

Precog Task Report - Analyzing Hateful Memes

C Swaroop - 2022101114

February 2024

1 Introduction

This report summarizes and lists down all of the findings, insights, analysis, and methodologies followed in the programming task.

2 Parts of the task I did:

1. Object detection to identify various elements and objects present in the images.
2. Cataloging the frequency and distribution of various objects present in the images.
3. Assessing the impact of hindrance of captions in memes while performing object detection.
4. Analyzing the effectiveness of object detection.
5. Meme-Classification system- A classifier which classifies whether an image is a meme or not.
6. Hateful Meme classification system- A computer vision-based classifier that classifies whether a meme is hateful or non-hateful.
7. NLP-based classification system-Classifies whether a meme is hateful or not only based on the text of the meme.
8. Percentage of hatefulness calculator for a text using NLP.
9. Text extractor system- Extracts text from memes using a pre-trained OCR model.

3 Parts of the task I did not do:

1. **Classification system based on the catalogue of distribution of objects and set of labels which tells meme is toxic obtained in the object detection task:**

I did not do this task since I have opted for another classifier which predicts whether a meme is hateful or not. I have also employed an NLP method which can calculate the toxicity score/hatefulness score based on the presence of certain words in the captions of the meme.

2. **Image processing techniques to filter text from images:**

I have stumbled across many pre-trained OCR models that can remove text from images coupled by inpainting technique, but I couldn't get the installation of the dependency libraries right. For instance, I couldn't set up TensorRT (NVIDIA-based library for inference optimization of deep learning models) with TensorFlow well enough for the OCR models to work. As a respite, I have found pre-trained models that can detect images from text and display them.

4 Report of Tasks:

4.1 Object Detection system and Caption Hinderance Assessment:

For this purpose, I have used a pre-trained model *YOLOV8* which is pre-trained on Microsoft's COCO dataset consisting of about 330,000 images with 80 object classes.

```
1 from ultralytics import YOLO
2 from PIL import Image, ImageDraw
3 from collections import Counter
4 import matplotlib.pyplot as plt
5 import os
6 import torch
7 from collections import defaultdict
```

These are the necessary libraries and dependencies needed for the task.

```
1 model=YOLO("yolov8n.pt")
2
3 results=model.train(data="coco128.yaml", epochs=100)
4
```

The above lines of code initializes a YOLO object detection model using the pre-trained weights specified in the file "yolov8n.pt". Model is instantiated and

stored in variable model. The results variable stores information or statistics about the training process, such as loss values, accuracy, or other metrics.

```
1  # Load the trained model
2  model = YOLO("runs/detect/train/weights/best.pt")
3
4  # Provide the path to the test image or directory containing test images
5  test_images_path = "dataset_hate/test/images/0"
6
7  # Perform object detection on the test images
8  results = model(test_images_path)
9  # Create a directory to save the results if it doesn't exist
10 save_dir = "results"
11 if not os.path.exists(save_dir):
12     os.makedirs(save_dir)
13
14 # Iterate through the results and save each image
15 for idx, r in enumerate(results):
16     im_array = r.plot() # plot a BGR numpy array of predictions
17     im = Image.fromarray(im_array[..., ::-1]) # RGB PIL image
18
19     # Save the image with a unique filename
20     save_path = os.path.join(save_dir, f"result_{idx+1}.jpg")
21     im.save(save_path) # save image
22
23 print("Results saved successfully.")
```

The above lines of code loads the trained model and path to the directory containing test images is provided and object detection is then applied. The results of the object detection is in a directory which consists of the test images with bounding boxes around each detected object label in the images.

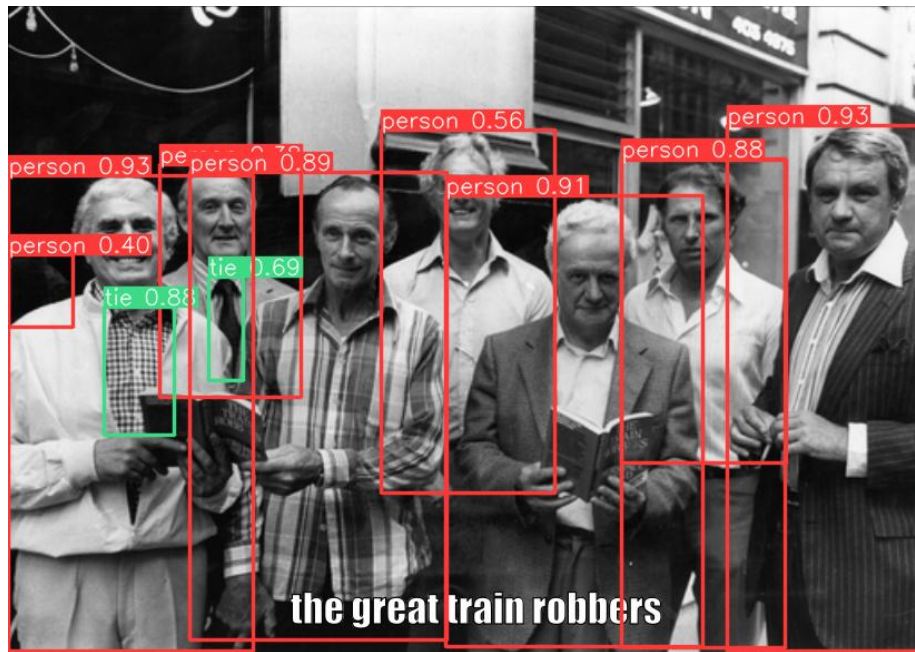


Figure 1: Example Result Image

```

1  class_confidences = defaultdict(float)
2  class_counts = defaultdict(int)
3
4
5  # Initialize a counter to store the frequency of each class label
6  class_counter = Counter()
7
8  for i in range(0, len(results)):
9      for box in results[i].boxes:
10         class_id=int(box.cls)
11         class_label=results[i].names[class_id]
12         class_counter.update([class_label])
13         confidence = float(box.conf)
14         class_confidences[class_label] += confidence
15         class_counts[class_label] += 1
16
17  labels, frequencies = zip(*class_counter.items())
18
19
20  # Calculate the average confidence for each class label
21  average_confidences = {label: total_confidence / class_counts[label] for label, total_confidence in class_confidences.items()}
22  overall_average_confidence = sum(class_confidences.values()) / sum(class_counts.values()) if sum(class_counts.values()) > 0 else 0
23  frequency_distribution = dict(class_counter)
24  print("Frequency distribution of all class labels:")
25  print(frequency_distribution)
26

```



```

8         class_label = results[i].names[class_id]
9         class_counter.update([class_label])
10        confidence = float(box.conf)
11        class_confidences[class_label] += confidence
12        class_counts[class_label] += 1
13
14    # Extract the class labels and their frequencies
15    labels, frequencies = zip(*class_counter.items())
16
17    # Calculate the average confidence for each class label
18    average_confidences = {label: total_confidence / class_counts[label] for label, total_confidence in class_confidences.items()}
19
20    overall_average_confidence = sum(class_confidences.values()) / sum(class_counts.values()) if sum(class_counts.values()) > 0 else 0
21
22    # Create a frequency distribution dictionary of all class labels
23    frequency_distribution = dict(class_counter)
24    print("Frequency distribution of class labels predicted with more than 0.5 certainty:")
25    print(frequency_distribution)

```

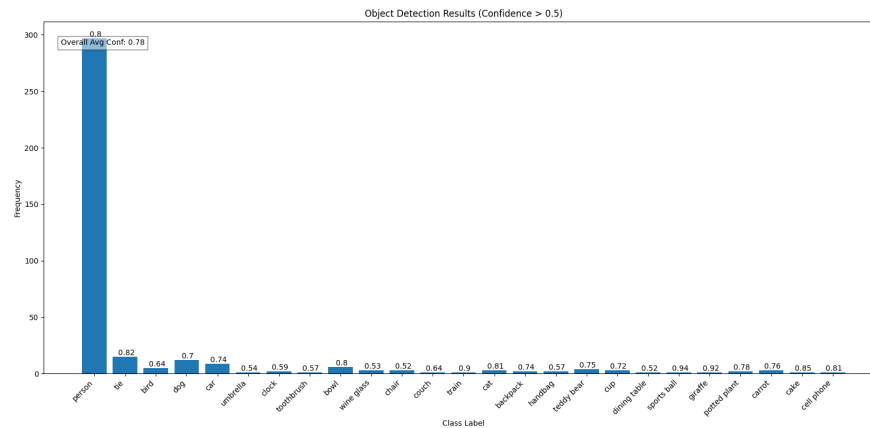


Figure 3: Catalogue plot with confidence values greater than 0.5

For the above plot I only considered plotting the instances of object detections which have a confidence value of greater than 0.5. The reason I went ahead with case of considering case of confidence values greater than 0.5 is that the model predicts with more certainty and the detections are more likely to be true here.

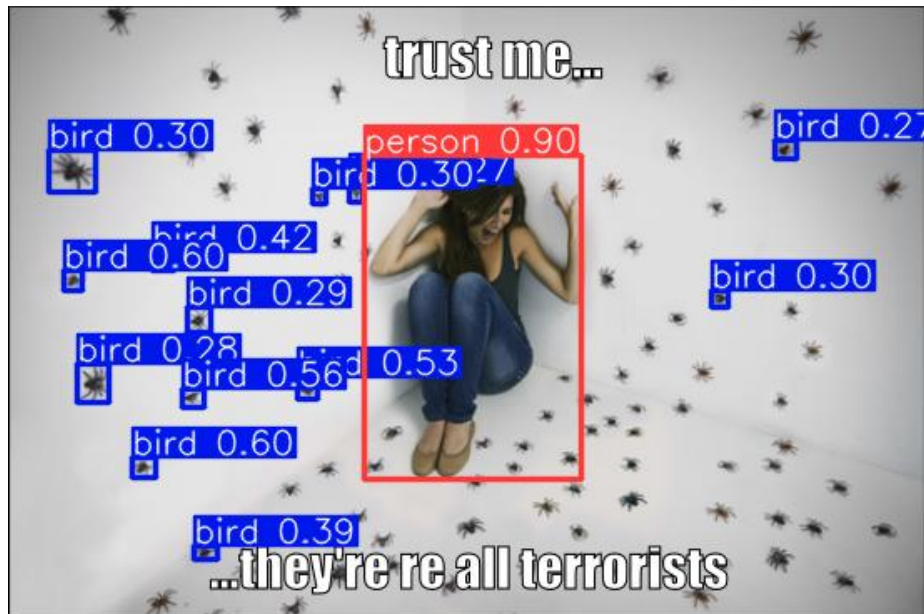


Figure 4:

For example in the above image, the objects detected as *bird* are detected with a low confidence values. In reality they are spiders and not birds!

I have also plotted catalogues of object classes distribution when I applied the object detection system to the same set of testing images before with their text removed.

