# Detail Report on Email Classification for Support Team

## Introduction to the Problem Statement

The need for privacy and security in handling sensitive information has become increasingly critical in today's digital landscape. Many applications and services collect, store, and process personally identifiable information (PII), which can pose serious risks if mishandled. In the context of support emails, this problem is even more pertinent as such emails often contain sensitive data such as phone numbers, email addresses, credit card information, and more.

The objective of this project is to build a system that:

1. Classifies support emails into specific categories, such as "Incident," "Request," "Problem," and "Change."

2. Ensures that any PII in these emails is masked, ensuring the data remains secure and private.

This system is important for compliance with privacy regulations (like GDPR) and to safeguard against data breaches, while also providing automated classification for efficient email management in support systems.

## Approach Taken for PII Masking and Classification

### PII Masking

The first challenge in the project is handling PII. Personally identifiable information must be identified and replaced with placeholder text to ensure privacy. This masking process uses regular expressions (regex) to detect and replace specific patterns of PII such as:

- Phone numbers

- Email addresses

- Full names

- Dates of birth

- Credit card details

- CVV numbers

The regex patterns are designed to be flexible and comprehensive, allowing for the identification of both local and international formats for these entities.

### Email Classification

The second challenge is classifying the support emails into predefined categories. We used a pre-trained **RoBERTa** (Robustly Optimized BERT Pretraining Approach) model for sequence classification. RoBERTa has shown strong performance on many NLP tasks, making it an excellent choice for text classification.

We fine-tuned this model on a custom dataset of support emails to categorize the emails into one of the following labels:

- **Incident**

- **Request**

- **Problem**

- **Change**

Here are some key code snippets from the implementation:

# PII Masking Utility (utils.py)

This code defines the function to mask different types of PII using regex patterns.

```python
import re

PII_PATTERNS = {
    "phone_number": [
        r"\b(?:\+?\d{1,3}[-.\s]?)?(?:\d{1,4}[-.\s]?){2,5}\d{2,4}\b"
    ],
    "full_name": [
        r"\b(?:Mr|Ms|Mrs|Dr)\.?\s+[A-Z][a-z]{1,20}\s+[A-Z][a-z]{1,20}\b",  # With titles
        r"\b[A-Z][a-z]{1,20}\s+[A-Z][a-z]{1,20}\b",  # Stricter name pattern
    ],
    "email": [
        r"\b[a-zA-Z0-9_.+-]+@[a-zA-Z0-9-]+\.[a-zA-Z0-9-.]+\b"
    ]
}

def mask_pii(text: str):
    entities = []
    matches = []
    for entity_type, patterns in PII_PATTERNS.items():
        for pattern in patterns:
            for match in re.finditer(pattern, text):
                start, end = match.span()
                matches.append({
                    "start": start,
                    "end": end,
                    "classification": entity_type,
                    "entity": match.group()
                })
    matches.sort(key=lambda x: (x['start'], -x['end']))
    non_overlapping = []
```

```
    last_end = -1
    for m in matches:
        if m['start'] >= last_end:
            non_overlapping.append(m)
            last_end = m['end']

    masked_text = text
    offset = 0
    entities = []
    for m in reversed(non_overlapping):
        start, end = m['start'], m['end']
        replacement = f"[{m['classification']}]"
        masked_text = masked_text[:start] + replacement + masked_text[end:]
        entities.append({
            "position": [start, start + len(replacement)],
            "classification": m['classification'],
            "entity": m['entity']
        })
    entities.reverse()
    return masked_text, entities
```

# Model Selection and Training Details

## Model Selection: RoBERTa

For the email classification task, we selected **RoBERTa** (Robustly optimized BERT pretraining approach), which is a variant of the BERT (Bidirectional Encoder Representations from Transformers) model. RoBERTa has been proven to outperform BERT in many Natural Language Processing (NLP) tasks. This model is well-suited for handling text classification tasks like email categorization, as it excels in understanding context, relationships, and nuances in language.

The main advantages of using RoBERTa for this task are:

1. **Pre-training on a large corpus**: RoBERTa has been pre-trained on massive amounts of data, enabling it to understand complex linguistic patterns.

2. **Contextual Understanding**: RoBERTa excels in understanding the relationship between words in a sentence, even if they are far apart, which is crucial for email classification where the relevant information might be spread throughout the email.

3. **Transfer Learning**: Since RoBERTa has been pre-trained, fine-tuning it on our task-specific dataset requires fewer data points and results in faster convergence compared to training a model from scratch.

