n{} : {:.2f}% Confidence'.format(eps, label,
30     plt.show()
31     return eps
```

In [127]:

```
1 j = 0
2 model_epsilonss = []
3 epsilonss = []
4 for i in range(0, 10):
5     epsilonss.append(get_lowest_epsilon_and_display(model_images[j][i], model_
6 model_epsilonss.append(epsilonss)
```

1/1 [=====] - 0s 84ms/step  
1/1 [=====] - 0s 92ms/step  
1/1 [=====] - 0s 94ms/step  
1/1 [=====] - 0s 88ms/step  
1/1 [=====] - 0s 86ms/step  
1/1 [=====] - 0s 85ms/step  
1/1 [=====] - 0s 86ms/step  
1/1 [=====] - 0s 81ms/step  
1/1 [=====] - 0s 83ms/step  
1/1 [=====] - 0s 83ms/step  
1/1 [=====] - 0s 89ms/step  
1/1 [=====] - 0s 91ms/step  
1/1 [=====] - 0s 86ms/step  
1/1 [=====] - 0s 80ms/step  
1/1 [=====] - 0s 83ms/step  
1/1 [=====] - 0s 85ms/step  
1/1 [=====] - 0s 81ms/step  
1/1 [=====] - 0s 81ms/step  
1/1 [=====] - 0s 86ms/step  
1/1 [=====] - 0s 84ms/step

In [128]:

```
1 print(model_epsilonss)
```

```
[[0.3500000000000014, 0.03, 0.02, 0.01, 0.6600000000000004, 1.020000000000000
7, 0.01, 0.02, 0.4100000000000002, 0.01]]
```

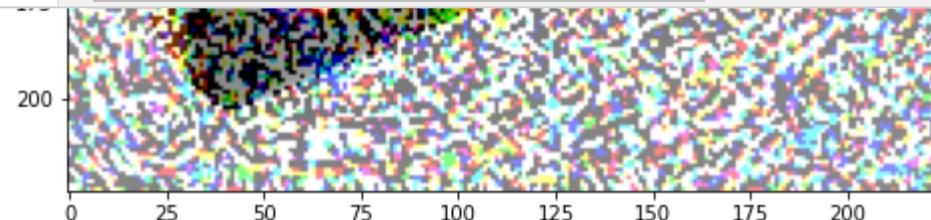
In [129]:

```
1 j = 1
2 epsilons = []
3 for i in range(0, 10):
4     epsilons.append(get_lowest_epsilon_and_display(model_images[j][i], model_
5 model_epsilons.append(epsilons)
```

1/1 [=====] - 0s 55ms/step  
1/1 [=====] - 0s 58ms/step  
1/1 [=====] - 0s 57ms/step  
1/1 [=====] - 0s 52ms/step  
1/1 [=====] - 0s 54ms/step  
1/1 [=====] - 0s 55ms/step  
1/1 [=====] - 0s 52ms/step  
1/1 [=====] - 0s 53ms/step  
1/1 [=====] - 0s 51ms/step  
1/1 [=====] - 0s 50ms/step  
1/1 [=====] - 0s 55ms/step  
1/1 [=====] - 0s 53ms/step  
1/1 [=====] - 0s 57ms/step  
1/1 [=====] - 0s 53ms/step  
1/1 [=====] - 0s 55ms/step  
1/1 [=====] - 0s 60ms/step  
1/1 [=====] - 0s 53ms/step  
1/1 [=====] - 0s 54ms/step  
1/1 [=====] - 0s 60ms/step  
1/1 [=====] - 0s 57ms/step

In [131]:

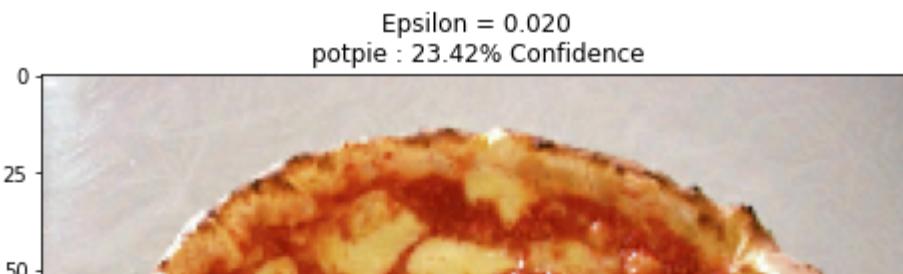
```
1 j = 2
2 epsilons = []
3 for i in range(0, 10):
4     epsilons.append(get_lowest_epsilon_and_display(model_images[j][i], model_
5 model_epsilons.append(epsilons)
```



1/1 [=====] - 0s 128ms/step  
1/1 [=====] - 0s 142ms/step

In [132]:

```
1 j = 3
2 epsilons = []
3 for i in range(0, 10):
4     epsilons.append(get_lowest_epsilon_and_display(model_images[j][i], model_
5 model_epsilons.append(epsilons)
```



In [133]:

