# Summary of ImageNet Classification with
# Deep Convolutional Neural Networks

## DIPANSHU PANDA(CE22B004)

The research paper "ImageNet Classification with Deep Convolutional Neural Networks" by Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton discusses the successful application of deep convolutional neural networks (CNNs) for image classification. Here is a summary of the key points from the paper:

## Abstract:

The authors trained a large, deep convolutional neural network (CNN) to classify 1.2 million high-resolution images from the ImageNet LSVRC-2010 contest into 1000 different classes. They achieved top-1 and top-5 error rates of 37.5% and 17.0%, respectively, which outperformed the previous state-of-the-art. The network consists of five convolutional layers followed by max-pooling layers and three fully-connected layers, ending with a 1000-way softmax. They used non-saturating neurons and efficient GPU implementations to speed up training. To combat overfitting, they employed dropout in the fully-connected layers. The model variant entered in the ILSVRC-2012 competition achieved a winning top-5 test error rate of 15.3%, compared to 26.2% by the second-best entry.

## Network Architecture:

The CNN described in the study has 60 million parameters and 650,000 neurons. The architecture includes:

- Five convolutional layers, some followed by max-pooling layers.
- Three fully-connected layers, with the final layer using a 1000-way softmax.
- Non-saturating neurons (Rectified Linear Units, or ReLUs) to enhance training speed.

## Methods to Reduce Overfitting:

Key techniques and innovations to reduce overfitting include:

- **Data Augmentation:** The easiest and most common method to reduce overfitting on image data is to artificially enlarge the dataset using label-preserving transformations. Additional training examples were generated through image translations, horizontal reflections, and intensity alterations of RGB channels. They employed two distinct forms of data augmentation, in which the transformed images are generated in Python code on the CPU while the GPU is training on the previous batch of images making the process computationally free-:

  The first form of data augmentation consisted of generating image translations and horizontal reflections. This increased the size of the training set by a factor of 2048, although the resulting training examples were, of course, highly interdependent. It is mentioned that without this scheme, the network suffers from substantial overfitting forcing the use of smaller networks.

  The second form of data augmentation consisted of altering the intensities of the RGB channels in training images. They performed PCA on the set of RGB pixel values throughout the ImageNet training set. To each training image, multiples of the found principal components were added with magnitudes proportional to the corresponding eigenvalues times a random variable drawn from a Gaussian with mean zero and standard deviation 0.1.

- **Dropout:** The authors employed the "dropout" technique for model combination, which consists of setting to zero for the output of each hidden neuron with a probability of 0.5. The neurons that are "dropped out" in this way do not contribute to the forward pass and do not participate in backpropagation. So every time an input is presented, the neural network samples a different architecture, but all these architectures share weights. This technique reduces the co-adaptations of neurons, since a neuron cannot rely on the presence of any other neuron particularly. It is, therefore,

forced to learn more robust features that are useful in conjunction with many different random subsets of the other neurons. They used dropout in the first two fully-connected layers of the network. It is also mentioned that without dropout the network exhibited substantial hence, making it necessary despite doubling the number of iterations to converge.

## Results

The CNN's performance was evaluated on the ImageNet LSVRC-2010 dataset, achieving a top-5 test error rate of 17%, significantly lower than the previous methods. The paper also demonstrates the network's effectiveness through qualitative evaluations, showing its ability to recognize objects in various settings and probing the network's visual knowledge through image similarity. The Euclidean distance measure was used over the L2, to account for objects placed in different poses. The depth of the network was crucial, as removing any convolutional layer resulted in a significant performance drop.

## Conclusion:

This research illustrates that deep CNNs can achieve state-of-the-art results in image classification tasks when trained with large datasets and using innovative techniques such as dropout. The network's success in the ILSVRC-2012 competition underscores the practical applicability of these models for real-world image recognition challenges. This paper is a seminal work in deep learning and computer vision, demonstrating the power of deep CNNs in handling large-scale image classification tasks.