# Project Report: Visualizing Convolutional Neural Networks

**Amey Agrawal & Jaikumar Balani**
{f2014148, f2014022}@pilani.bits-pilani.ac.in

## 1   Abstract

Convolutional neural networks are becoming increasingly popular and have been the centre of the deep learning revolution but yet training a large CNN model has remained a mystery. We are attempting to build a framework to visualise a CNN model, using multiple methods. Including deconvolution nets proposed by Zeiler & Fergus (2014), activation based method proposed by Erhan et al. (2009) and t-SNE.

## 2   Introduction

Convolutional neural networks have become a popular technique in all kinds of computer vision task since they showed exceptional results on ImageNet challenge in 2011. Today CNNs are used for all kinds of applications including sentiment analysis of short sentences, generating voice from text and creating 3D models from 2D images.
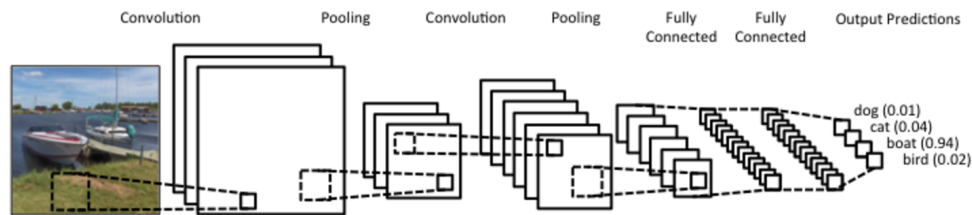


Figure 1: Genral feature extraction pipe of convolutional neural networks

Visualising the activations on each layer is the simplest and the most intuitive way of visualising convolutional neural networks but it has its obvious shortcomings. As the depth of filters is more than 3 it is generally hard to make sense beyond the first convolution layer.

## 3   Problem Statment

To build a framework to visualise a CNN model, using deconvolution nets proposed by Zeiler & Fergus (2014), activation based method proposed by Erhan et al. (2009) and t-SNE.

## 4   Related work

Multiple techniques were suggested by Erhan et al. (2009) for visualisation of CNNs based upon activations of individual units. Simonyan et al. (2013) established a connection between gradient based methods of CNN visualisation and deconvolution networks. Zeiler  Fergus (2014) introduced a way to visualise CNNs using deconvolution nets.

Following is a brief review of the major works on which we have developed our toolkit.

## 4.1   METHOD 1: DECONVOLUTIONAL NETWORK

Deconvolution networks reveal the part of the input that excites individual feature maps at any layer in the model. Deconvolutional Network with switched pooling layers is used to project the activations back to the input pixel space. We use 2 different variations of this method which differ only in the fact that one stores the positions of the max values during pooling and reproduces similarly at the time of unpooling while the other only stores the maximum values and then reconstructs after unpooling.
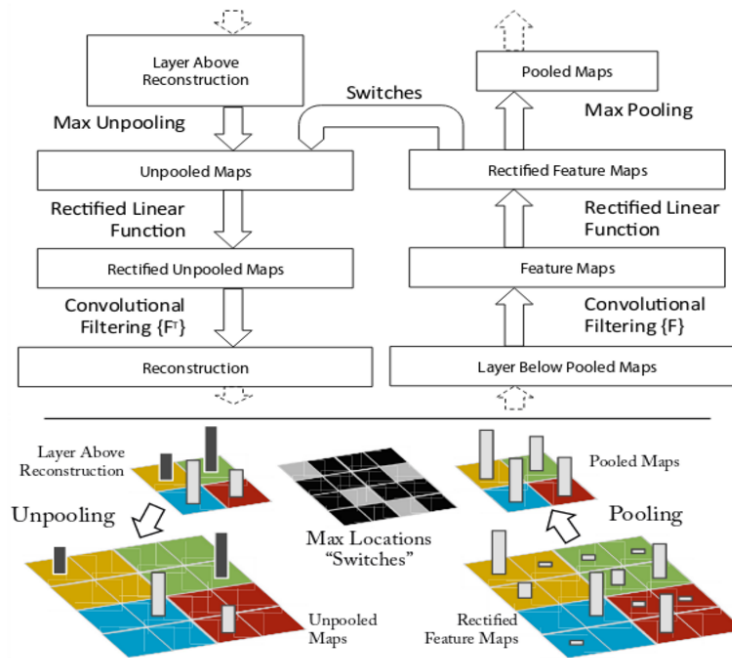