## Model Architecture

We used the **RoBERTaForSequenceClassification** architecture from Hugging Face's **Transformers** library. This model is specifically designed for text classification tasks, and it outputs a prediction for the class label of the given input text.

The key components of the RoBERTa model are:

- **Input Layer**: The model takes tokenized input text, typically a batch of sentences or a single sentence.

- **Transformer Layers**: RoBERTa uses multiple layers of attention mechanisms to compute contextualized representations of the input text.

- **Classification Layer**: A fully connected layer that maps the output of the transformer layers to a prediction for each class.

## Training Process

The training of RoBERTa for the classification task was done with the following steps:

1. **Data Preparation**:

   o We used **masked email data** to train the model. Masking the PII before training ensures that the model learns to classify the text based on the content and not sensitive information.

   o The email data was preprocessed into input features (input_ids and attention_masks) that are fed into the RoBERTa model.

   o Labels for the emails were prepared based on their categories (such as "Incident," "Request," "Problem," "Change").

2. **Fine-Tuning**:

   o Fine-tuning a pre-trained model allows us to adapt the model to our specific email classification task.

   o The model was fine-tuned using the **masked emails** as inputs and the **predefined categories** as labels.

   o The **RoBERTa tokenizer** was used to tokenize the input email text into tokens that the model can understand. The text was padded to a maximum length of 512 tokens to ensure consistent input size.

3. **Model Hyperparameters**:

   o **Batch Size**: 16 for training and 64 for evaluation.

   o **Learning Rate**: Set to optimize the model's performance during fine-tuning.

   o **Epochs**: 3 epochs were used, which provided a good balance between training time and model accuracy.

   o **Max Length**: We set the maximum length of the input sequence to 512 tokens to ensure that longer emails could be processed effectively.

4. **Training Loop**:
   The model was trained using the **Trainer** class from the Hugging Face **Transformers** library, which simplifies training and evaluation. The Trainer class takes care of the backward propagation and gradient updates during training.

5. **Loss Function**:
   The loss function used for this classification task is the **cross-entropy loss**, which is standard for multi-class classification tasks. This loss function measures the difference between the predicted class probabilities and the true labels, and the model adjusts its weights to minimize this loss.

6. **Validation**:
   We used a **validation dataset** to check the model's performance during training. This helps in detecting overfitting and adjusting the model's hyperparameters accordingly.

## Code for Model Training

The following snippet shows how the model was trained using the **Trainer** class:

```python
from transformers import RobertaForSequenceClassification, RobertaTokenizer, Trainer, TrainingArguments
import torch
from torch.utils.data import Dataset, DataLoader

# Prepare dataset for training
class EmailDataset(Dataset):
    def __init__(self, encodings, labels):
        self.encodings = encodings
        self.labels = labels

    def __getitem__(self, idx):
        item = {key: torch.tensor(val[idx]) for key, val in self.encodings.items()}
        item['labels'] = torch.tensor(self.labels[idx])
        return item

    def __len__(self):
        return len(self.labels)

# Initialize tokenizer and prepare tokenized data
tokenizer = RobertaTokenizer.from_pretrained("roberta-base")
train_encodings = tokenizer(df_masked_emails['masked_email'].tolist(), truncation=True, padding=True, max_length=512)
train_labels = df_masked_emails['category'].map(label_map).tolist()

train_dataset = EmailDataset(train_encodings, train_labels)

# Initialize the model
model = RobertaForSequenceClassification.from_pretrained("roberta-base", num_labels=len(label_map))

# Set up the Trainer with training arguments
training_args = TrainingArguments(
```

```
    output_dir='./fine_tuned_roberta',        # Output directory for saving model
    num_train_epochs=3,                 # Number of epochs
    per_device_train_batch_size=16,          # Training batch size
    per_device_eval_batch_size=64,           # Evaluation batch size
    logging_steps=100,                  # Log every 100 steps
    save_steps=500,                     # Save the model every 500 steps
)

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset
)

# Train the model
trainer.train()

# Save the fine-tuned model
trainer.save_model('./fine_tuned_roberta')
tokenizer.save_pretrained('./fine_tuned_roberta')
```

## Model Evaluation

After training, the model was evaluated on a validation set to assess its performance. The evaluation metrics typically include accuracy, loss, and F1 score. The Trainer class in Hugging Face also facilitates this evaluation with built-in functions.

```
# Evaluate the model
eval_results = trainer.evaluate()
print(f"Evaluation Results: {eval_results}")
```

