# DOCUMENT CLASSIFICATION USING VGG-16 WITH ATTENTION
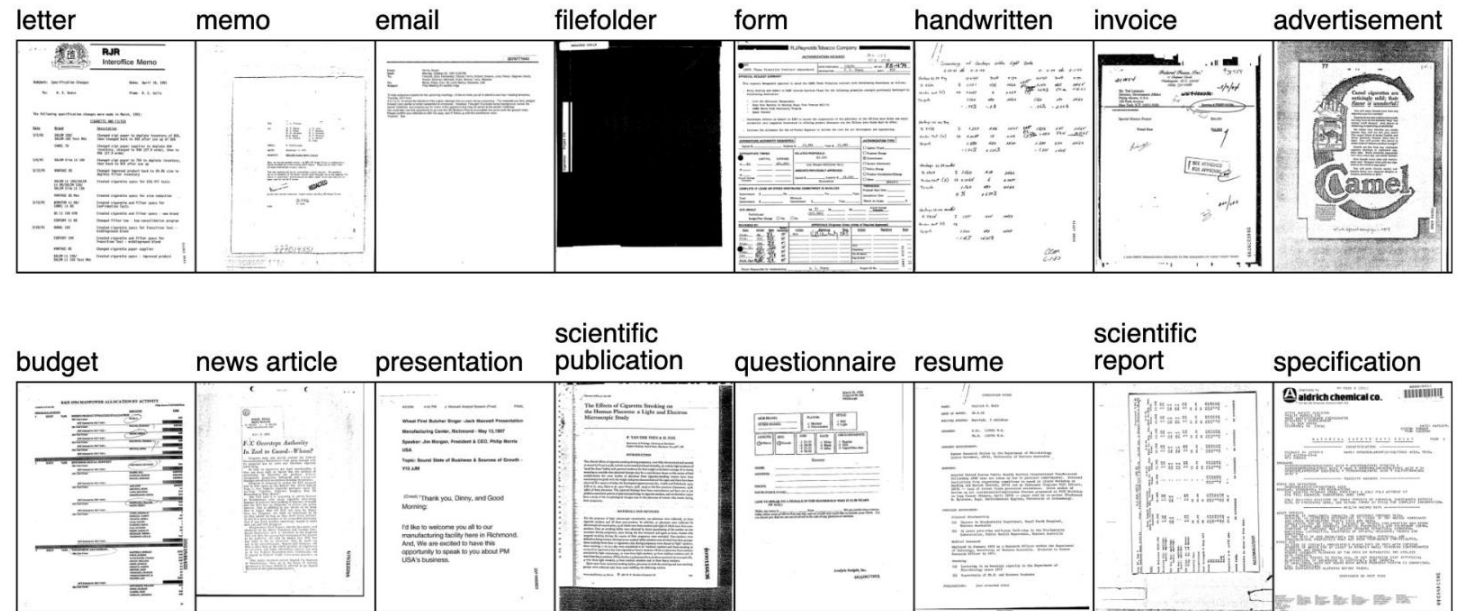
Submitted By Rishab Tomar (224161007)  Atul Bhagat(224151002)