Figure 2: Deconvolution network with switched pooling



Figure 3: 16th Filter of layer 17th VGGnet visualized using deconvnet

## 4.2    METHOD 2: GRADIENT ASCENT ON ACTIVATION

In this method, we define the loss function as sum activations of all the neurones in the chosen filter. This loss is the maximised using gradient accent which enables us to visualise what excites a the given filter in general.
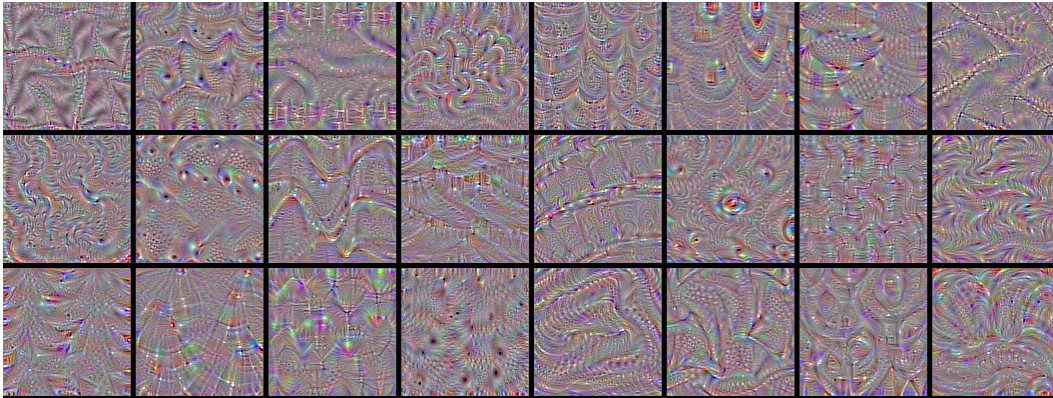


Figure 4: Visulalization of some filters of VGGnet using method 2.

## 4.3    METHOD 3: T-SNE VISUALISATION OF DENSE LAYERS

Fully connected layers from CNNs have been used in multiple transfer learning application concluding image captioning and sentiment analysis. Hence visualising the last fully connected layer on a CNN can give us an understanding of what the network has learned overall. Maaten & Hinton (2008) proposed an optimisation based technique to visualise high dimensional vectors in lower dimensions called t-SNE. t-SNE tries preserve distance between two points between the two vector spaces.
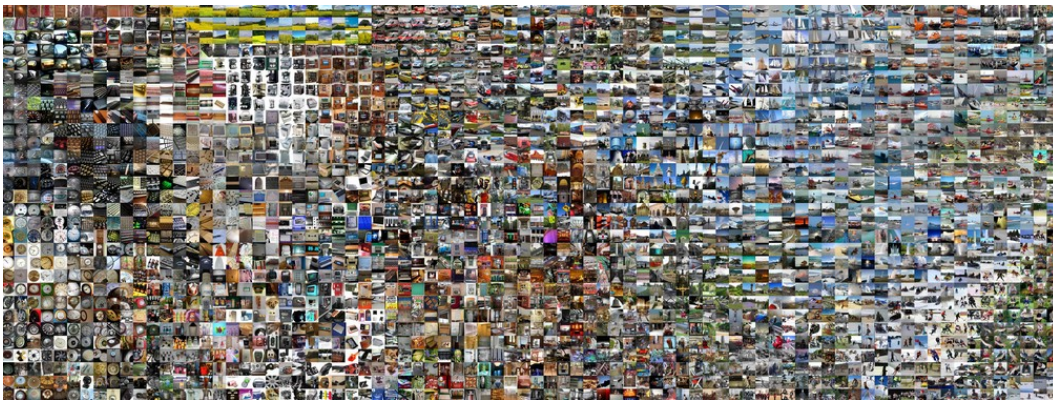


Figure 5: A visualization done using t-SNE from work of Stanford's CS 231n

t-SNE is a computationally extensive algorithm. Our implementation is a single threaded CPU based upon Scikit-learn. Hence while trying to visualise last layers of large models like VGGnet if we cannot use a large number of images.

## 5    DESIGN OF SOLUTION

The framework is built in python with Keras. Keras has been gaining popularity as deep learning framework due to its simplicity and flexibility. Keras can use both TensorFlow and Theano backends.