```
1
2
3 for i in range(0, 10):
4     epsilons.append(get_lowest_epsilon_and_display(model_images[j][i], model_perturbations
5 model_epsilons.append(epsilons)
```



```
In [135]: 1 print(model_epsilon)
```

```
In [139]: 1 mobileneteps = [0.17,0.01,0.03,0.01,0.01,0.44,0.01,0.01,0.01,0.28,0.01]
 2 print(np.mean(model_averages, axis=1))
 3 print(np.mean(mobileneteps))
 4
```

[0.254 0.145 0.268 0.278 0.266]  
0.098

In general, the ResNet, NASNet and DenseNet models were more robust to the FGSM attack than MobileNet. This is shown by the fact that the average required epsilon to change the class of an image was higher for all the tested models compared to MobileNet. This is probably the case because of the higher complexity of the other models as compared to MobileNet. As one can see in the code above, the other models have a lot more weights (on the order of  $10^8$  as opposed to MobileNet's approximately  $10^4$  weights). Thus, because the models have more complex functions, their loss functions are also more complex, making it harder to find and navigate to a decision boundary for a given image.

Image Class	MobileNet	ResNet101V2	ResNet50V2	ResNet152V2	NasNetMobile	DenseNet201
Goldfinch	0.17	0.35	0.27	0.38	0.27	0.59
Boa Constrictor	0.01	0.03	0.03	0.14	0.01	0.07
Scorpion	0.03	0.02	0.27	0.33	0.34	0.06
Pizza	0.01	0.01	0.01	0.01	0.02	0.04
Buckle	0.01	0.66	0.42	0.48	0.73	0.35
Digital clock	0.44	1.02	0.03	0.91	0.96	1.23
Dutch oven	0.01	0.01	0.01	0.01	0.01	0.15
Scabbard	0.01	0.02	0.01	0.01	0.01	0.04
Wooden Spoon	0.28	0.41	0.39	0.40	0.38	0.11
Geyser	0.01	0.01	0.01	0.01	0.05	0.02

**BONUS: Can you provide a better attack?**

Can you design a better attack that lowers the epsilon required for the images?

### Task:

- Design another attack.
- Compare the epsilon values on 10 images.
- Does it perform better than the FGSM attack? That is, does it have lower epsilon values?

Write the code and provide your answers below.

In [179]:

```

1 def create_adversarial_another(input_image, input_label, model):
2     #Add random noise to the input image (RAND-FGSM)
3
4     with tf.GradientTape() as tape:
5         tape.watch(input_image)
6         prediction = model(input_image + 0.001 * np.random.normal(size=(1, 224,
7             loss = loss_object(input_label, prediction)
8
9     # Get the gradients of the loss w.r.t to the input image.
10    gradient = tape.gradient(loss, input_image)
11    # Get the sign of the gradients to create the perturbation
12    signed_grad = tf.sign(gradient)
13    return signed_grad

```

In [180]:

```

1 probs = []
2 images = []
3 model = models[3]
4 model_ref = model_refs[3]
5
6 for i in range(0, 10):
7     image_raw = tf.io.read_file('image'+ str(i) + '.jpg')
8     image = tf.image.decode_image(image_raw)
9     image = preprocess_v(image, model_ref)
10    images.append(image)
11    image_probs = model.predict(image)
12    probs.append(image_probs)
13
14    # display_images(image, "")
15 model_images.append(images)
16 model_images_probs.append(probs)

```

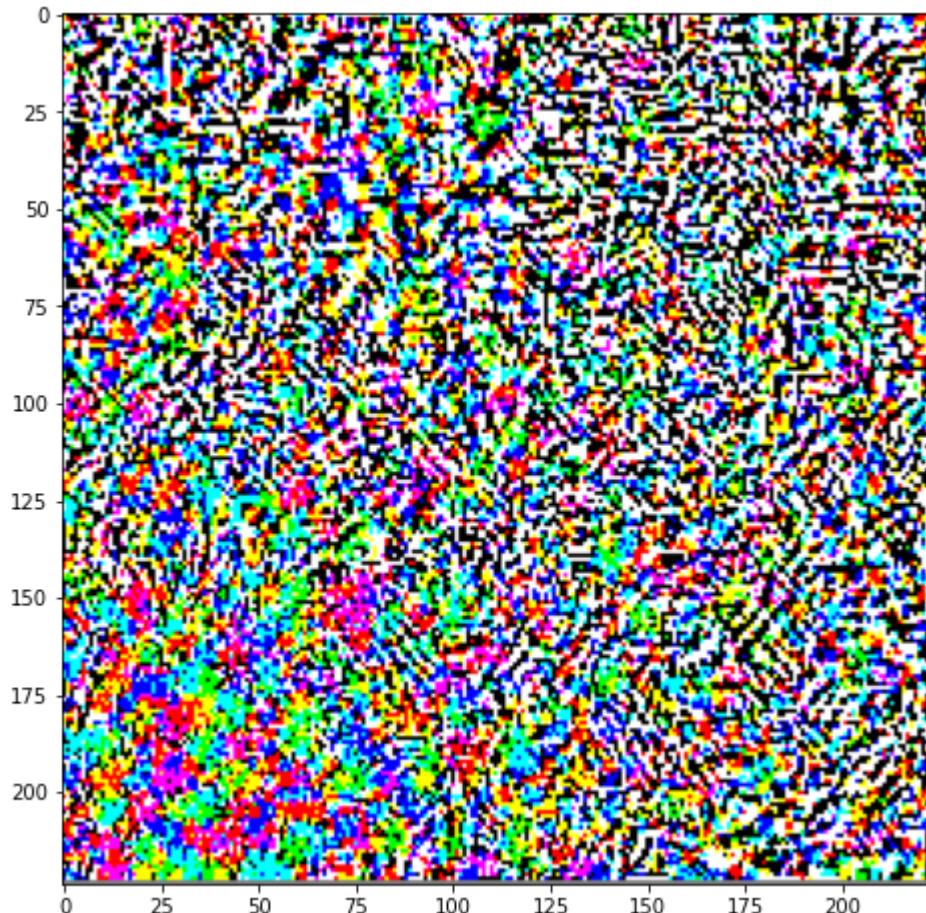
```

