# Submission 2

Problem Statement:

## Data Anonymizer

Name: *Megha Dhanya*

SRN: *PES2UG23CS336*

Section: *F*

## Abstract

This project implements a basic data anonymization system using Natural Language Processing techniques. A pre-trained Named Entity Recognition (NER) model is used to detect person names from unstructured text. Once identified, these names are anonymized by replacing them with the tag [REDACTED]. Such anonymization techniques are useful in privacy-sensitive domains where personal information must be protected before further data processing or analysis.

## Short Documentation

In this project, Named Entity Recognition (NER) is used to identify personal names within the raw text data. Used a pre-trained BERT-based NER model available through the Hugging Face Transformers library. The model is fine-tuned on the CoNLL-2003 dataset, which supports entity types such as Person (PER), Organization (ORG), and Location (LOC).

The implementation steps:

1. Load the NER pipeline using the Transformers library.
2. Provide a sample input text.
3. Detect named entities present in the text.
4. Filter entities of type PERSON (PER).
5. Replace detected person names with the tag [REDACTED] using string manipulation.

This approach shows how machine learning models can be combined with simple text processing techniques to achieve data anonymization.

# Sample Output / Screenshots

**Input Text:**

John signed the agreement on behalf of AlphaTech Solutions in Singapore.

**Detected Entity:**

- John = PERSON

**Anonymized Output:**

[REDACTED] signed the agreement on behalf of AlphaTech Solutions in Singapore.

(Screenshots of the NER output and anonymized result are included.)

## Sample Input Text

```
[3]: text = "John signed the agreement on behalf of AlphaTech Solutions in Singapore."
```

## Named Entity Recognition Output

```
[4]: entities = ner_pipeline(text)
     entities
```

```
[4]: [{'entity_group': 'PER',
       'score': np.float32(0.9986154),
       'word': 'John',
       'start': 0,
       'end': 4},
      {'entity_group': 'ORG',
       'score': np.float32(0.9988221),
       'word': 'AlphaTech Solutions',
       'start': 39,
       'end': 58},
      {'entity_group': 'LOC',
       'score': np.float32(0.9997867),
       'word': 'Singapore',
       'start': 62,
       'end': 71}]
```

## Anonymization Process

Detected entities are replaced with their corresponding labels.

```
[5]: anonymized_text = text

     for ent in entities:
         if ent["entity_group"] == "PER":
             anonymized_text = anonymized_text.replace(ent["word"], "[REDACTED]")

     print("Original Text:")
     print(text)
     print("\nAnonymized Text:")
     print(anonymized_text)
```

```
Original Text:
John signed the agreement on behalf of AlphaTech Solutions in Singapore.

Anonymized Text:
[REDACTED] signed the agreement on behalf of AlphaTech Solutions in Singapore.
```

## GitHub Link:

https://github.com/PES2UG23CS336/Gen-AI.git