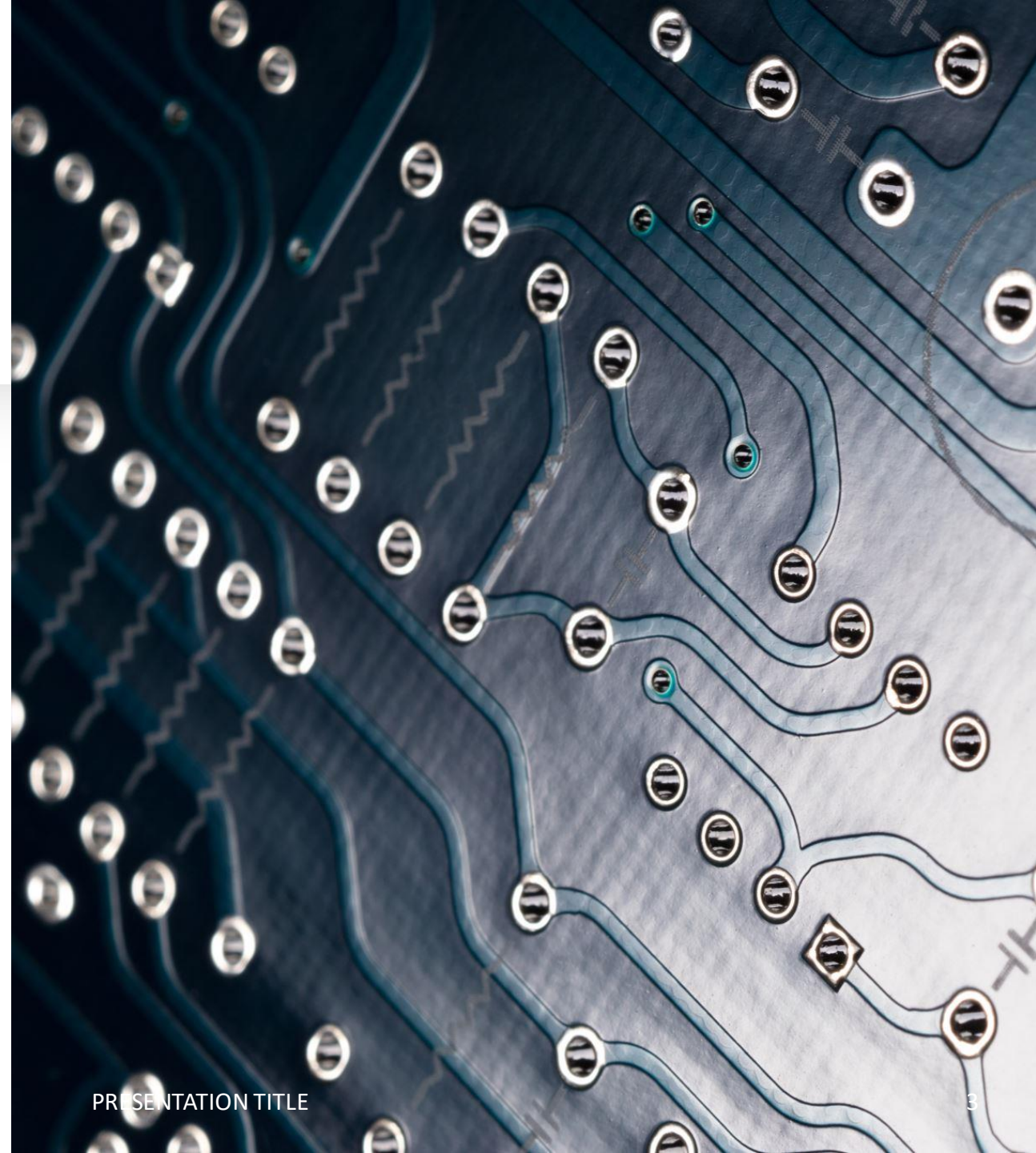