## Inference

Once the model was trained and saved, we used the following code to classify new emails:

```
def predict_email_type(email_text, model, tokenizer, label_map):
    """Predicts the type of an email using the fine-tuned RoBERTa model."""
    inputs = tokenizer(email_text, return_tensors="pt", padding=True, truncation=True,
max_length=512)
    inputs = inputs.to(device)
    outputs = model(**inputs)
    predicted_label_id = torch.argmax(outputs.logits).item()

    # Reverse the label mapping to get the original label
    reverse_label_map = {v: k for k, v in label_map.items()}
    predicted_label = reverse_label_map[predicted_label_id]
```

```
    return predicted_label
```

# Email Classification Logic (app.py)

This code snippet shows how the pre-trained RoBERTa model is used to classify emails.

```python
from transformers import RobertaForSequenceClassification, RobertaTokenizer
import torch

class EmailClassifier:
    def __init__(self, model_path: str = "fine_tuned_roberta", categories: Optional[List[str]] =
None):
        self.device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
        self.model =
RobertaForSequenceClassification.from_pretrained(model_path).to(self.device)
        self.tokenizer = RobertaTokenizer.from_pretrained(model_path)
        self.model.eval()
        self.categories = categories or ["Incident", "Request", "Problem", "Change"]

    def classify(self, text: str) -> str:
        inputs = self.tokenizer([text], truncation=True, padding='longest', max_length=512,
return_tensors="pt")
        inputs = {k: v.to(self.device) for k, v in inputs.items()}
        with torch.no_grad():
            outputs = self.model(**inputs)
        logits = outputs.logits
        pred_idx = torch.argmax(logits, dim=1).item()
        return self.categories[pred_idx]
```

# FastAPI Integration (api.py)

This code snippet integrates the classification and PII masking functionality into an API using FastAPI.

```python
from fastapi import FastAPI, HTTPException
from pydantic import BaseModel
from typing import Dict
from app import EmailClassifier
from utils import mask_pii

app = FastAPI()

classifier = EmailClassifier()

class EmailRequest(BaseModel):
    email_body: str

@app.post("/classify", response_model=Dict)
```

```
async def classify_email(request: EmailRequest):
    try:
        masked_email, entities = mask_pii(request.email_body)
        category = classifier.classify(masked_email)
        return {
            "input_email_body": request.email_body,
            "list_of_masked_entities": [{"position": e['position'], "classification": e['classification'],
"entity": e['entity']} for e in entities],
            "masked_email": masked_email,
            "category_of_the_email": category
        }
    except Exception as e:
        raise HTTPException(status_code=500, detail=str(e))
```

# Deployment of Email Classification System with PII Masking on Hugging Face Spaces

## 1. Introduction to Deployment

In this section, we detail the deployment process of the **Email Classification System with PII Masking** to **Hugging Face Spaces** using **Docker SDK**. Hugging Face Spaces provides a platform for easily deploying machine learning models and APIs, and with the help of Docker, the application is containerized for efficient and scalable deployment.

## 2. Deployment Process

The deployment process involves pushing the code to **Hugging Face Spaces**, configuring the environment using Docker, and ensuring that the API is accessible from the provided space URL.

# 2.1 Preparing the Code for Deployment

1. **Code Components**:

   o **app.py**: Contains the main logic for email classification using the **RoBERTa** model.

   o **api.py**: Contains the FastAPI application, including the endpoints to classify emails and mask PII.

   o **utils.py**: Contains utility functions to mask PII using regular expressions.

   o **Dockerfile**: Contains instructions to build the Docker image that encapsulates the environment for running the application.

2. **Dockerfile Configuration**:

The following **Dockerfile** ensures the environment is set up for running the **FastAPI** application using **Uvicorn**:

```
FROM python:3.12-slim

# Upgrade system packages to mitigate vulnerabilities
RUN apt-get update && apt-get upgrade -y && apt-get clean

# Create a new user to avoid running as root
RUN useradd -m -u 1000 user
USER user

# Set environment path
ENV PATH="/home/user/.local/bin:$PATH"

# Set the working directory
WORKDIR /app

# Copy requirements and install dependencies
COPY --chown=user ./requirements.txt requirements.txt
RUN pip install --no-cache-dir --upgrade -r requirements.txt

# Copy the application code
COPY --chown=user . /app

# Command to run the application using Uvicorn
CMD ["uvicorn", "api:app", "--host", "0.0.0.0", "--port", "7860"]
```

- o **Base Image**: The image uses python:3.12-slim to minimize the image size.
- o **System Packages**: It ensures that the system is up-to-date and cleans up unnecessary files to reduce image size.
- o **User Creation**: The Docker container runs as a non-root user for security purposes.
- o **Dependencies**: The required Python packages are installed by copying requirements.txt and using pip.
- o **FastAPI Server**: The application is run with Uvicorn, which serves the FastAPI app on port 7860.