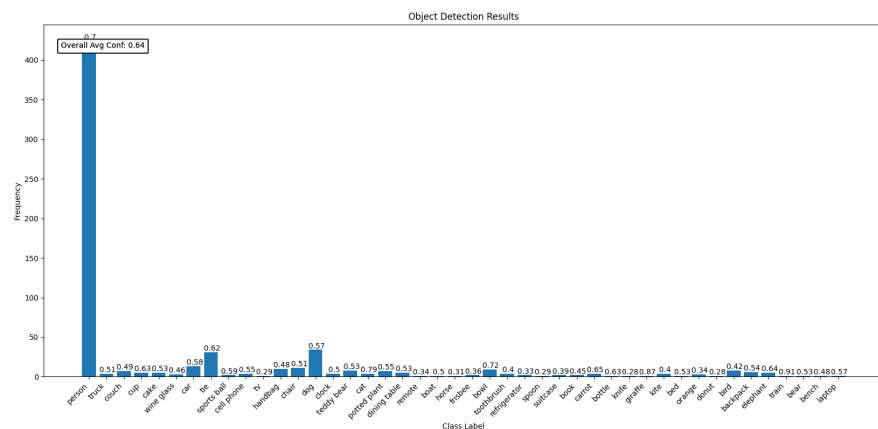


Figure 5: Catalogue of images with their text removed

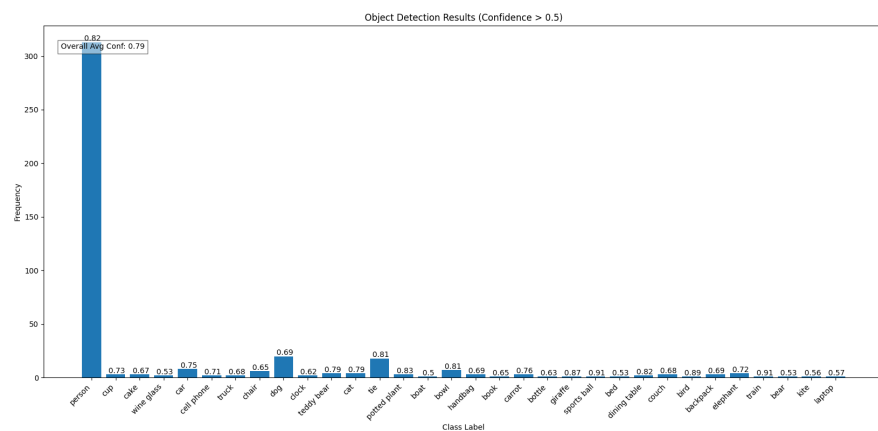


Figure 6: Catalogue of images with their text removed and confidence values greater than 0.5

It can be seen from the above plots that the frequency of distribution of various object classes have changed after removing the text from the images. New classes have also been detected this time. As predicted, the average confidence values and the overall average confidence of all classes have increased with the removing of text from the images.

Caption Hinderance Assessment: Clearly from the above plots, the text present in the images hinders the process of object detection done on them as evident by the rise of average confidence values and increase of cases of object detections and presence of new class labels. We shall consider one example here:

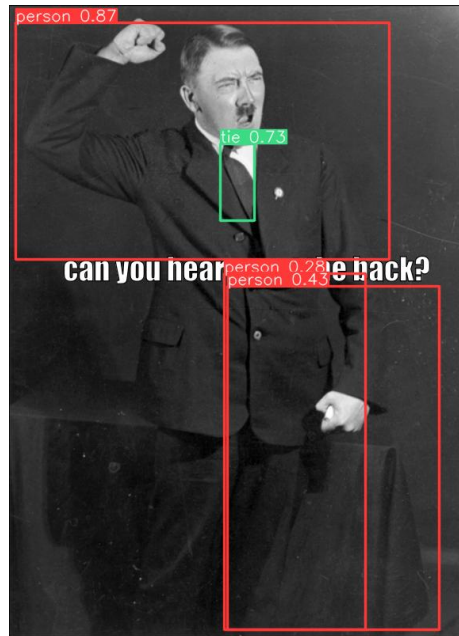


Figure 7:

From the above figure it's evident that only one person is present, but the object detector makes a claim of three people which is evidently wrong. Note that the extra two claims of people are made with less confidence values-0.28 and 0.43.

Now consider the case when the text is removed from the picture and object detection is then performed:

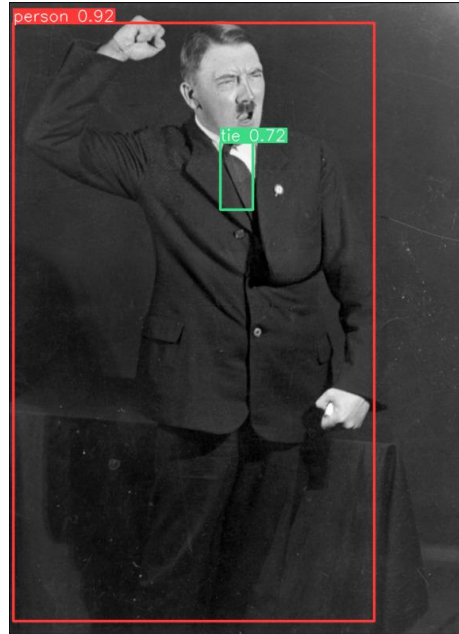


Figure 8:

Now the object detector correctly predicted the presence of only one person in the image. In the original meme the text right at the middle of the image caused ambiguities for the object detector system which wrong detected the presence of three people in the image.

Thus the shortcomings of the object detection system applied to the original set of memes with text can be attributed to the presence of captions which hinder the object detection process.

4.2 Meme-Classifier:

This classifier predicts whether or not an image is a meme or not a meme.

```
1 import tensorflow as tf
2 import os
3 from tensorflow import keras
```

```

4 from tensorflow.keras import layers, optimizers, Sequential
5 from tensorflow.keras.layers import Activation , Dropout , Conv2D, MaxPooling2D, Dense, Flatten
6 import matplotlib.pyplot as plt
7 from tensorflow.keras.models import load_model
8 from tensorflow.keras.preprocessing.image import load_img, img_to_array
9 import numpy as np
10 from tensorflow.keras.preprocessing.image import ImageDataGenerator

```

Above are the required libraries and dependencies required for the task.

```

1 model = Sequential()
2 model.add(Conv2D(32, (3,3), input_shape = (64,64,3), activation = 'relu'))
3 model.add(MaxPooling2D(pool_size=(2,2)))
4 model.add(Conv2D(64, (3, 3)))
5 model.add(Activation('relu'))
6 model.add(MaxPooling2D(pool_size=(2, 2)))
7 model.add(Conv2D(64, (3, 3)))
8 model.add(Activation('relu'))
9 model.add(MaxPooling2D(pool_size=(2, 2)))
10 model.add(Flatten())
11 model.add(Dense(256, activation='relu'))
12 model.add(Dropout(0.2))
13 model.add(Dense(128, activation='relu'))
14 model.add(Dropout(0.2))
15 model.add(Dense(1, activation='sigmoid'))
16 model.summary()

```

This code defines a convolutional neural network (CNN) model using TensorFlow's Keras API.

```

1 train_datagen = ImageDataGenerator(
2     rescale=1. / 255,
3     rotation_range=30,
4     zoom_range = 0.15,
5     width_shift_range=0.10,
6     height_shift_range=0.10,
7     horizontal_flip=True)
8 test_datagen = ImageDataGenerator(rescale=1./255)
9 training_set = train_datagen.flow_from_directory(
10     'meme_classifier/training_set',
11     target_size=(64,64),
12     batch_size=15,
13     class_mode='binary')
14 test_set = test_datagen.flow_from_directory(
15     'meme_classifier/test_set',
16     target_size=(64,64),
17     batch_size=15,
18     class_mode='binary')
19
20 with tf.device('/GPU:0'):
21     history = model.compile( optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
22     model.fit(
23         training_set,

```

```

24         steps_per_epoch=991//15,
25         epochs=30,
26         validation_data=test_set,
27         validation_steps=296//15
28     )

```

This code is responsible for training a convolutional neural network (CNN) model for a binary classification task using TensorFlow and Keras.

Note that The training process is performed on a GPU ('/GPU:0'), which can significantly speed up the training process for deep neural networks.

The paths to the training dataset and the testing dataset are also mentioned in the code - `meme_classifier/training_set` and `meme_classifier/test_set`.

```

1 model.save('meme_classifier/model.h5')

```

This line of code saves the model as `model.h5` inside the `meme_classifier` directory.

```

1 classifier = load_model('meme_classifier/model.h5')
2 path = 'dataset_hate/test/images/0/05438.png'
3 img_original = load_img(path)
4 img = load_img(path, target_size = (64,64))
5 img_tensor = img_to_array(img)
6 img_tensor = np.expand_dims(img_tensor, axis = 0)
7 img_tensor/=255.0
8 pred = classifier.predict(img_tensor)
9 print(pred)
10 if pred<.5: str = '-----Meme-----'
11 else: str = '-----Non meme-----'
12 plt.imshow(img_original)
13 plt.axis('off')
14 plt.title(str)
15 plt.show()

```

The above lines of code loads the meme classifier model and path to image for which prediction is to be done is loaded. The next few lines does preprocessing on the image and prediction is then performed on it telling whether the image is a meme or not a meme.

Below are some example outputs of the meme-classifier-



Figure 9:



Figure 10:

```

1 loaded_model = load_model('meme_classifier/model.h5')
2
3 test_data_dir = 'meme_classifier/test_set'
4
5 # Load the test dataset
6 test_datagen = ImageDataGenerator(rescale=1./255)
7 test_set = test_datagen.flow_from_directory(
8     test_data_dir,
9     target_size=(64, 64),
10    batch_size=15, # Adjust batch size as needed
11    class_mode='binary')
12
13 # Evaluate the model on the test data
14 evaluation = loaded_model.evaluate(test_set)
15
16
17 # Print the performance metrics
18 print("Test Accuracy:", evaluation[1])

```

The above lines of code loads the trained model and tests the model for performance and accuracy. Here it assumes that the test dataset is organized in subdirectories where each subdirectory represents a class, i.e meme or not a meme. It automatically infers the classes from the subdirectories and assigns binary labels (0 or 1) based on the classmode parameter. The evaluate method computes the loss and accuracy metrics of the model on the provided dataset.

```

1 Test Accuracy: 0.8484849333763123

```

The accuracy of 84 percent for the meme classifier is a pretty good measure of the model's performance.

Summarizing: The meme classifier uses a special type of Deep Learning technique called a Convolutional Neural Network (CNNs) to figure out if a picture is a meme or not. It looks at lots of example pictures that have already been labeled, learning what memes usually look like. During training, it's taught to recognize important patterns in the pictures that help it make the right guesses. After training, it's tested on new pictures to see how well it can tell memes apart from other pictures. Once it's trained and tested, it can quickly classify new pictures as memes or non-memes, helping to sort content efficiently. Techniques of Data augmentation and Image preprocessing such as Techniques such as rotation, zooming, shifting, and horizontal flipping are carried out to images in the process of classification. After training, the model's performance is evaluated on a separate testing dataset. The model's accuracy and loss metrics are computed to assess its performance on unseen data.

4.3 Hateful meme classifier using CV and NLP:

4.3.1 CV:

```
1 import tensorflow as tf
2 import os
3 from tensorflow import keras
4 from tensorflow.keras import layers, optimizers, Sequential
5 from tensorflow.keras.layers import Activation, Dropout, Conv2D, MaxPooling2D, Dense, Flatten
6 import numpy as np
7 import pandas as pd
8 from sklearn.model_selection import train_test_split
9 from tensorflow.keras.preprocessing.image import img_to_array, load_img
10 from tensorflow.keras.preprocessing.image import ImageDataGenerator
```

These are the required libraries and dependencies required for the task.

```
1 train_df = pd.read_json('train.jsonl', lines=True)
```

The above line of code puts the training dataset which is referenced in the train.jsonl file into a pandas DataFrame df which can be used for further processing and analysis.

```
1 {"id":42953,"img":"img/42953.png","label":0,"text":"its their character not their color that matters"}
2 {"id":23058,"img":"img/23058.png","label":0,"text":"don't be afraid to love again everyone is not like your ex"}
```

The train.jsonl file is structured such that each line represents a JSON object. Each object contains the following elements:

- id: The identifier of the meme-image.
- path: The path to the meme image.
- label: Signifies whether the meme is hateful or not, with 0 representing Non-Hateful and 1 representing Hateful.
- text: The text associated with the meme.