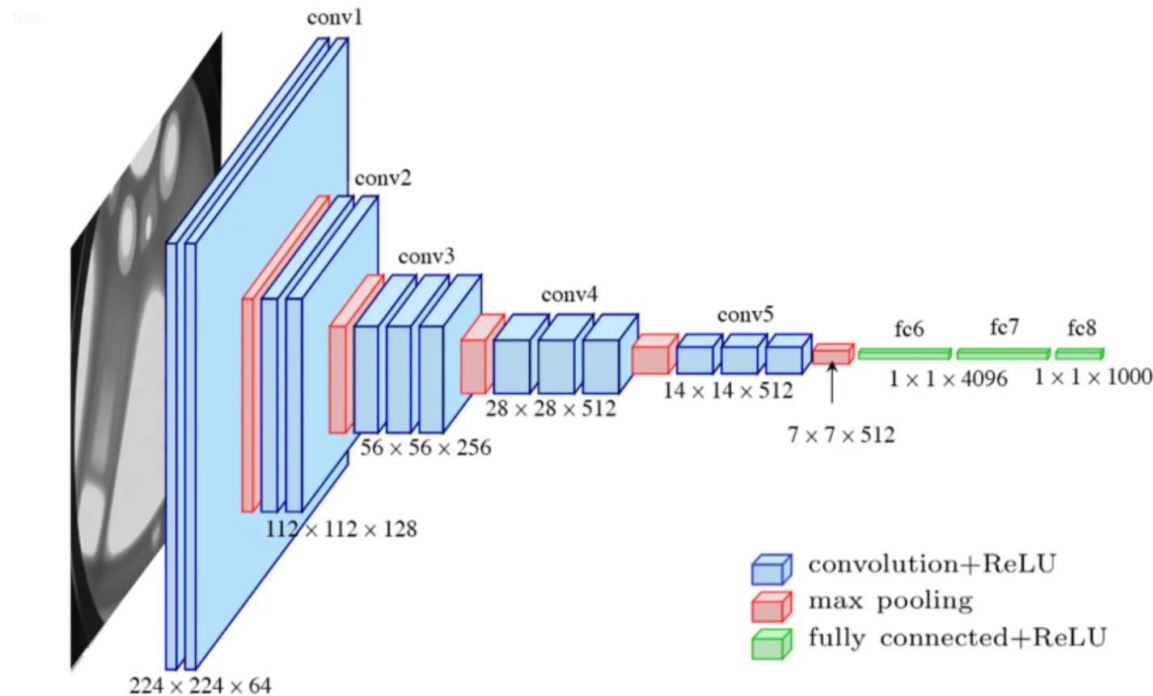Submitted To  Dr. Debanga Raj Neog