1/1 [=====] - 0s 33ms/step
1/1 [=====] - 0s 32ms/step
1/1 [=====] - 0s 32ms/step
1/1 [=====] - 0s 31ms/step
1/1 [=====] - 0s 34ms/step
1/1 [=====] - 0s 32ms/step
1/1 [=====] - 0s 32ms/step
1/1 [=====] - 0s 32ms/step
1/1 [=====] - 0s 31ms/step
1/1 [=====] - 0s 32ms/step

```

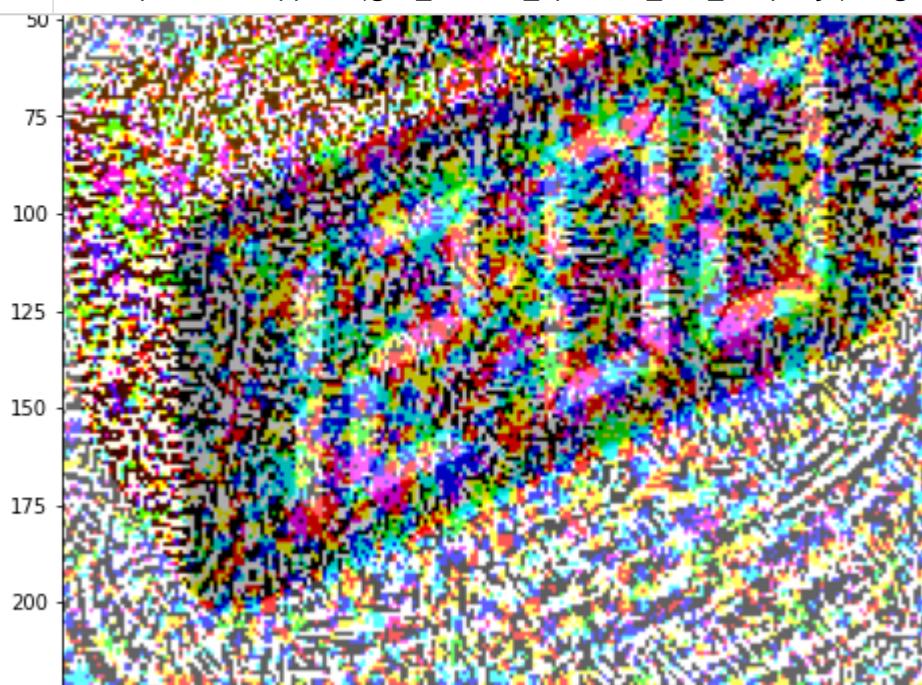
In [181]:

```
1 perbs = []
2
3 for i in range(0, 10):
4     label = tf.one_hot(chosen_indexes[i], probs[i].shape[-1])
5     label = tf.reshape(label, (1, probs[i].shape[-1]))
6
7     perturbations = create_adversarial_another(images[i], label, model)
8     plt.imshow(perturbations[0] * 0.5 + 0.5); # To change [-1, 1] to [0,1]
9     perbs.append(perturbations)
```



In [182]:

```
1  epsilons = []
2  for i in range(0, 10):
3      epsilons.append(get_lowest_epsilon_and_display(images[i], perbs[i], im_1
```



In [183]:

```
1  print(epsilons)
2  print(np.mean(epsilons))
```

```
[0.2200000000000006, 0.01, 0.3100000000000001, 0.02, 0.7000000000000004, 1.180
000000000008, 0.01, 0.01, 0.3700000000000016, 0.0600000000000005]
0.2890000000000015
```

I attempted to add random variation when updating the adversarial samples-- this is to help mask the gradient and overcome any adversarial training the model might have. We end up with worse performance, as the required epsilons are overall higher. This is because while adding the random mask might help with any existing robustness measures, it actively harms our attempts to get to the direction of the optimal loss to get to the nearest decision boundary. Thus, it takes more epsilon to actually affect the image.