Note that the training dataset consists of 8000 objects.

```
1 # Preprocess the data
2 def preprocess_data(df):
3     image_data = []
4     labels = []
5     for index, row in df.iterrows():
6         img = load_img(row['img'], target_size=(64, 64))
7         img = img_to_array(img) / 255.0 # Normalize the pixel values
8         image_data.append(img)
9         label = 1 if row.get('label') == 1 else 0 # If label exists, use it; otherwise, default to 0
10        labels.append(label)
11    return np.array(image_data), np.array(labels)
```

The above lines of code describe a function which takes a DataFrame containing image data and labels, loads and preprocesses the images, and returns the preprocessed data ready for training a machine learning model.

```
1 # Preprocess the data
2 X, y = preprocess_data(train_df)
```

The array of image data and labels from the training dataframe after preprocessing are stored in X and y respectively.

```
1 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=42)
2
```

This line of code splits the preprocessed data into training and validation sets in the ratio of 80:20 for training and validation respectively. Here, `X_train` and `y_train` contain the training image data and corresponding labels, respectively. `X_val` and `y_val` contain the validation image data and corresponding labels, respectively.

```
1 # Define your model
2 model = Sequential()
3 model.add(Conv2D(32, (3, 3), input_shape=(64, 64, 3), activation='relu'))
4 model.add(MaxPooling2D(pool_size=(2, 2)))
5 model.add(Conv2D(64, (3, 3)))
6 model.add(Activation('relu'))
7 model.add(MaxPooling2D(pool_size=(2, 2)))
8 model.add(Conv2D(64, (3, 3)))
9 model.add(Activation('relu'))
10 model.add(MaxPooling2D(pool_size=(2, 2)))
11 model.add(Flatten())
12 model.add(Dense(256, activation='relu'))
13 model.add(Dropout(0.2))
14 model.add(Dense(128, activation='relu'))
15 model.add(Dropout(0.2))
16 model.add(Dense(1, activation='sigmoid'))
17 model.summary()
```

This code defines a convolutional neural network (CNN) model and its architecture using TensorFlow's Keras API.

```
1 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
```

This line of code configures the model for training, specifying the optimization algorithm, loss function, and evaluation metrics to be used during the training process.

```
1 history = model.fit(X_train, y_train, epochs=50, batch_size=15, validation_data=(X_val, y_val))
```

This line of code trains the model on the training data, using the specified number of epochs, batch size, and validation data, while storing the training history for later analysis.

```
1 import matplotlib.pyplot as plt
2
3 # Get training history
4 train_loss = history.history['loss']
5 val_loss = history.history['val_loss']
6 train_acc = history.history['accuracy']
7 val_acc = history.history['val_accuracy']
8 epochs = range(1, len(train_loss) + 1)
9
10 # Plot training and validation loss
11 plt.figure(figsize=(12, 6))
12 plt.subplot(1, 2, 1)
13 plt.plot(epochs, train_loss, 'b', label='Training loss')
14 plt.plot(epochs, val_loss, 'r', label='Validation loss')
15 plt.title('Training and Validation Loss')
16 plt.xlabel('Epochs')
17 plt.ylabel('Loss')
18 plt.legend()
19
20 # Plot training and validation accuracy
21 plt.subplot(1, 2, 2)
22 plt.plot(epochs, train_acc, 'b', label='Training accuracy')
23 plt.plot(epochs, val_acc, 'r', label='Validation accuracy')
24 plt.title('Training and Validation Accuracy')
25 plt.xlabel('Epochs')
26 plt.ylabel('Accuracy')
27 plt.legend()
28
29 plt.tight_layout()
30 plt.show()
```

This code segment uses Matplotlib to visually represent the performance of the model across training epochs. It extracts and plots training and validation loss, as well as training and validation accuracy, on separate subplots. The visualization helps understand how the model's performance evolves during training and validation processes, facilitating analysis and potential improvements especially in choosing the number of epochs for training to avoid the cases of underfitting and overfitting.

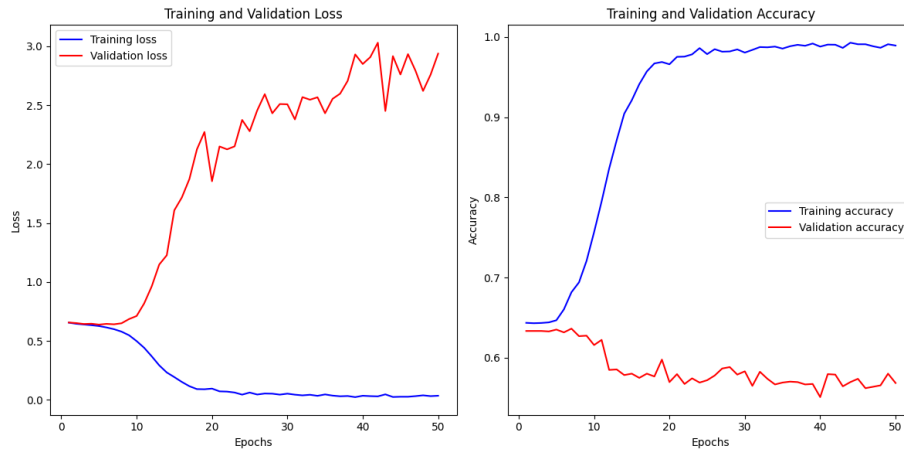


Figure 11:

The provided plots depict how the model's performance evolves over time, showcasing both loss and accuracy metrics during the training process.

Over time the model's training accuracy increases by learning more of the training data.

However the number of epochs chosen for training process should be chosen in a way so as to ensure the validation accuracy doesn't fall down drastically or remain a plateau over time.

Increasing the number of epochs to a really high number might cause issues of overfitting where the model completely memorises the training data and doesn't generalize well to unseen data and also having less number of epochs might result in underfitting where the model didn't learn enough in the training process resulting in high bias.

There should be a sweet spot of number of epochs to be chosen for training to simultaneously avoid the problems of both underfitting and overfitting.

There are various techniques to choose for selecting a good value of number of epochs for training such as Validation Curve Analysis, Early Stopping, Cross-Validation and even Experimentation and Iteration would work fine. Just choose the number of epochs such that the validation accuracy doesn't fall down drastically over time.

Note:For our task, it didn't really cause much of difference in accuracies in altering number of epochs anywhere between 10 to 50.

```

1 # Save the model
2 model.save('model.h5')

```

Model is saved into the directory as *model.h5*

```
1 # Load the saved model
2 model = tf.keras.models.load_model('model.h5')
```

Model is loaded and is ready for use.

```
1 # Load your test data from test.jsonl
2 test_df = pd.read_json('test_seen.jsonl', lines=True)
```

This code reads the test data from the specified JSONL file and stores it in the `test_df` DataFrame, allowing for further analysis and processing of the test dataset in Python.

```
1 {"id": "16395", "img": "img/16395.png", "label": 1, "text": "handjobs sold seperately"}
2 {"id": "37405", "img": "img/37405.png", "label": 1, "text": "introducing fidget spinner for women"}
```

The `test_seen.jsonl` file is structured such that each line represents a JSON object. Each object contains the following elements:

- id: The identifier of the meme-image.
- label: Signifies whether the meme is hateful or not, with 0 representing Non-Hateful and 1 representing Hateful.
- path: The path to the meme image.
- text: The text associated with the meme.

Note that the testing dataset consists of 1000 objects. The json object structure of the training dataset jsonl file is identical that of the testing dataset jsonl file.

```
1 # Preprocess the test data
2 def preprocess_test_data(df):
3     image_data = []
4     image_paths = []
5     for index, row in df.iterrows():
6         img = load_img(row['img'], target_size=(64, 64))
7         img = img_to_array(img) / 255.0 # Normalize the pixel values
8         image_data.append(img)
9         image_paths.append(row['img'])
10    return np.array(image_data), image_paths
```

This function preprocesses the test data by loading and normalizing images and returns the preprocessed image data along with their respective paths for further analysis or inference.

```
1 # Preprocess the test data
2 X_test, image_paths = preprocess_test_data(test_df)
```

This line of code preprocesses the test data and stores the preprocessed image data in `X_test` and the corresponding image paths in `image_paths` for subsequent use.

```
1 # Make predictions on the test data
2 predictions = model.predict(X_test)
```

This line of code generates predictions on the test data using the trained model and stores the predicted outputs in the `predictions` variable for further analysis

```
1 # Decode predictions
2 predicted_labels = ['Hateful' if pred > 0.5 else 'Non-Hateful' for pred in predictions]
```

This line of code translates the model predictions into human-readable labels, making it easier to interpret and analyze the model's output.

```
1 # Print predictions
2
3
4 output_file = "meme_predictions.txt"
5
6
7 # Print predictions to the external file
8 with open(output_file, "w") as f:
9     for i in range(len(image_paths)):
10         f.write(f"Image: {image_paths[i]}, Predicted Label: {predicted_labels[i]}\n")
11
12 print(f"Predictions saved to {output_file}")
```

Predictions of the model are saved to a file `meme_predictions.txt`.

```
1 Image: img/16395.png, Predicted Label: Non-Hateful
2 Image: img/37405.png, Predicted Label: Non-Hateful
3 Image: img/94180.png, Predicted Label: Non-Hateful
4 Image: img/54321.png, Predicted Label: Hateful
```

Glimpse of contents of `meme_predictions.txt`

Summarizing: This covers the process of training a convolutional neural network (CNN) model to classify memes as either 'Hateful' or 'Non-Hateful'. It starts by loading and preparing the training data, then splits it into training and validation sets. The model architecture is defined and trained using the training data, with progress displayed using Matplotlib. The trained model is saved for later use. Test data is loaded, processed, and used to make predictions with the trained model. Predictions are converted to readable labels and saved to a text file.

Food for thought: This model predicts whether a meme is hateful or not. It is trained with sufficient examples given their labels. But don't you think this classifier is biased towards one modal-Image and leaves out the text? i.e A meme is a multi-modal source of data containing two forms of data-Image and the text captions. This model classifies by having a bias to the image while training. The classifier fails to be a multi-modal model and instead behaves like a uni-modal model on image classification. Let's look at the NLP classifier version:

4.3.2 NLP:

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.svm import LinearSVC
5 from sklearn.metrics import classification_report
6 import nltk
7 from nltk.corpus import stopwords
8 from nltk.tokenize import word_tokenize
9 from nltk.stem import WordNetLemmatizer
10 import re
11 import json
12 from sklearn.metrics import accuracy_score
```

The required libraries and dependencies needed for the task.

```
1 # Download NLTK resources
2 nltk.download('stopwords')
3 nltk.download('punkt')
4 nltk.download('wordnet')
```

The NLTK resources such as stopwords, tokenizers, and lemmatizers are downloaded for text preprocessing.

```
1 # Preprocessing setup
2 stop_words = set(stopwords.words('english'))
3 lemmatizer = WordNetLemmatizer()
4 # Preprocessing function
5 def preprocess_text(text):
6     text = re.sub(r'[\W\s]', '', text)
7     tokens = word_tokenize(text)
8     tokens = [word.lower() for word in tokens if word.isalpha()]
9     tokens = [lemmatizer.lemmatize(word) for word in tokens if word not in stop_words]
10    return ' '.join(tokens)
```

Stop words are removed, and words are lemmatized using NLTK. A function `preprocess_text` is defined to preprocess the text data by removing punctuation, tokenizing, converting to lowercase, removing stop words, and lemmatizing.

```

1  # Load JSON lines file for testing
2  with open("test_seen.jsonl", "r") as f:
3      test_json_lines = f.readlines()
4  # Parse JSON lines for testing
5  test_data = []
6  for line in test_json_lines:
7      item = json.loads(line)
8      test_data.append(item['text'])
9  # Load JSON lines file
10 with open("train.jsonl", "r") as f:
11     train_json_lines = f.readlines()
12 # Parse JSON lines for training
13 train_data = []
14 for line in train_json_lines:
15     item = json.loads(line)
16     train_data.append((item['text'], item['label']))
17
18 # Convert training data to DataFrame
19 train_df = pd.DataFrame(train_data, columns=['text', 'label'])
20
21 # Apply preprocessing to training data
22 train_df['clean_text'] = train_df['text'].apply(preprocess_text)

```

Training and testing data are loaded from JSONL files and parsed. Training data is further converted into a DataFrame, and the text is preprocessed using the defined function.