# ABOUT DATASET



- RVL-CDIP Dataset containing a total of 16 document categories including scientific articles ,handwritten letters. newspapers and more.

- It Contains a total of 400000 images 320000 as Training and 40000 as Validation and 40000 as Testing, but we have used only 16000 images with 1000 images per class to maintain a balance class distribution among classes.

# PROBLEM STATEMENT

- We aim to explore and evaluate the impact of incorporating an attention block into a VGG-16 network compared to a VGG-16 network without attention.

- The RVL-CDIP dataset used contains a diverse collection of images, and the task is to classify them into specific categories
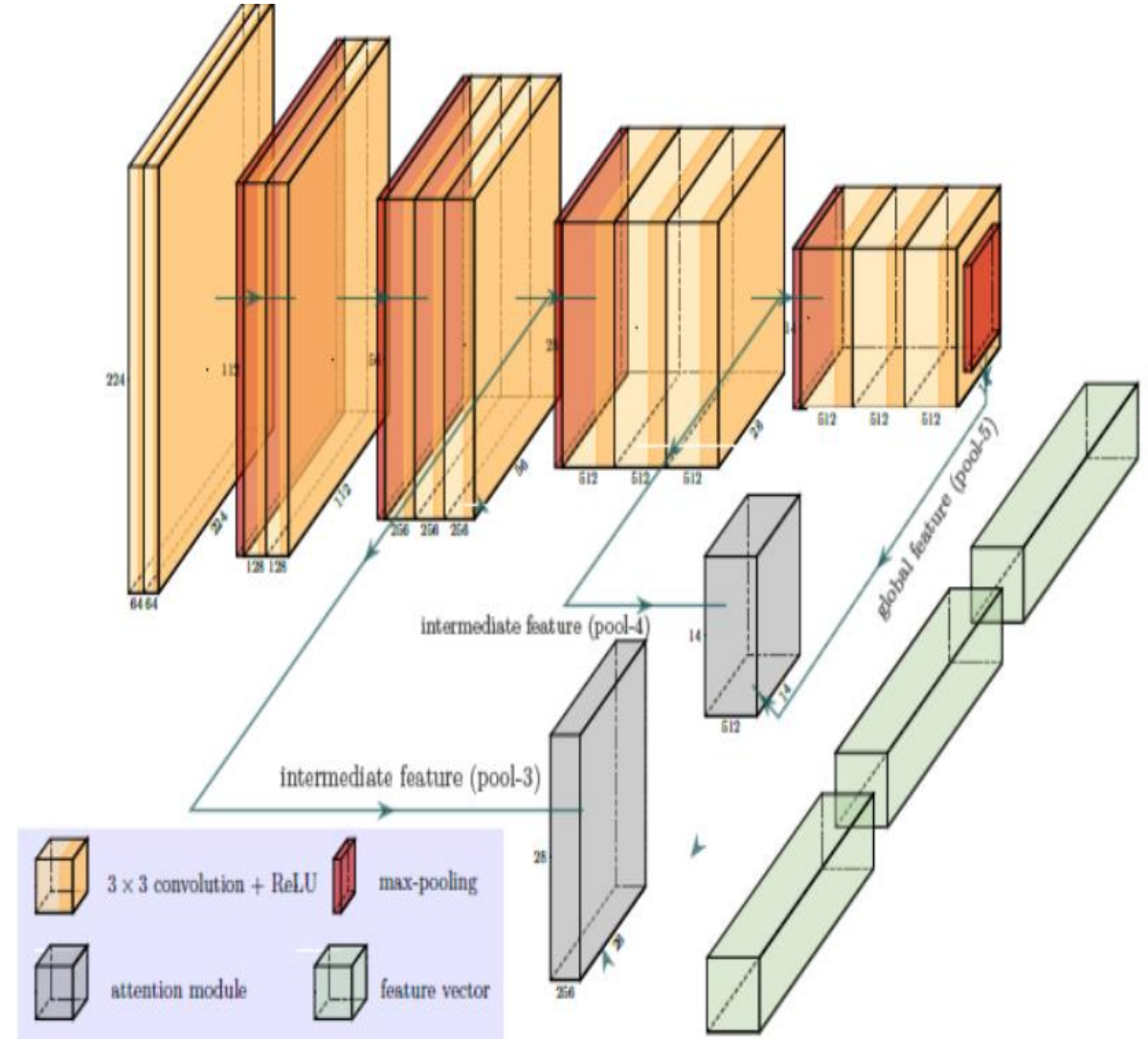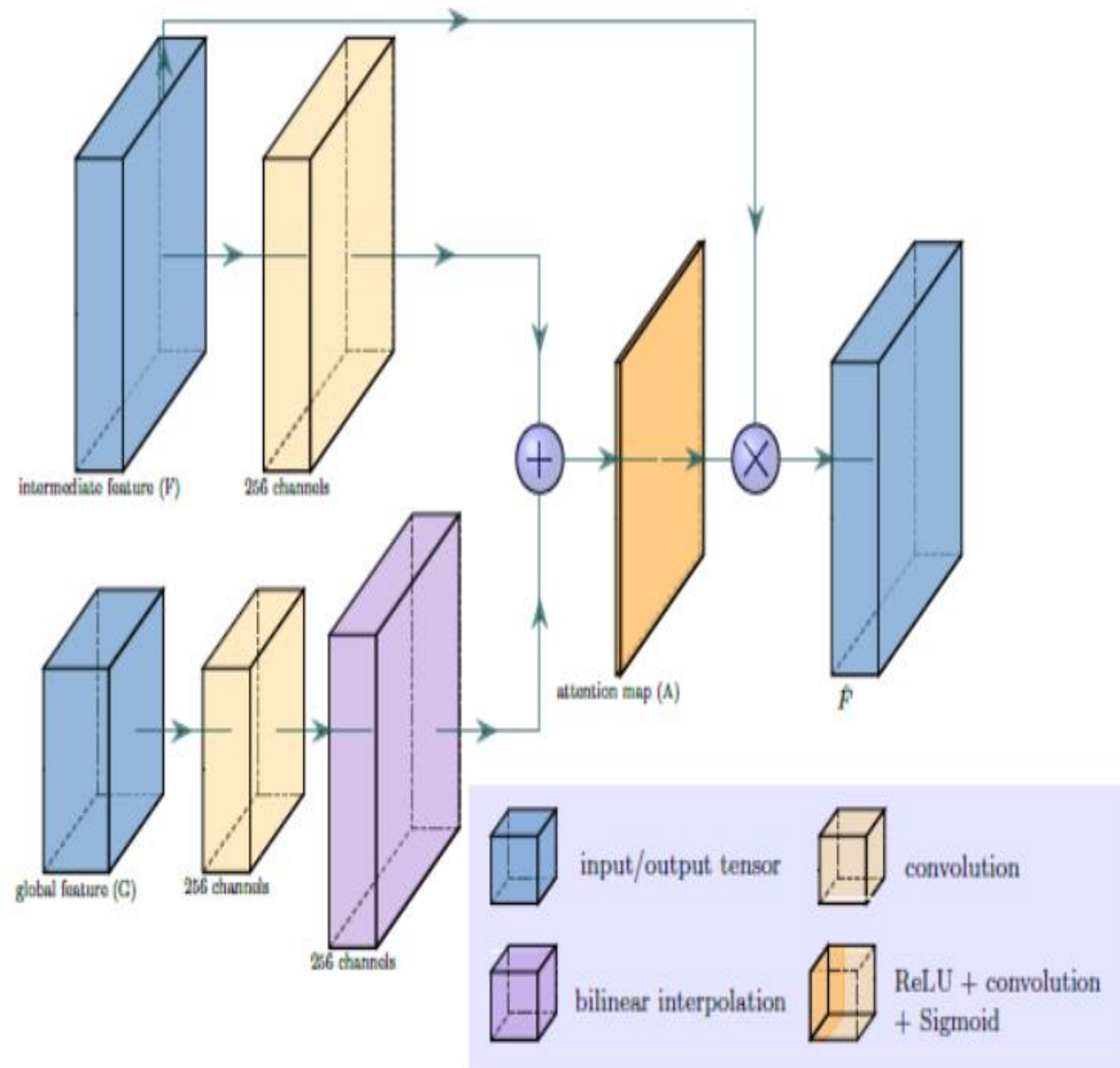
# MODEL-VGG16



