# Configuration Manual
## Classification of PII and Non PII files using Machine learning (NER) for DLP

MSc Research Project
MSc Cyber Security

# Shivraj Prithviraj Mohite
Student ID: x21171327

School of Computing
National College of Ireland

Supervisor:     Prof. Apurva Vangujar

| | | | |
|---|---|---|---|
| **Student Name:** | Shivraj Prithviraj Mohite | | |
| **Student ID:** | x21171327 | | |
| **Programme:** | MSc Cyber Security | **Year:** | 2022-2023 |
| **Module:** | Research Project | | |
| **Lecturer:** | Prof. Apurva Vangujar | | |
| **Submission Due Date:** | 14th August 2023 | | |
| **Project Title:** | Classification of PII and Non PII files using Machine learning (NER) for DLP | | |
| **Word Count:** | 1407 | **Page Count:** 11 | |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

**Signature:** Shivraj Prithviraj Mohite

**Date:** 14th August 2023

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies) | □ |
| **Attach a Moodle submission receipt of the online project submission,** to each project (including multiple copies). | □ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | □ |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual
Classification of PII and Non PII files Using Machine learning (NER) for DLP

## Shivraj Prithviraj Mohite
## X21171327

# 1  Introduction

Welcome to the Configuration Manual for the PII and Non PII files Classifier Application. This guide empowers you to efficiently set up and utilize this tool, simplifying the process of analysing security vulnerabilities. In this manual we will guide you on how to train and finetune BERT model according to your dataset and detect PII and Non PII data when input is provided and later using the finetuned model to create a basic python application which could read a .txt file and determine whether it's a PII and Non PII file. In this configuration manual we will show from scratch how to install all the libraries needed, the software's and everything.

# 2  Configurations

## 2.1  Hardware

- Operating System: Windows 11
- Processor: Intel i7
- Architecture: 64bit
- Storage: 1TB SSD
- Memory: 16GB
- GPU: Nvidia RTX 3060

## 2.2  Software

- Visual Studio Code
- Python
- Google collab

## 2.3  Model Used

- **BERT**

# 3  Implementation

Prerequisites: Installing the latest python version along with Visual Studio Code IDE which is required for our implementation

Step 1: Creation of Dataset: Dataset consisted of 18lakh entries of finely labelled entities. The creation of the Dataset involved two steps, one was finding the format in which we wanted to

create a dataset and then was to generate a dataset as we could not find one because PII data is not available on internet. A fake PII data generator code was developed using Faker library in python which created our data then we used AI tool to put that created data into a format which was suitable for BERT to be trained on. This finely labelled dataset of Non PII and PII data was used to train BERT model and achieve good accuracy.. Ill share you the dataset file where we have the content and their respective which are required to train or finetune the BERT model. Below is the image of dataset and its format:

| A | B | C | D |
|---|---|---|---|
| Sentence # | Word | POS | Tag |
| Sentence: 1 | Thousands | NNS | O |
| | of | IN | O |
| | demonstrators | NNS | O |
| | have | VBP | O |
| | marched | VBN | O |
| | through | IN | O |
| | London | NNP | B-geo |
| | to | TO | O |
| | protest | VB | O |
| | the | DT | O |
| | war | NN | O |
| | in | IN | O |
| | Iraq | NNP | B-geo |
| | and | CC | O |
| | demand | VB | O |
| | the | DT | O |
| | withdrawal | NN | O |
| | of | IN | O |
| | British | JJ | B-gpe |
| | troops | NNS | O |
| | from | IN | O |
| | that | DT | O |
| | country | NN | O |
| | . | . | O |

Step 2: Open Google collab and import the Jupyter notebook file which I have uploaded: The name of the file is **Model_training.ipynb** which contains all the below code. This code is used to train of Pretrained BERT model to get good accuracy.



Step 3 : Upload the labelled dataset which I have uploaded in artefact on to the google collab notebook and then run all the commands in the Jupypter Notebook. When you run all the

commands given below it will start training the model for 8 epoch which means it will train the model 8 times over the same dataset.



Step 4 : Once the model is trained , you can pass the text input using model.predict function: To give a random text to the model which the model will classify and give you output. In my case I have given the text as "My name is Irene Adler i am from London and my birthdate is 07/30/1997 and SSN is 444-55-333" which is a random text and is not associated with any human.



As seen in the image above we gave a sentence to the model through model.predict function and the output was stored in prediction. When we displayed the output using prediction function it showed us how the model classifies each word and label its with an identifier. In the above scenario words like My, name are labelled as O which represents Others, London is labelled as B-Geo which represents Geo location, Irene Adler is tagged as Name(PII) , SSN is

tagged as Social Security number(PII), 07/30/1997 is tagged as Birthdate(PII) these all labels show PII data has been classified with great accuracy which is our main goal.

Step 5: Importing the BERT model and classifier into our local machine so we could use it further in our application or use it for any other purpose.



Step 6: Open the a new jupypter notebook and import the file named **final_gui.ipynb** which is code which creates a application where it imports the finetuned BERT model and gives us a user interface to interact with the model.



Here in this code snippet we can see we have imported the finetuned BERT model and its tokenizer by giving the path on your local computer where the files are present.

Step 7: Executing the code directly once the path of BERT model and tokenizer is set (Note: All the code has been explained in implementation part of the report)

This is the basic GUI of the application now you can use the load file button to upload a text file which the application will read and then will classify the content of the file using BERT and then give alert according to the contents present in the file.

Step 8: Running the test cases to see whether the model and the application is running correctly or not.

Test case 1 : We will load a text file which will have Non sensitive PII information
First ill show you the content of the file which we are going to load:



You can see that this file does not have Sensitive PII but has Non Sensitive PII, we will now check if it works correctly.

It gives alert that Non sensitive PII detected as the file contains name, birthdate etc. Hence the application is working fine.



It will show the detected entities as Birthdate(PII), B-Geo, Name (PII) and wont display the content of the file. In our scenario above 3 labels are classified as Non-Sensitive PII according to that the application is working accurately.

Test Case 2: Loading Sensitive File

Lets first show you content of the file we are going to load:

Here you can clearly see Sensitive PII information which is SSN number is present. Hence we will now test by uploading the file.



It shows that the file we loaded has Sensitive PII information present in it like social security number.

Here we can see in detected entities Social security number is detected which is considered as Sensitive PII which indicates our model is working fine.

Test Case 3 : Loading file which contains Non PII information.



This is the content of the file and as you can see there are is no sensitive information in that. And after loading it up on our application it should show that the file is non PII and show display its content

As you can see when we load the file the alert pops up saying its an Non PII file. And in the next screenshot you can see the contents of the file are displayed and classified.



Here we can actually see the content of loaded file along with its classification.