```

1  # Vectorization
2  vectorizer = TfidfVectorizer(max_features=5000)
3  X_train_vec = vectorizer.fit_transform(train_df['clean_text'])
4  y_train = train_df['label']

```

```

1  # Vectorization
2  vectorizer = TfidfVectorizer(max_features=5000)
3  X_train_vec = vectorizer.fit_transform(train_df['clean_text'])
4  y_train = train_df['label']

```

The preprocessed text data is vectorized using TF-IDF vectorization, which converts text data into numerical feature vectors. For the uninitiated, TF-IDF (Term Frequency-Inverse Document Frequency) vectorization is a popular technique used in natural language processing to convert textual data into numerical feature vectors. This method learns the importance of words in sentences and depending on the sentence between hateful or not assigns certain weights to words. Essentially the code identifies words (based on their frequencies with which they appear in respective labels of sentences) that are strongly associated with each class label. These words are considered indicative of the content that the classifier identifies as hateful or non-hateful.

```

1 classifier = LinearSVC()
2 classifier.fit(X_train_vec, y_train)

```

A Linear Support Vector Classifier (LinearSVC) is created and trained on the TF-IDF vectorized features (X_train_vec) of the training text data (y_train). This classifier can now be used for predictions.

```

1 from wordcloud import WordCloud
2 import matplotlib.pyplot as plt
3 # Get the feature names from the TfidfVectorizer
4 feature_names = vectorizer.get_feature_names_out()
5
6 # Get the coefficients (weights) from the trained LinearSVC model
7 coefficients = classifier.coef_[0]
8
9 # Create a dictionary with feature names as keys and coefficients as values
10 word_coef = dict(zip(feature_names, coefficients))
11
12 # Separate words with positive coefficients (indicating association with hateful class)
13 hateful_words = {word: coef for word, coef in word_coef.items() if coef > 0}
14
15 # Separate words with negative coefficients (indicating association with non-hateful class)
16 non_hateful_words = {word: coef for word, coef in word_coef.items() if coef < 0}
17
18 # Generate word clouds for hateful and non-hateful words
19 hateful_cloud = WordCloud(width=140, height=140, background_color='white').generate_from_frequencies(hateful_words)
20 non_hateful_cloud = WordCloud(width=140, height=140, background_color='white').generate_from_frequencies(non_hateful_words)

```

This code snippet analyzes the coefficients learned by a LinearSVC model to identify words associated with hateful and non-hateful classes. It retrieves feature names from the TF-IDF vectorization process, extracts coefficients from the model, and separates words based on their coefficients into dictionaries for hateful and non-hateful classes. Finally, word clouds are generated to visually represent the most influential words for each class. Feature Names consist of : 'abandon', 'abdul', ... 'zoo', 'zuckerberg'.

```

1 # Plot the word clouds
2 plt.figure(figsize=(15, 6))
3 plt.subplot(1, 2, 1)
4 plt.imshow(hateful_cloud, interpolation='bilinear')
5 plt.title('Hateful Words')
6 plt.axis('off')
7
8 plt.subplot(1, 2, 2)
9 plt.imshow(non_hateful_cloud, interpolation='bilinear')
10 plt.title('Non-Hateful Words')
11 plt.axis('off')
12
13 plt.show()

```

The wordclouds containing influential words pertaining to hateful and non-hateful labels are showed:



Figure 12:

```

1 def classify_text(input_text):
2     # Preprocess the input text
3     preprocessed_text = preprocess_text(input_text)
4
5     # Vectorize the preprocessed text
6     text_vector = vectorizer.transform([preprocessed_text])
7
8     # Predict the label using the trained classifier
9     predicted_label = classifier.predict(text_vector)
10
11     # Map the predicted label to its corresponding category
12     if predicted_label[0] == 1:
13         return "Hateful"
14     else:
15         return "Non-Hateful"
16
17 # Example usage:
18 input_text = "retarded"
19 classification = classify_text(input_text)
20 print("Classification:", classification)

```

For a single input text, the classify function above predictions whether or not the input text is Hateful or Not. For example,

```

1 Input text:'i am not racist some of my best slaves are black'
2 Output:Hateful
3 Input text:'we talked about it i agreed dinner would be ready on time in the future'
4 Output:Non-Hateful

```

```

1  # Load JSON lines file for testing
2  with open("test_seen.jsonl", "r") as f:
3      test_json_lines = f.readlines()
4
5  # Parse JSON lines for testing
6  test_data = []
7  for line in test_json_lines:
8      item = json.loads(line)
9      test_data.append((item['text'], item['id']))
10
11 # Predict labels for test data and save results to an external file
12 output_file = "text_predictions.txt" # Define the output file name
13
14 with open(output_file, "w") as f:
15     for text, image_id in test_data:
16         prediction = classify_text(text)
17         f.write(f"Text: {text} | Image ID: {image_id} | Prediction: {prediction}\n")
18
19 print(f"Predictions saved to {output_file}")

```

This code snippet loads text data from a JSON lines file named "test_seen.jsonl" for testing. It then parses the JSON lines to extract text and corresponding IDs. Next, it predicts labels for the test data using a function called `classify_text`, which applies the classifier to classify the text as hateful or non-hateful. Predictions, along with text and image IDs, are written to an external file named "text_predictions.txt".

```

1  Text: handjobs sold seperately | Image ID: 16395 | Prediction: Non-Hateful
2  Text: introducing fidget spinner for women | Image ID: 37405 | Prediction: Hateful
3  Text: happy pride month let's go beat up lesbians | Image ID: 94180 | Prediction: Non-Hateful

```

Glimpse of text_predictions.txt file.

```

1  def calculate_toxicity_score(input_text):
2      # Preprocess the input text
3      preprocessed_text = preprocess_text(input_text)
4
5      # Tokenize the preprocessed text
6      tokens = word_tokenize(preprocessed_text)
7
8      # Count the number of hateful words
9      hateful_word_count = sum(classifier.predict(vectorizer.transform([word])) == 1 for word in tokens)
10
11     # Calculate toxicity score
12     total_words = len(tokens)
13     toxicity_score = (hateful_word_count / total_words) * 100
14
15     return toxicity_score
16
17 # Example usage:
18 input_text = "retarded"

```

```
19 toxicity_score = calculate_toxicity_score(input_text)
20 print("Toxicity Score:", toxicity_score)
```

The above code is a simple toxicity score calculator for a given input text, which gives a score between 0-100. How the function works is by counting the number of hateful words and then dividing that by the total number of words in the sentence which is then multiplied by 100 to obtain percentage. Example,

```
1 Input text:'we talked about it i agreed dinner would be ready on time in the future'
2 Output:14.285
3 Input txt:'imagine being so disgusting there have to be laws to try to stop normal people from hating you'
4 Output:37.5
```

Drawback of method: One drawback is that there is no real contextual learning/understanding of texts labelled with hateful/non-hateful happening. It's all word by word analysis of the text and learning. This method may not be accurate in some cases: For example: Input txt:'muslim' Output:100 percent toxic/hateful. Here there's no real context to the text being hateful yet the classifier predicts hateful. This happens because in the training of the classifier it was fed with the notion that the word 'muslim' appeared a lot in hateful labelled texts and hence associated it to be hateful word which influences the label of a sentence.

We've explored both meme-based and text-based models, each exhibiting a bias toward one modality—either image or text. Both models generate predictions, which are then stored in separate output files for further analysis. We shall analyse each of their individual predictions and combined predictions and make a new model which takes the best out of the two models.

4.3.3 Comparison and analysis of the Meme-based model and Text-based model:

```
1 import re
2 import json
3 import matplotlib.pyplot as plt
```

The required libraries needed for this task.

```
1 # Function to parse predictions from meme_predictions.txt
2 def parse_meme_predictions(file_path):
3     predictions = []
4     with open(file_path, 'r') as file:
```

```

5         for line in file:
6             match = re.search(r'Image: (.+), Predicted Label: (.+)', line)
7             if match:
8                 image_id = match.group(1).split('/')[-1].split('.')[0] # Extract image ID from the file path
9                 predicted_label = match.group(2)
10                predictions.append((image_id, predicted_label))
11            return predictions
12
13    # Function to parse predictions from text_predictions.txt
14    def parse_text_predictions(file_path):
15        predictions = []
16        with open(file_path, 'r') as file:
17            for line in file:
18                match = re.search(r'Text: (.+) \| Image ID: (\d+) \| Prediction: (.+)', line)
19                if match:
20                    text = match.group(1)
21                    image_id = match.group(2)
22                    predicted_label = match.group(3)
23                    predictions.append((image_id, predicted_label))
24            return predictions
25
26    # Function to parse true labels from test_seen.jsonl
27    def parse_true_labels(file_path):
28        true_labels = {}
29        with open(file_path, 'r') as file:
30            for line in file:
31                data = json.loads(line)
32                image_id = data['id']
33                label = data['label']
34                true_labels[image_id] = label
35            return true_labels
36
37    # Function to save image IDs into a text file
38    def save_image_ids(image_ids, file_path):
39        with open(file_path, 'w') as file:
40            for image_id in image_ids:
41                file.write(image_id + '\n')
42

```

These functions serve to parse predictions and true labels from respective text files and JSON files, and save image IDs into a text file. `parse_meme_predictions()` and `parse_text_predictions()` extract predictions from text files containing predictions for memes and text, respectively, by employing regular expressions to capture relevant information such as image IDs and predicted labels. `parse_true_labels()` retrieves true labels from a JSON file, storing them in a dictionary with image IDs as keys. Finally, `save_image_ids()` writes a list of image IDs to a specified text file. These functions facilitate the extraction, parsing, and storage of predictions and true labels, which can then be used for further analysis or evaluation of the models' performance.

```

1
2    # Parse predictions from meme_predictions.txt and text_predictions.txt
3    meme_predictions = parse_meme_predictions('meme_predictions.txt')
4    text_predictions = parse_text_predictions('text_predictions.txt')

```

```

5
6 # Parse true labels from test_seen.jsonl
7 true_labels = parse_true_labels('test_seen.jsonl')

```

The predictions from each of the text files are stored in the respective predictions lists and the true labels is also extracted into a dictionary which can be used for further analysis.

```

1 import matplotlib.pyplot as plt
2
3 # Calculate the total number of predictions
4 total_predictions = len(true_labels)
5
6 # Count the number of hateful and non-hateful cases
7 hateful_count = sum(1 for label in true_labels.values() if label == 1)
8 non_hateful_count = sum(1 for label in true_labels.values() if label == 0)
9
10 # Calculate the percentage of hateful and non-hateful cases
11 hateful_percentage = (hateful_count / total_predictions) * 100
12 non_hateful_percentage = (non_hateful_count / total_predictions) * 100
13
14 # Plotting the pie chart
15 labels = ['Hateful', 'Non-Hateful']
16 sizes = [hateful_percentage, non_hateful_percentage]
17 colors = ['red', 'blue']
18 explode = (0.1, 0) # explode the 1st slice (Hateful)
19
20 plt.figure(figsize=(8, 6))
21 plt.pie(sizes, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%', shadow=True, startangle=140)
22 plt.title('Percentage of Actual Hateful and Non-Hateful Cases')
23 plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle
24 plt.show()

```

This plots a pie-chart which shows the percentage of actual Hateful and Non-hateful cases in the test dataset sourced from the true labels. Note that there

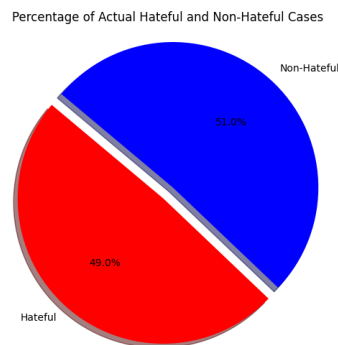


Figure 13:

are 1000 memes to be tested in the test dataset. From the pie-chart it's evident that 490 memes are actually Hateful and the rest 510 memes are Non-Hateful.