- VGG-16, short for the Visual Geometry Group 16-layer network, is a deep convolutional neural network architecture designed for image classification and object recognition tasks.

- VGG-16 consists of 16 layers, including 13 convolutional layers and 3 fully connected layers. It is characterized by its use of small 3x3 convolutional filters and max-pooling layers, which result in a highly expressive and deep network

# MODEL- VGG16 WITH ATTENTION

- Two attention modules are applied (the gray blocks). The output of intermediate feature maps(pool-3 and pool-4) are used to infer attention maps. Output of pool-5 serves as a form of global-guidance because the last stage feature contains the most abstract and compressed information over the entire image.

- The three feature vectors (green blocks) are computed via global average pooling and are concatenated together to form the final feature vector, which serves as the input to the classification layer
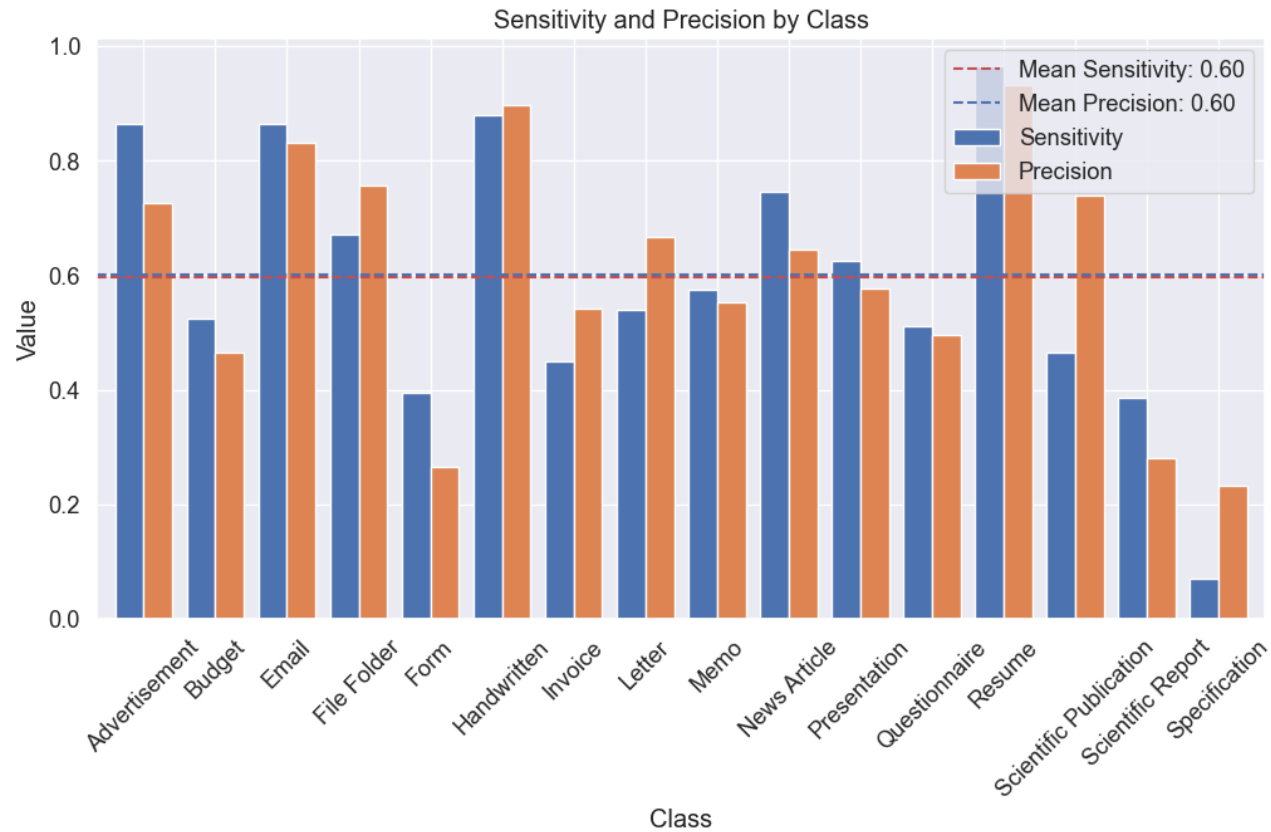
- The intermediate feature vector(F) is the output of pool-3 or pool-4 and the global feature vector (output of pool-5) is fed as input to the attention layer.

- Both the feature vectors pass through a convolution layer.

- If spatial size of global and intermediate features are different, feature upsampling is done via bilinear interpolation.

- The *up_factor* determines by what factor the convoluted global feature vector has to be upscaled.

- After that an element wise sum is done followed by a convolution operation that just reduces the 256 channels to 1.

- This is then fed into a Softmax layer, which gives us a normalized Attention map (A). Each scalar element in A represents the degree of attention to the corresponding spatial feature vector in F.

- The new feature vector $\hat{F}$ is then computed by *pixel-wise* multiplication. That is, each feature vector f is multiplied by the attention element a

- So, the attention map A and the new feature vector $\hat{F}$ are the outputs of the Attention Layer.

- The architecture of VGG16 is kept mostly the same except the Dense layers are removed.

- We pass pool-3 and pool-4 through the attention layer to get $\hat{F}3$ and $\hat{F}4$ .