We have developed an extensible driver module with command interactive line interface. It can take any keras model stored in hdf5 format as input. If the user does not provide an input model VGG16 checkpoint is used by default.

Keras is used for automatic differentiation in method 1. While custom keras layers are made to create the deconvolution counterpart of given model. Scikit-learn is used to find the t-SNE embeddings of the last layer of the network. And the images are plotted on a 2D graph on their corresponding coordinates.

Matplotlib is used to plot and store images. Sometimes the output plotted by matplotlib looks very noisy though it is stored perfectly.

## 6 SCREESHOTS AND VISUALS

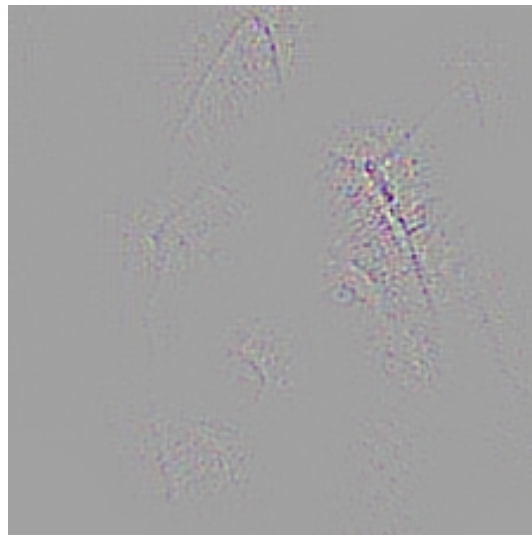Following are few filters visualised with method 1.



Figure 6: Deconvnet: Output 1

Following are some visualisation using method 2.

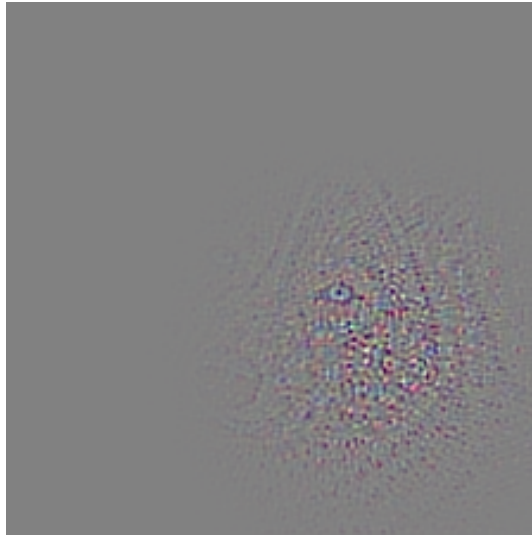Following is few screenshots from the application:
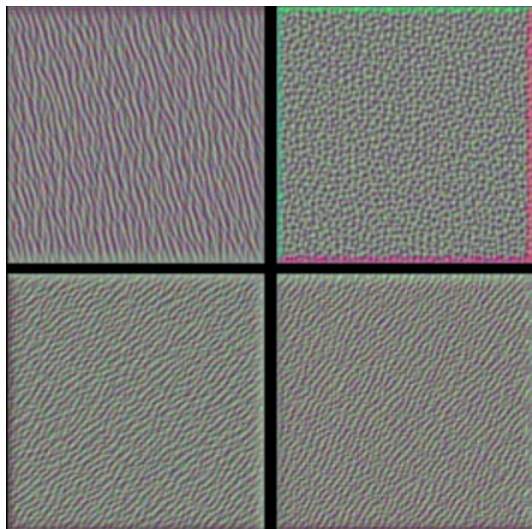
Figure 7: Deconvnet: Output 2
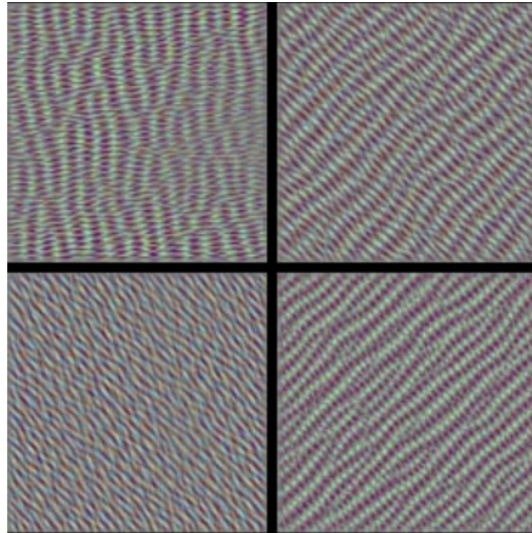