```

1  def plot_percentage_hateful_non_hateful(meme_predictions, text_predictions, true_labels):
2      # Count the total number of memes
3      total_meme_count = len(meme_predictions)
4      total_text_count = len(text_predictions)
5
6      # Count the number of hateful and non-hateful memes predicted by each model
7      meme_hateful_count = sum(1 for _, label in meme_predictions if label == 'Hateful')
8      meme_non_hateful_count = total_meme_count - meme_hateful_count
9
10     text_hateful_count = sum(1 for _, label in text_predictions if label == 'Hateful')
11     text_non_hateful_count = total_text_count - text_hateful_count
12
13     # Calculate percentages
14     meme_hateful_percentage = meme_hateful_count / total_meme_count * 100
15     meme_non_hateful_percentage = meme_non_hateful_count / total_meme_count * 100
16
17     text_hateful_percentage = text_hateful_count / total_text_count * 100
18     text_non_hateful_percentage = text_non_hateful_count / total_text_count * 100
19
20     # Plotting
21     labels = ['Meme Model', 'Text Model']
22     hateful_percentages = [meme_hateful_percentage, text_hateful_percentage]
23     non_hateful_percentages = [meme_non_hateful_percentage, text_non_hateful_percentage]
24
25     x = range(len(labels))
26
27     plt.figure(figsize=(8, 6))
28     plt.bar(x, hateful_percentages, color='red', width=0.4, label='Hateful')
29     plt.bar(x, non_hateful_percentages, color='blue', width=0.4, label='Non-Hateful', bottom=hateful_percentages)
30
31     plt.xlabel('Model')
32     plt.ylabel('Percentage')
33     plt.title('Percentage of Hateful and Non-Hateful Memes Predicted by Each Model')
34     plt.xticks(x, labels)
35     plt.legend()
36     plt.ylim(0, 100)
37
38     # Add annotations to the bars
39     for i, (hateful_percentage, non_hateful_percentage) in enumerate(zip(hateful_percentages, non_hateful_percentages)):
40         plt.text(i, hateful_percentage / 2, f'{hateful_percentage:.2f}%', ha='center', va='bottom')
41         plt.text(i, hateful_percentage + non_hateful_percentage / 2, f'{non_hateful_percentage:.2f}%', ha='center', va='bottom')
42
43     plt.show()
44
45 plot_percentage_hateful_non_hateful(meme_predictions, text_predictions, true_labels)

```

The above lines of code plots the percentage of Hateful and Non-Hateful memes predicted by each model-Meme based and the text based model. Analysis of this is done with the help of `meme_predictions.txt` and `text_predictions.txt`.

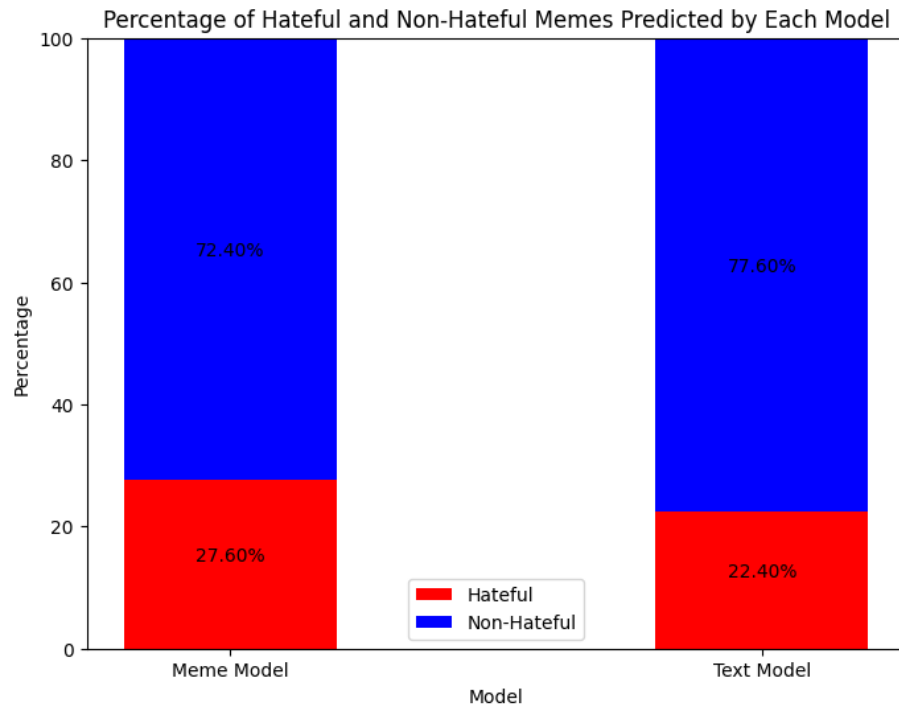


Figure 14:

```

1  from sklearn.metrics import confusion_matrix
2  import seaborn as sns
3
4  # Function to plot confusion matrix
5  def plot_confusion_matrix(true_labels, predicted_labels, model_name):
6      # Calculate confusion matrix
7      cm = confusion_matrix(true_labels, predicted_labels)
8
9      # Plot confusion matrix using seaborn heatmap
10     plt.figure(figsize=(8, 6))
11     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False)
12     plt.xlabel('Predicted Label')
13     plt.ylabel('True Label')
14     plt.title(f'Confusion Matrix for {model_name}')
15     plt.show()
16
17     # Calculate true labels and predicted labels for each model
18     true_labels_text = [true_labels.get(image_id) for image_id, label in text_predictions]
19     predicted_labels_text = [1 if label == 'Hateful' else 0 for image_id, label in text_predictions]
20
21     true_labels_meme = [true_labels.get(image_id) for image_id, label in meme_predictions]
22     predicted_labels_meme = [1 if label == 'Hateful' else 0 for image_id, label in meme_predictions]
23
24     # Plot confusion matrix for the text model
25     plot_confusion_matrix(true_labels_text, predicted_labels_text, 'Text Model')
26
27     # Plot confusion matrix for the image model
28     plot_confusion_matrix(true_labels_meme, predicted_labels_meme, 'Meme Model')
29

```

The above lines of code plots the *Confusion matrix* for both the text based model and the meme based model.

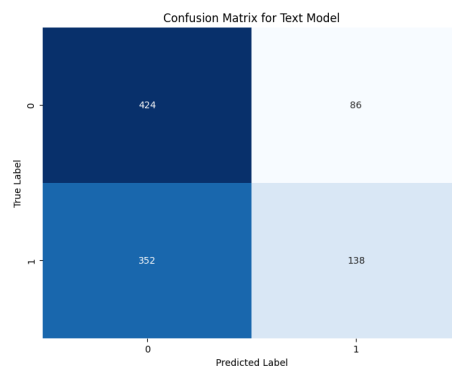


Figure 15:

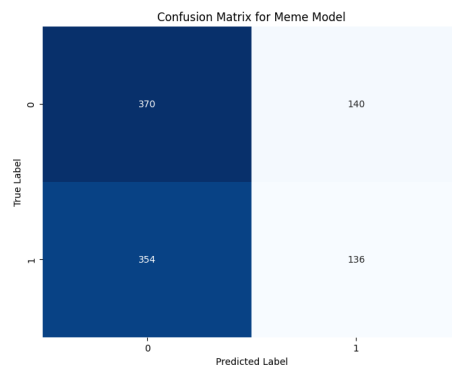


Figure 16:

```

1  # Function to plot percentage of correct predictions for each label type (Hateful and Non-Hateful)
2  def plot_correct_percentage(predictions, true_labels, model_name):
3      hateful_correct = 0
4      non_hateful_correct = 0
5      total_hateful = 0
6      total_non_hateful = 0
7
8      for pred in predictions:
9          image_id, label = pred
10         true_label = true_labels.get(image_id)
11         if true_label is not None:
12             if true_label == 1: # Hateful
13                 total_hateful += 1
14                 if label == 'Hateful':

```



```

15         hateful_correct += 1
16     elif true_label == 0: # Non-Hateful
17         total_non_hateful += 1
18         if label == 'Non-Hateful':
19             non_hateful_correct += 1
20
21     # Calculate percentages
22     hateful_percentage = hateful_correct / total_hateful * 100 if total_hateful != 0 else 0
23     non_hateful_percentage = non_hateful_correct / total_non_hateful * 100 if total_non_hateful != 0 else 0
24
25     # Plotting
26     labels = ['Hateful', 'Non-Hateful']
27     percentages = [hateful_percentage, non_hateful_percentage]
28
29     plt.figure(figsize=(6, 4))
30     bars = plt.bar(labels, percentages, color=['red', 'blue'])
31     plt.title(f'Percentage of Correct Predictions by {model_name}')
32     plt.xlabel('Meme Type')
33     plt.ylabel('Percentage')
34     plt.ylim(0, 100)
35
36     # Add annotations to the bars
37     for bar, percentage in zip(bars, percentages):
38         height = bar.get_height()
39         plt.text(bar.get_x() + bar.get_width() / 2, height, f'{percentage:.2f}%', ha='center', va='bottom')
40
41     plt.show()
42
43     # Plot percentage of correct predictions for each label type (Hateful and Non-Hateful) for the text model
44     plot_correct_percentage(text_predictions, true_labels, 'Text Model')
45
46     # Plot percentage of correct predictions for each label type (Hateful and Non-Hateful) for the meme model
47     plot_correct_percentage(meme_predictions, true_labels, 'Meme Model')

```

The code above plots the percentage of correct predictions for each of the Hateful/Non-Hateful predictions made by both the models.

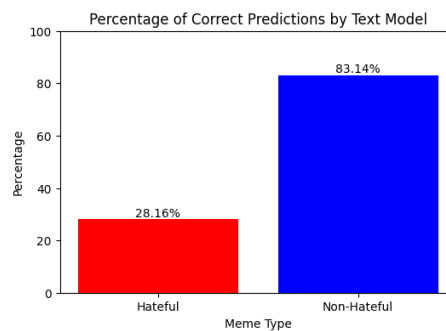


Figure 17:

1 # Function to calculate values for each quadrant in the confusion matrix

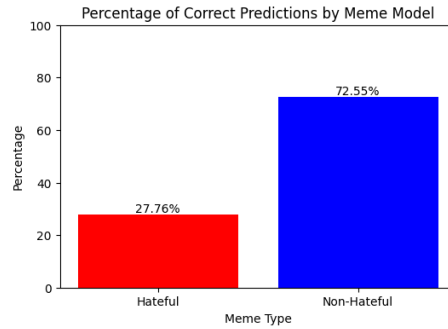


Figure 18:

```

2  def calculate_confusion_matrix_values(true_labels, predicted_labels):
3      true_hateful = true_non_hateful = false_hateful = false_non_hateful = 0
4
5      for true_label, predicted_label in zip(true_labels, predicted_labels):
6          if true_label == 1:
7              if predicted_label == 1:
8                  true_hateful += 1
9              else:
10                 false_non_hateful += 1
11          else: # True label is Non-Hateful
12              if predicted_label == 0:
13                  true_non_hateful += 1
14              else:
15                  false_hateful += 1
16
17      return true_hateful, true_non_hateful, false_hateful, false_non_hateful
18
19  # Calculate values for each quadrant in the confusion matrix for the text model
20  text_true_hateful, text_true_non_hateful, text_false_hateful, text_false_non_hateful = calculate_confusion_matrix(
21
22  # Calculate values for each quadrant in the confusion matrix for the image model
23  meme_true_hateful, meme_true_non_hateful, meme_false_hateful, meme_false_non_hateful = calculate_confusion_matrix(
24
25  # # Print the values for each quadrant for both models
26  # print("Text Model Confusion Matrix:")
27  # print("True Hateful:", text_true_hateful)
28  # print("True Non-Hateful:", text_true_non_hateful)
29  # print("False Hateful:", text_false_hateful)
30  # print("False Non-Hateful:", text_false_non_hateful)
31
32  # print("\nImage Model Confusion Matrix:")
33  # print("True Hateful:", image_true_hateful)
34  # print("True Non-Hateful:", image_true_non_hateful)
35  # print("False Hateful:", image_false_hateful)
36  # print("False Non-Hateful:", image_false_non_hateful)

```

The values for each of the quadrants in the confusion matrices are stored in relevant variables for calculation of each model's performance metrics such as

Accuracy, Precision, Recall and F1 score.