- $\hat{F}3$ ,$\hat{F}4$ and G(pool-5) are concatenated and fed into the final classification layer.
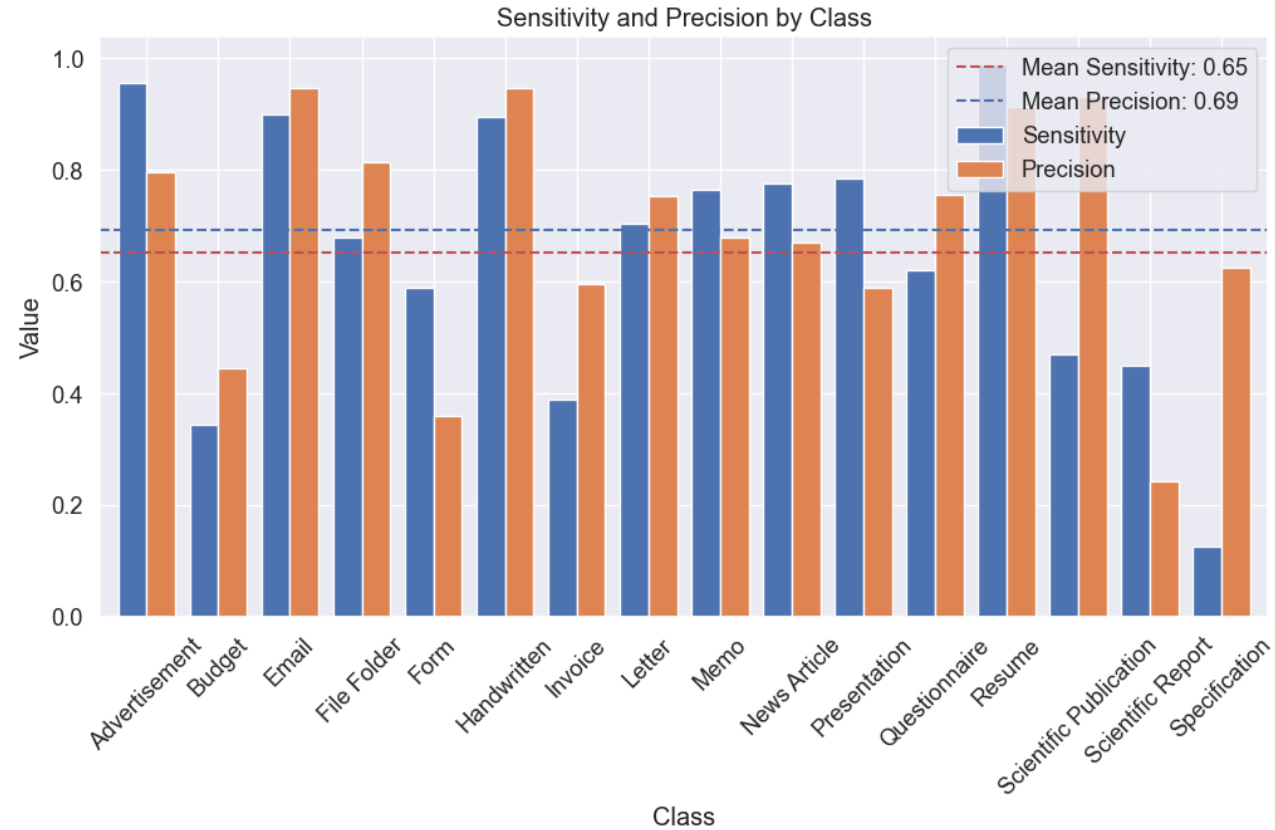
- The whole network is trained end-to-end.

# RESULTS

1. VGG-16



Sensitivity and Precision by Class



Confusion Matrix

# 2. VGG16 WITH ATTENTION

# RESULTS DISCUSSION

- VGG16(Model-1) shows relatively good sensitivity and precision for most classes, with an average sensitivity of 0.5956 and an average precision of 0.6004.

- VGG16 With Attention(Model-2), on the other hand, generally has higher sensitivity and precision values across the classes, resulting in a slightly better performance with an average sensitivity of 0.6525 and an average precision of 0.6918.

- The Percentage Increase in the mean sensitivity is 9.56%, and the percentage increase in the mean precision is 15.24% from Model-1 to Model-2.

- The Attention layer in Model 2, used in conjunction with VGG16, enhances performance by selectively emphasizing important image regions, improving precision, reducing overfitting, and providing interpretability to the model's predictions. It optimizes feature extraction and object recognition