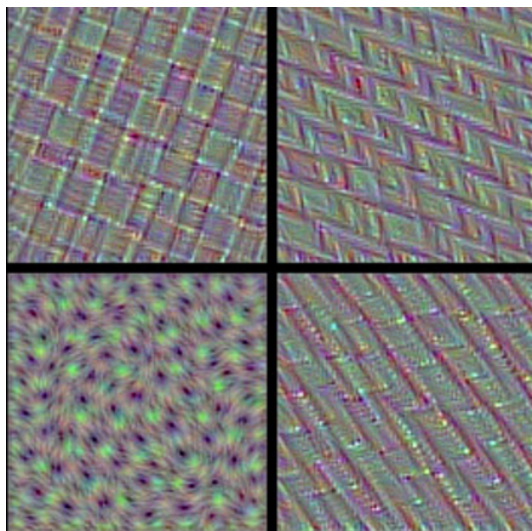Figure 8: Gradient acent: Output 1

Figure 9: Gradient acent: Output 2



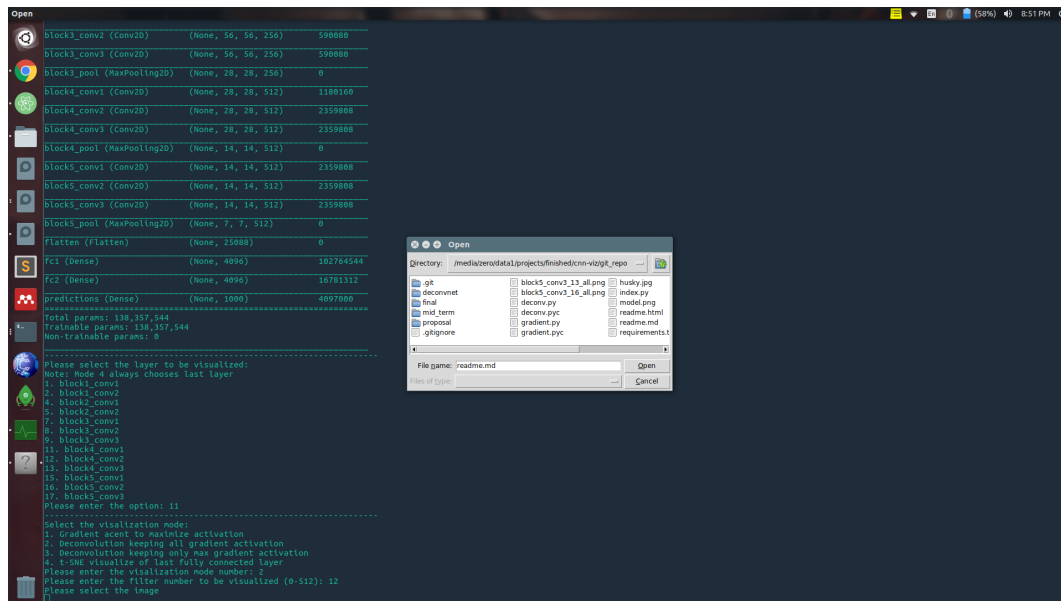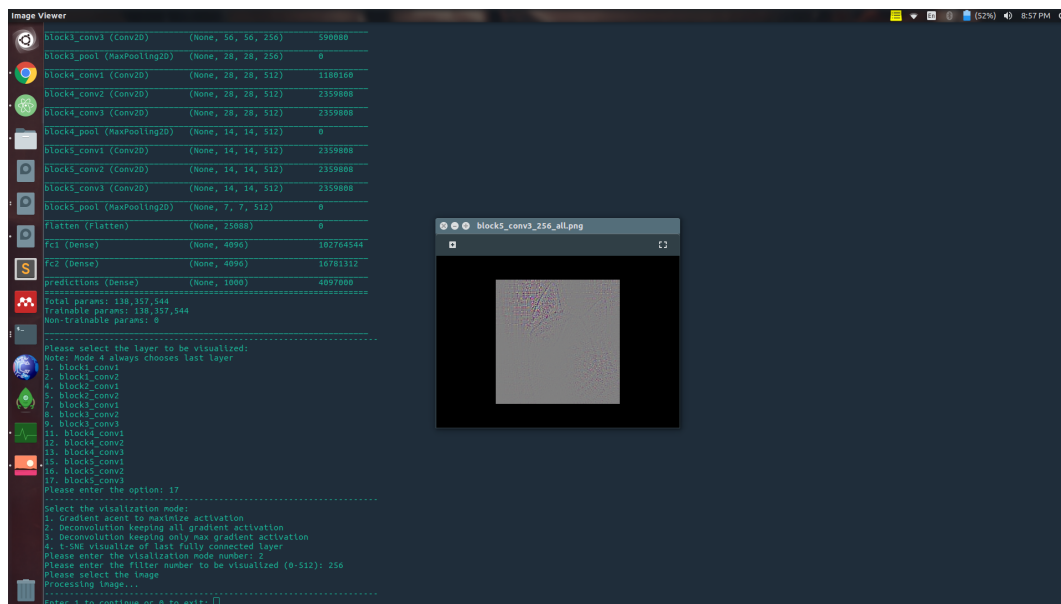Figure 10: Gradient acent: Output 3