```
1
2
3 # Calculate accuracy for the text model
4 text_total = len(text_predictions)
5 text_correct = sum(1 for text_pred in text_predictions if true_labels.get(text_pred[0]) == int(text_pred[1] == 'H
6 text_accuracy = text_correct / text_total * 100
7
8 # Calculate accuracy for the image model
9 meme_total = len(meme_predictions)
10 meme_correct = sum(1 for meme_pred in meme_predictions if true_labels.get(meme_pred[0]) == int(meme_pred[1] == 'H
11 meme_accuracy = meme_correct / meme_total * 100
12
13 # Calculate precision, recall, and F1 score for the text model
14 text_precision = text_true_hateful/(text_true_hateful + text_false_hateful)
15 text_recall = text_true_hateful/(text_true_hateful + text_false_non_hateful)
16 text_f1 = 2 * (text_precision * text_recall) / (text_precision + text_recall)
17
18 # Calculate precision, recall, and F1 score for the image model
19 meme_precision = meme_true_hateful / (meme_true_hateful + meme_false_hateful)
20 meme_recall = meme_true_hateful / (meme_true_hateful + meme_false_non_hateful)
21 meme_f1 = 2 * (meme_precision * meme_recall) / (meme_precision + meme_recall)
22
```

The performance metrics for each model is calculated and printed.

```
1 Text Model Accuracy: 56.20%
2 Text Model Precision: 0.62
3 Text Model Recall: 0.28
4 Text Model F1 Score: 0.39
```

```
1 Meme Model Accuracy: 50.60%
2 Meme Model Precision: 0.49
3 Meme Model Recall: 0.28
4 Meme Model F1 Score: 0.36
```

Reflecting on the performance metrics of the models, the text based model performed better on all fronts compared to the meme based model. Notably the less than average F1 scores for each of the models can be attributed to the inefficiency of the models in predicting labels. Various parameters have to be hyper-tuned to achieve a good F1 score. F1 score is the metric which gives a clear perspective about a model's performance. Accuracy is not really the best parameter to judge a model's performance since the testing dataset can be biased towards a single label. Each of the precision and recall parameters individually are not good enough for judging performance of a model. F1 combines both Precision and Recall to achieve a good metric for judging a model's performance.

```
1 # Define scenario labels
2 scenarios = ['Both Hateful', 'Both Non-Hateful', 'Meme Hateful, Text Non-Hateful', 'Meme Non-Hateful, Text Hateful']
```

```

3 scenario_counts = [0] * len(scenarios)
4 # Compare predictions and categorize cases
5 for meme_pred, text_pred in zip(meme_predictions, text_predictions):
6     image_id_meme, label_meme = meme_pred
7     image_id_text, label_text = text_pred
8     if label_meme == 'Hateful' and label_text == 'Hateful':
9         scenario_counts[0] += 1 # Both Hateful
10    elif label_meme == 'Non-Hateful' and label_text == 'Non-Hateful':
11        scenario_counts[1] += 1 # Both Non-Hateful
12    elif label_meme == 'Hateful' and label_text == 'Non-Hateful':
13        scenario_counts[2] += 1 # Meme Hateful, Text Non-Hateful
14    elif label_meme == 'Non-Hateful' and label_text == 'Hateful':
15        scenario_counts[3] += 1 # Meme Non-Hateful, Text Hateful
16
17 # Calculate the total number of predictions
18 total_predictions = len(meme_predictions)
19
20 # Calculate the percentage for each scenario
21 scenario_percentages = [count / total_predictions * 100 for count in scenario_counts]

```

Now the analysis has shifted towards comparing both of the output files of the models in a combined manner. In the provided code segment, scenario labels are defined, including categories such as "Both Hateful," "Both Non-Hateful," "Meme Hateful, Text Non-Hateful," and "Meme Non-Hateful, Text Hateful." A list is initialized to count occurrences of each scenario. Then, the predictions from both meme-based and text-based models are compared and categorized into these scenarios based on the predicted labels. Counts are incremented accordingly for each scenario. Finally, the total number of predictions is calculated, and the percentage of occurrences for each scenario is computed relative to the total number of predictions. This analysis offers insights into the agreement and disagreement between the predictions of the two models. For example, for one meme, one model predicts hateful and the other predicts non-hateful, and cases like these are recorded in the code segment. Cases of agreement are also recorded in the code segment.

```

1 # Plot the comparison results with percentages and annotations
2 plt.figure(figsize=(10, 6))
3 bars = plt.bar(scenarios, scenario_percentages, color='skyblue')
4 plt.xlabel('Scenario')
5 plt.ylabel('Percentage of Predictions')
6 plt.title('Comparison of Predictions between Text-based and Meme-based Models')
7 plt.xticks(rotation=45, ha='right')
8 plt.tight_layout()
9 # Add annotations to the bars
10 for bar, percentage in zip(bars, scenario_percentages):
11     height = bar.get_height()
12     plt.text(bar.get_x() + bar.get_width() / 2, height, f'{percentage:.2f}%', ha='center', va='bottom')
13
14 plt.show()

```

Cases of Agreement and Disagreement between the models are plotted. From

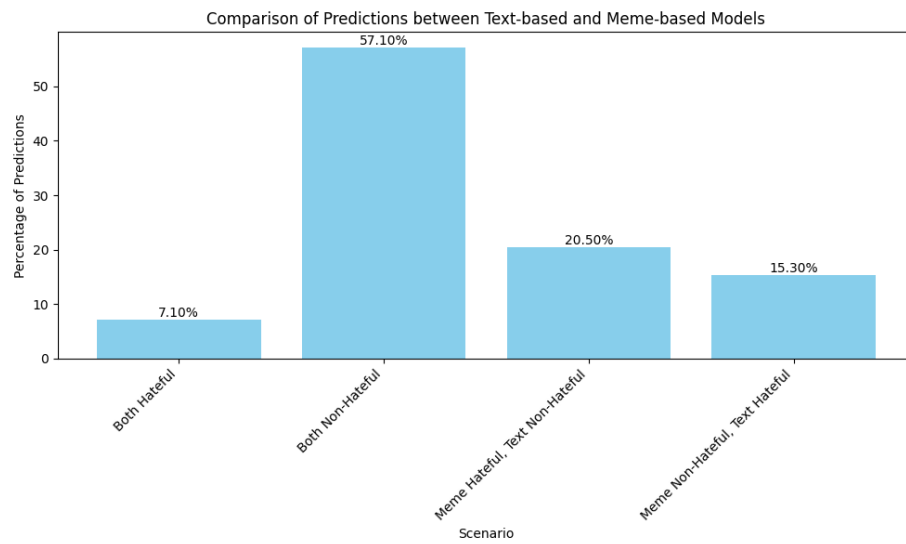


Figure 19:

the plot its clear that to a good extent both of the models agree with each other's predictions on the test dataset.

```

1  meme_hateful_text_hateful_both_correct =0
2
3  for meme_pred, text_pred in zip(meme_predictions, text_predictions):
4      image_id_meme, label_meme = meme_pred
5      image_id_text, label_text = text_pred
6
7      true_label = true_labels.get(image_id_meme)
8      if true_label is not None:
9          if label_meme == 'Hateful' and label_text == 'Hateful' and true_label == 1:
10             meme_hateful_text_hateful_both_correct += 1
11
12
13  total_meme_hateful_text_hateful=scenario_counts[0]
14  meme_hateful_text_hateful_both_correct_percentage = meme_hateful_text_hateful_both_correct / total_meme_hateful_t
15  print(f"Percentage of cases where both meme and text predict hateful and their prediction is right: {meme_hateful_t
16
17  meme_non_hateful_text_non_hateful_both_correct =0
18
19  for meme_pred, text_pred in zip(meme_predictions, text_predictions):
20      image_id_meme, label_meme = meme_pred
21      image_id_text, label_text = text_pred
22
23      true_label = true_labels.get(image_id_meme)
24      if true_label is not None:
25          if label_meme == 'Non-Hateful' and label_text == 'Non-Hateful' and true_label == 0:
26             meme_non_hateful_text_non_hateful_both_correct += 1
27
28
29  total_meme_non_hateful_text_non_hateful=scenario_counts[1]
30  meme_non_hateful_text_non_hateful_both_correct_percentage = meme_non_hateful_text_non_hateful_both_correct / tota
31  print(f"Percentage of cases where both meme and text predict non-hateful and their prediction is right: {meme_non
32
33  both_same_correct=0
34  for meme_pred, text_pred in zip(meme_predictions, text_predictions):
35      image_id_meme, label_meme = meme_pred
36      image_id_text, label_text = text_pred
37
38      true_label = true_labels.get(image_id_meme)
39      if true_label is not None:
40          if (label_meme == 'Hateful' and label_text == 'Hateful' and true_label == 1) or (label_meme == 'Non-Hatef
41             both_same_correct+=1
42
43  total_both_same=scenario_counts[0]+scenario_counts[1]
44  both_same_correct_percentage = both_same_correct/total_both_same * 100
45  print(f"Percentage of cases where both meme and text predict same and their prediction is right: {both_same_corre
46
47  # Initialize counter for cases where meme predicts hateful and text predicts non-hateful and meme is correct
48  meme_hateful_text_non_hateful_meme_correct = 0
49
50  # Iterate through both sets of predictions and compare with true labels
51  for meme_pred, text_pred in zip(meme_predictions, text_predictions):
52      image_id_meme, label_meme = meme_pred
53      image_id_text, label_text = text_pred
54
55      true_label = true_labels.get(image_id_meme)
56      if true_label is not None:

```

```

57         # Check if text predicts hateful and meme predicts non-hateful
58         if label_meme == 'Hateful' and label_text == 'Non-Hateful' and true_label == 1:
59             meme_hateful_text_non_hateful_meme_correct += 1
60
61     # Calculate the percentage of correct predictions by the text model
62     total_meme_hateful_text_non_hateful = scenario_counts[2]
63     meme_hateful_text_non_hateful_meme_correct_percentage = (meme_hateful_text_non_hateful_meme_correct / total_meme_
64
65     print(f"Percentage of cases where meme predicts hateful and text predicts non-hateful, and meme prediction is rig
66     # Calculate the percentage of cases where the prediction made by the meme model is correct
67     meme_hateful_text_non_hateful_text_correct_percentage = 100 - meme_hateful_text_non_hateful_meme_correct_percenta
68
69     print(f"Percentage of cases where meme predicts hateful and text predicts non-hateful, and text prediction is rig
70
71     # Initialize counter for cases where meme predicts non-hateful and text predicts hateful and text is correct
72     meme_non_hateful_text_hateful_text_correct = 0
73
74     # Iterate through both sets of predictions and compare with true labels
75     for meme_pred, text_pred in zip(meme_predictions, text_predictions):
76         image_id_meme, label_meme = meme_pred
77         image_id_text, label_text = text_pred
78
79         true_label = true_labels.get(image_id_meme)
80         if true_label is not None:
81             # Check if meme predicts hateful and text predicts non-hateful
82             if label_meme == 'Non-Hateful' and label_text == 'Hateful' and true_label == 1:
83                 meme_non_hateful_text_hateful_text_correct += 1
84
85     # Calculate the percentage of correct predictions by the meme model
86     total_meme_non_hateful_text_hateful = scenario_counts[3]
87     meme_non_hateful_text_hateful_text_correct_percentage = (meme_non_hateful_text_hateful_text_correct/total_meme_no
88
89     print(f"Percentage of cases where meme predicts non-hateful and text predicts hateful, and text prediction is rig
90     # Calculate the complementary percentage for cases where meme predicts hateful and text predicts non-hateful
91     meme_non_hateful_text_hateful_meme_correct_percentage = 100 - meme_non_hateful_text_hateful_text_correct_percenta
92
93     print(f"Percentage of cases where meme predicts non-hateful and text predicts hateful, and meme prediction is rig
94

```

The provided code snippet facilitates a comprehensive comparison between the predictions of the meme-based and text-based models against the true labels of the test memes. It computes various percentages to analyze agreement and disagreement between the models' predictions and the true labels. Specifically, it calculates the percentage of cases where both models predict hateful and are correct, both predict non-hateful and are correct, both predict the same and are correct, meme predicts hateful while text predicts non-hateful and meme's prediction is correct, and vice versa. This analysis offers a good understanding of the performance and alignment between the models' predictions and the ground truth labels.