3. **requirements.txt**:
   The requirements.txt file should list all the necessary dependencies. Here is an example:

```
fastapi
uvicorn
transformers
```

```
torch
pydantic
```

4. **Model Files**:
   The pre-trained or fine-tuned model (fine_tuned_roberta) and tokenizer must be
   stored in the application directory or fetched from Hugging Face if not included in
   the code repository.

---

# 2.2 Deploying on Hugging Face Spaces

1. **Create a Hugging Face Space**:

   o   Navigate to the [Hugging Face Spaces](#) page and create a new space.

   o   Choose **Docker** as the SDK for deployment. You will be prompted to push
       the code to the space.

2. **Push Code to Hugging Face Space**:

   o   Clone the space repository provided by Hugging Face:

       ```
       git clone
       https://huggingface.co/spaces/sachinmosambe/email_classification
       cd email_classification
       ```

   o   Copy your **app.py**, **api.py**, **utils.py**, **Dockerfile**, **requirements.txt**, and
       model files into the cloned repository.

   o   Commit and push the changes to Hugging Face:

       ```
       git add .
       git commit -m "Initial commit of Email Classification API with PII
       Masking"
       git push
       ```

3. **Building and Deployment**:
   Once the code is pushed, **Hugging Face Spaces** will automatically build a Docker
   image based on the Dockerfile. This process includes:

   o   Installing the dependencies.

   o   Running the FastAPI app using **Uvicorn**.

   o   Exposing the API on port 7860.

4. **Accessing the Deployed API**:
   After the deployment is successful, you will receive a URL for your Hugging Face
   Space (e.g., https://sachinmosambe-email-classification.hf.space). You can access
   your API documentation at:

5. You can test the API and interact with it directly from the Swagger UI.

---

# 2.3 Metadata for Hugging Face Space

To ensure that the Hugging Face Space has the correct metadata and styling, include the following **metadata** in the Space settings:

```
license: mit
title: email_classification
sdk: docker
emoji: 🏆
colorFrom: gray
colorTo: indigo
```

- **license**: Specifies the project license (MIT in this case).

- **title**: The name of the project, i.e., email_classification.

- **sdk**: Defines the deployment method, which is Docker here.

- **emoji**: A visual identifier for your space, represented by a trophy emoji.

- **colorFrom** and **colorTo**: Defines the color gradient for the space's theme (from gray to indigo).

---

# 3. Accessing the Deployed API

After deployment, you can access your API by visiting the provided URL. For example:

https://sachinmosambe-email-classification.hf.space

This URL will provide access to:

- **API Endpoints**: Access and test the /classify endpoint.

- **Swagger UI**: An interactive documentation interface at /docs.

- **Health Check**: Use the /health endpoint to check the API's status.

## Challenges Faced and Solutions Implemented

1. **Challenge 1: PII Extraction Accuracy**

   o  **Problem**: Ensuring accurate extraction and masking of PII across different formats.

- - **Solution**: We employed multiple regex patterns for each type of PII (e.g., phone numbers, emails) and ensured that the regex patterns were optimized for different regional formats. This allowed the system to identify a wide range of PII accurately.

2. **Challenge 2: Overlapping PII**

   - **Problem**: In some cases, PII could overlap in the text, making it difficult to mask without affecting other entities.

   - **Solution**: We handled overlapping matches by sorting the identified PII entities by their start position and ensuring that only non-overlapping entities were kept. This ensured that masking did not interfere with the remaining text.

3. **Challenge 3: Balancing Privacy and Classification Accuracy**

   - **Problem**: Ensuring that PII masking did not impact the model's classification performance.

   - **Solution**: We masked the PII in such a way that the core context of the email remained intact. This ensured that the classification model could focus on the actual content of the email, without being influenced by the sensitive data.

4. **Challenge 4: Model Overfitting**

   - **Problem**: The model might overfit on the small dataset.

   - **Solution**: We used techniques like cross-validation and data splitting to prevent overfitting. We also carefully monitored the validation loss to ensure the model generalized well to unseen data.

---

## Deployment Information for Email Classification System with PII Masking

### 1. Hugging Face Repository

**Link:** https://huggingface.co/spaces/sachinmosambe/email_classification

### 2. API Endpoint URL for testing

**Link:** https://sachinmosambe-email-classification.hf.space/docs

### 3. GitHub Repository

**Link:** https://github.com/SachinMosambe/email_classification_system