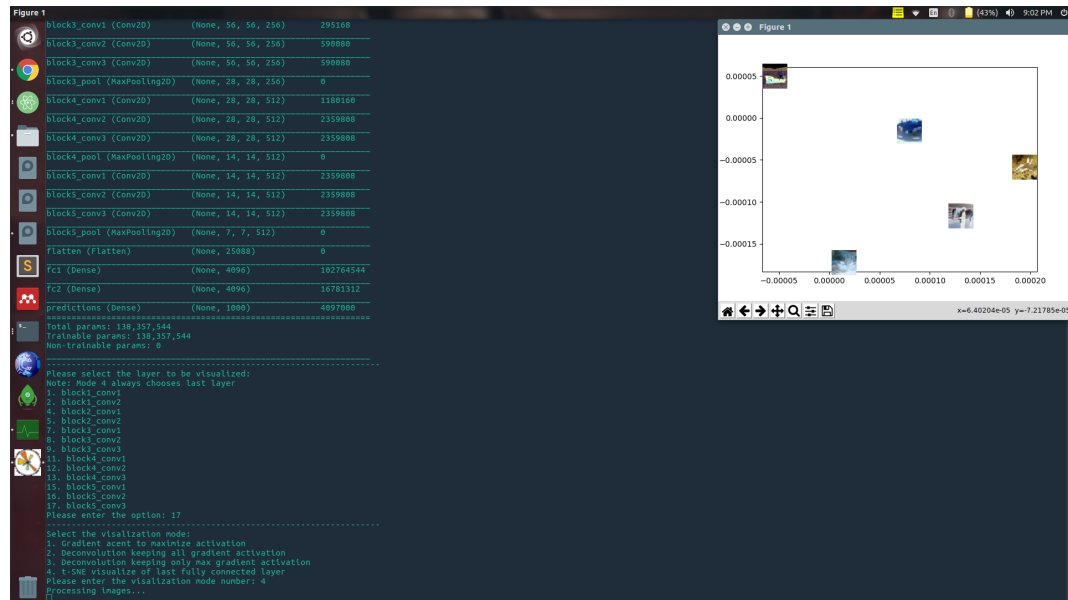Figure 11: Screenshot 1



Figure 12: Screenshot 2

Figure 13: Screenshot 3

## 7 CONCLUSION AND FUTURE REMAINING

The toolkit provides a simple interface for machine learning researchers to visualise any CNN model using the method suited best for their application.

Though we currently have deconvolution network counterpart for most of the commonly used layers, implementations of many more layers could be added to visualise more complicated models. Also, the t-SNE visualisation can be scaled up using a preprocessing the data with PCA or developing a multiprocessing algorithm.

## 8 CONTRIBUTION

Amey Agrawal: Development of custom keras layers for method 1 and implementation of method 2 in keras along with generation of t-SNE embedding in method 3.

Jaikumar Balani: Plotting and storing generated images using matplotlib.

## REFERENCES

Erhan, Dumitru, Bengio, Yoshua, Courville, Aaron, and Vincent, Pascal. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341:3, 2009.

Maaten, Laurens van der and Hinton, Geoffrey. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.

Simonyan, Karen, Vedaldi, Andrea, and Zisserman, Andrew. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.

Zeiler, Matthew D and Fergus, Rob. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pp. 818–833. Springer, 2014.