```

1  import matplotlib.pyplot as plt
2
3  # Define categories and their respective percentages
4  categories = [

```

```

5     "Both Meme and Text Predict Hateful (Both Correct)",
6     "Both Meme and Text Predict Non-Hateful (Both Correct)",
7     "Both Meme and Text Predict Same (Both Correct)",
8     "Meme Predicts Hateful, Text Predicts Non-Hateful (Meme Correct)",
9     "Meme Predicts Hateful, Text Predicts Non-Hateful (Text Correct)",
10    "Meme Predicts Non-Hateful, Text Predicts Hateful (Text Correct)",
11    "Meme Predicts Non-Hateful, Text Predicts Hateful (Meme Correct)"
12 ]
13 percentages = [
14     meme_hateful_text_hateful_both_correct_percentage,
15     meme_non_hateful_text_non_hateful_both_correct_percentage,
16     both_same_correct_percentage,
17     meme_hateful_text_non_hateful_meme_correct_percentage,
18     meme_hateful_text_non_hateful_text_correct_percentage,
19     meme_non_hateful_text_hateful_text_correct_percentage,
20     meme_non_hateful_text_hateful_meme_correct_percentage
21 ]
22
23 # Plotting the bar chart
24 plt.figure(figsize=(10, 6))
25 plt.barh(categories, percentages, color='skyblue')
26 plt.xlabel('Percentage')
27 plt.ylabel('Categories')
28 plt.title('Comparison of Correct Predictions between Meme and Text Models')
29 plt.xlim(0, 100)
30
31 # Adding percentage labels to the bars
32 for index, value in enumerate(percentages):
33     plt.text(value, index, f'{value:.2f}%', va='center')
34
35 plt.show()

```

A plot is generated which predicts the percentage of cases of Agreement and Disagreement of predictions of the two models compared with True ground labels. Based on insights derived from this plot, a model which takes the best

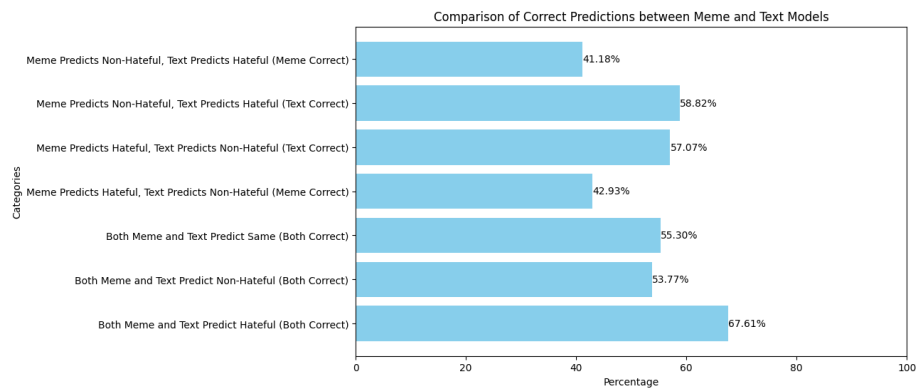


Figure 20:

of the two model's predictions can be built.

Food for thought: What do you think the combined model can be in this case? How much aspects of each of the models should the combined model imbibe in itself? A careful look at the plot above can lead to a very good guess of the details of the combined model.

Insights: From the above plot it's clear that whenever both the models predict the same thing, its more likely to be correct as evident by 67.21 percent of correct cases where both predict Hateful and 53.77 percent of correct cases where both predict Non-Hateful leading to 55.30 percent of correct cases where both agree in general. Now when the models disagree with each other, like in the cases where Meme model predicts Non-Hateful and text model predicts Hateful, the text model emerges victorious in most of the cases-58.62 percent. Even in cases where Meme model predicts Hateful and text model predicts Non-Hateful, the text model emerges victorious in most of the cases- 57.07 percent.

So a good combined model can have the following predictions dictated by the following rules:

- Whenever both of the models predict the same thing, put the final prediction of the combined model to be the same prediction made by both the models.
- Whenever both the models disagree with each other, put the final prediction of the combined model to be the prediction made by the text based model.

```
1  # Initialize a list to store the final combined predictions
2  combined_predictions = []
3
4  # Iterate through both sets of predictions and apply the combining rules
5  for meme_pred, text_pred in zip(meme_predictions, text_predictions):
6      image_id_meme, label_meme = meme_pred
7      image_id_text, label_text = text_pred
8
9      # Rule 1: Both models predict 'Hateful'
10     if label_meme == 'Hateful' and label_text == 'Hateful':
11         combined_predictions.append((image_id_meme, 'Hateful'))
12     # Rule 2: Both models predict 'Non-Hateful'
13     elif label_meme == 'Non-Hateful' and label_text == 'Non-Hateful':
14         combined_predictions.append((image_id_meme, 'Non-Hateful'))
15     # Rule 3: Meme model predicts 'Hateful' and text model predicts 'Non-Hateful'
16     elif label_meme == 'Hateful' and label_text == 'Non-Hateful':
17         combined_predictions.append((image_id_meme, 'Non-Hateful'))
18     # Rule 4: Meme model predicts 'Non-Hateful' and text model predicts 'Hateful'
19     elif label_meme == 'Non-Hateful' and label_text == 'Hateful':
20         combined_predictions.append((image_id_meme, 'Hateful'))
```

Code snippet which makes the final predictions for the combined model based on the rules described above.

p.s: Incase if it's not clear, the combined model essentially mirrors the behavior of the text-based model in its entirety. It incorporates the final pre-

dictions generated by the text-based model into its own predictions, relying heavily on the insights and decisions made by the text-based model throughout its prediction process. This also makes sense as previously we noted that the text-based model performed better on all performance parameters than the meme-based model and undoubtedly using the text-based model in its entirety for the final predictions in the combined model make sense.

Note: From the performance of the two models and the agreement vs disagreement analysis of the two models coupled with the true labels, the combined model effectively simulates the text based model in its entirety. For this particular test dataset this has happened to be the case. In some cases the combined model can combine aspects of the meme based model with the text based model depending on the dataset.

Food for Thought: Was analysing the agreement/disagreement cases of the two models coupled with the true labels actually necessary to choose the combined model to simulate the text based model in its entirety? Or the fact that the text based models performed better on all performance parameters is sufficient to make the decision of choosing the combined model to simulate the text based model in its entirety?

Some plots of the combined model: Essentially in this case, the plots of the combined model will be identical to that of the text model.

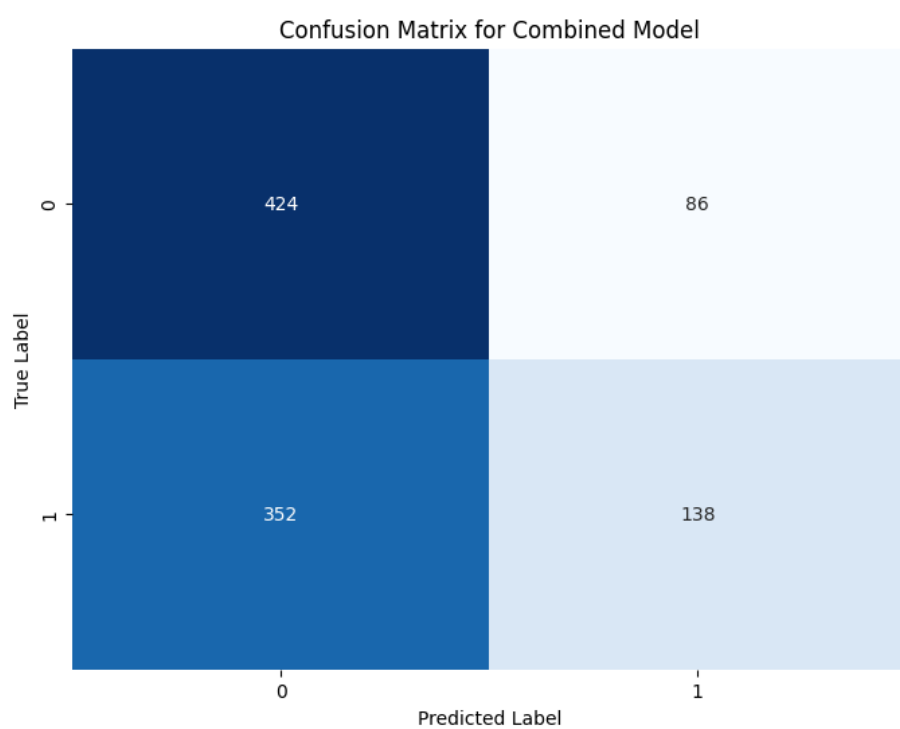


Figure 21:

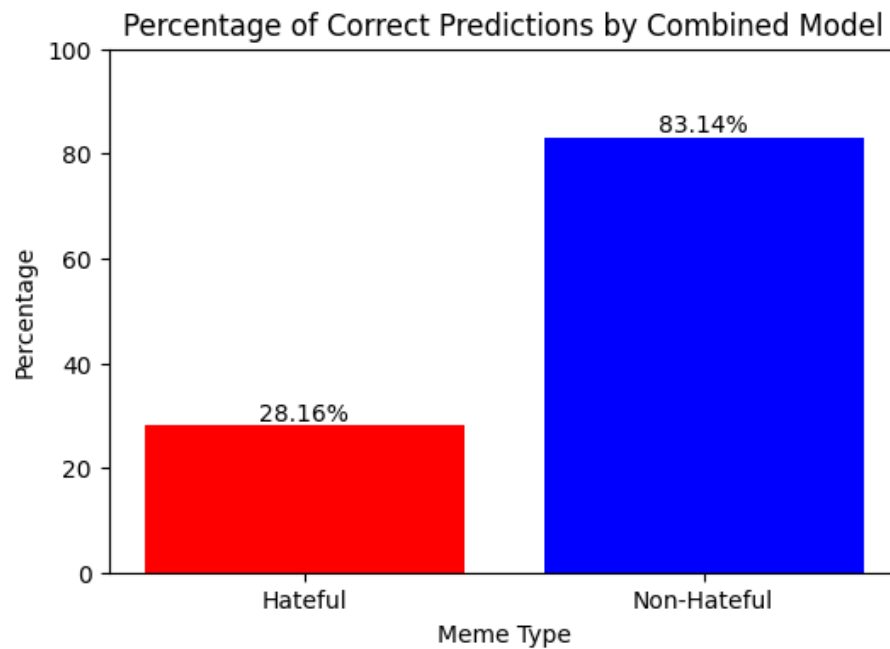


Figure 22:

The performance metrics of the Combined model would be identical to that of the text model in this case:

```
1 Combined Model Accuracy: 56.20
2 Combined Model Metrics:
3 Precision: 0.62
4 Recall: 0.28
5 F1 Score: 0.39
```

Finally a comparison between the meme-based model, text-based model and combined model.

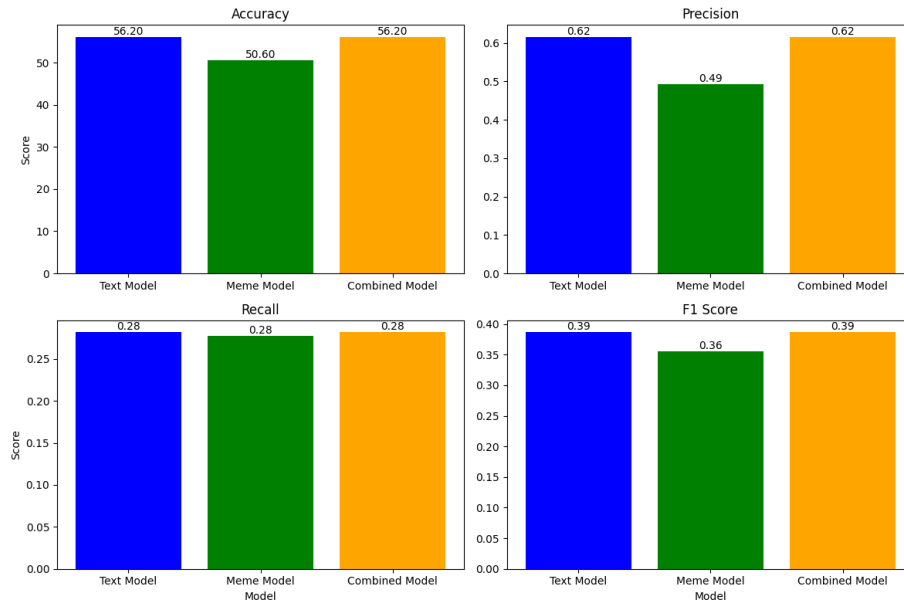


Figure 23:

```

1
2  # Save IDs of hateful and non-hateful memes predicted by each model into separate text files
3  meme_hateful_ids = [image_id for image_id, label in meme_predictions if label == 'Hateful']
4  meme_non_hateful_ids = [image_id for image_id, label in meme_predictions if label == 'Non-Hateful']
5  text_hateful_ids = [image_id for image_id, label in text_predictions if label == 'Hateful']
6  text_non_hateful_ids = [image_id for image_id, label in text_predictions if label == 'Non-Hateful']
7  combined_hateful_ids = [image_id for image_id, label in combined_predictions if label == 'Hateful']
8  combined_non_hateful_ids = [image_id for image_id, label in combined_predictions if label == 'Non-Hateful']
9
10 save_image_ids(meme_hateful_ids, 'meme_hateful_ids.txt')
11 save_image_ids(meme_non_hateful_ids, 'meme_non_hateful_ids.txt')
12 save_image_ids(text_hateful_ids, 'text_hateful_ids.txt')
13 save_image_ids(text_non_hateful_ids, 'text_non_hateful_ids.txt')
14 save_image_ids(combined_hateful_ids, 'combined_hateful_ids.txt')
15 save_image_ids(combined_non_hateful_ids, 'combined_non_hateful_ids.txt')
16
17
18 # Initialize lists to store image IDs for each scenario
19 both_hateful = []
20 both_non_hateful = []
21 text_hateful_meme_non_hateful = []
22 text_non_hateful_meme_hateful = []
23
24
25 # Compare predictions and categorize cases, storing image IDs accordingly
26 for meme_pred, text_pred in zip(meme_predictions, text_predictions):
27     image_id_meme, label_meme = meme_pred

```

```

28     image_id_text, label_text = text_pred
29     if label_meme == 'Hateful' and label_text == 'Hateful':
30         both_hateful.append(image_id_meme)
31     elif label_meme == 'Non-Hateful' and label_text == 'Non-Hateful':
32         both_non_hateful.append(image_id_meme)
33     elif label_meme == 'Hateful' and label_text == 'Non-Hateful':
34         text_hateful_meme_non_hateful.append(image_id_meme)
35     elif label_meme == 'Non-Hateful' and label_text == 'Hateful':
36         text_non_hateful_meme_hateful.append(image_id_meme)
37
38     # Save image IDs into separate text files
39     save_image_ids(both_hateful, 'both_hateful.txt')
40     save_image_ids(both_non_hateful, 'both_non_hateful.txt')
41     save_image_ids(text_hateful_meme_non_hateful, 'text_hateful_meme_non_hateful.txt')
42     save_image_ids(text_non_hateful_meme_hateful, 'text_non_hateful_meme_hateful.txt')
43
44     # Initialize lists to store image IDs for correct predictions of each model
45     meme_correct_hateful = []
46     meme_correct_non_hateful = []
47     text_correct_hateful = []
48     text_correct_non_hateful = []
49     combined_correct_hateful = []
50     combined_correct_non_hateful = []
51
52     # Compare predictions with true labels and categorize cases, storing image IDs accordingly
53     for meme_pred, text_pred, combined_pred in zip(meme_predictions, text_predictions, combined_predictions):
54         image_id_meme, label_meme = meme_pred
55         image_id_text, label_text = text_pred
56         image_id_combined, label_combined = combined_pred
57
58         true_label = true_labels.get(image_id_meme)
59         if true_label is not None:
60             if true_label == 1: # Hateful
61                 if label_meme == 'Hateful':
62                     meme_correct_hateful.append(image_id_meme)
63                 if label_text == 'Hateful':
64                     text_correct_hateful.append(image_id_meme)
65                 if label_combined == 'Hateful':
66                     combined_correct_hateful.append(image_id_meme)
67             elif true_label == 0: # Non-Hateful
68                 if label_meme == 'Non-Hateful':
69                     meme_correct_non_hateful.append(image_id_meme)
70                 if label_text == 'Non-Hateful':
71                     text_correct_non_hateful.append(image_id_meme)
72                 if label_combined == 'Non-Hateful':
73                     combined_correct_non_hateful.append(image_id_meme)
74
75     # Save image IDs of memes correctly matched with true labels for each model
76     save_image_ids(meme_correct_hateful, 'meme_correct_hateful.txt')
77     save_image_ids(meme_correct_non_hateful, 'meme_correct_non_hateful.txt')
78     save_image_ids(text_correct_hateful, 'text_correct_hateful.txt')
79     save_image_ids(text_correct_non_hateful, 'text_correct_non_hateful.txt')
80     save_image_ids(combined_correct_hateful, 'combined_correct_hateful.txt')
81     save_image_ids(combined_correct_non_hateful, 'combined_correct_non_hateful.txt')

```

The above lines of code saves the respective results of each into the relevant files.

4.3.4 Text Extractor from Meme:

```
1 import matplotlib.pyplot as plt
2 import cv2
3 import easyocr
4 from pylab import rcParams
5 from IPython.display import Image
6 rcParams['figure.figsize'] = 8, 16
```

The libraries and the dependencies required for the task. Note that this OCR model is pre-trained for this task.

```
1 reader = easyocr.Reader(['en'])
```

The model is storied in reader variable and is configured to work only on english texts.

```
1 Image("img/01245.png")
```



Figure 24:

Below is the original image before image-processing.

```
1 output = reader.readtext("img/01245.png")
2 output
```

Prints the detected text in the form of an array consisting of each of the detected text and their bounding-box coordinates.

```
1 ([[109, 0], [415, 0], [415, 35], [109, 35]],
2  'and that was the last',
3  0.7935187872980576),
4  ([[81, 24], [450, 24], [450, 72], [81, 72]],
5  'nativity play IV SOn was',
6  0.3808463801228905),
7  ([[97, 61], [437, 61], [437, 105], [97, 105]],
8  'invited t0 take partin',
9  0.4336272431372592)
```

The output for the given image.

```
1 cord = output[-1][0]
2 x_min, y_min = [int(min(idx)) for idx in zip(*cord)]
3 x_max, y_max = [int(max(idx)) for idx in zip(*cord)]
4 image = cv2.imread("img/01245.png")
5 cv2.rectangle(image, (x_min, y_min), (x_max, y_max), (0, 0, 255), 2)
6 plt.imshow(cv2.cvtColor(image, cv2.COLOR_BGR2RGB))
```

The final image where the given text to be detected from output array has a bounding box around the text.



Figure 25:

4.4 Difficulties faced while doing the tasks:

- Since Tensorflow makes use of GPU for computation, it needs sufficient memory in the gpu for a particular task. If there's insufficient storage in the GPU for computation, then manually kill some of the processes which consume the GPU at that point of time for Tensorflow to be able to work on the particular task.
- Could'nt configure TensorRT with Tensorflow which made me miss the task of removing text from images using Image Processing techniques.
- I don't really understand why the predictions made by the meme-based model keep changing all the time and hence affect the combined model as well. If I change the number of epochs while training the meme-based model, it's performance is bound to change which is expected, But even without changing the number of epochs the predictions change as evident by the change of the performance parameters of the model everytime I train and load and make predictions with the meme-based model. It's been discovered that in the report, the combined model mirrors the text-based model. This may not hold true all the time and this doesn't even have to do with the training dataset. This stems from the fact that the meme-based model give different predictions everytime. In the report, its been reported that the text-based model beats the meme-based model on all the performance parameters. This may not be true while considering this issue. The meme-based model at times can perform better than the text-based one and can have repercussions on the combined model in terms of how much of aspects of each model should the combined model incorporate in itself. However the logic and the workflow holds good, just that the fact that the combined model mirrors the text-based model everytime wont hold true.
- There was a problem in the site:[Link to Dataset download](#).The download button wasn't working for me and I instead took the dataset from facebook meme challenge on kaggle which is identical to that of the hatefulmemeschallenge site.

4.5 Areas for Potential Improvement:

- Instead of using the nltk based classifier for the text-based model which doesn't consider the context of the text while training, Bert model could be used. Bert is a deep-learning model which is more suitable for this purpose of text-classification since it allows for a contextual understanding of text unlike nltk. Although using Bert could be computationally more expensive.
- The proposed combined model, integrating elements from both the text-based and meme-based models, could be enhanced by leveraging multi-

modal models capable of processing both text and image modalities impartially.

- A better system of toxicity calculator can be built instead of just calculating the number of hateful words with total number of words in the sentence lacking contextual understanding.

4.6 Links to relevant sources:

- More info on various performance metrics - <https://towardsdatascience.com/a-look-at-precision-recall-and-f1-score-36b5fd0dd3ec>
- Object Detection Task-YOLOV8 documentation - <https://docs.ultralytics.com/modes/predict/key-features-of-predict-mode>
- YOLOV8-Object Detection - <https://yolov8.com/>
- Helped me tackle a problem regarding YOLOV8 - <https://github.com/ultralytics/ultralytics/issues/2868>
- Meme-classifier and Hateful-Meme classifier - <https://valueml.com/meme-classification-using-cnn-in-python/>
- Link to Dataset download - <https://www.kaggle.com/datasets/parthplc/facebook-hateful-meme-dataset/discussion/453370>
- Removing text from images using OCR techniques and inpainting - <https://towardsdatascience.com/removing-text-from-images-using-cv2-and-keras-ocr-24e7612ae4f4>
- Working with easy-ocr for text detection in images - <https://medium.com/@adityamahajan.work/easyocr-a-comprehensive-guide-5ff1cb850168>

GenAI tools like ChatGPT was used for aspects of Text-based model for Hateful meme classifier and to generate code for certain aspects of the tasks using suitable prompts.