# README

June 5, 2021

## 1 Introduction

The topic of our mini project is building a text generator model, which is one of the Natural Language Processing fields. To d o s o, w e u sed a d ataset t hat c ontains d ifferent b log p osts by different people, which we preprocessed, separated according to age groups, and later used to built individual neural networks. As multiple files a re i nvolved i n e ach s tep, t his n otebook is intended to lay out the structure of the folder system of the project in hope of helping our readers guide through the files more conveniently.

### 1.1 Note:

As some of the file sizes involved are relatively large, they are not included as part of the submission file that we sent via email. Instead, we have put them in this Google Drive link, which has the same folder structure as the submission and all the related files that we used to run the notebooks. Therefore, if one is interested, please find/download them using this link.
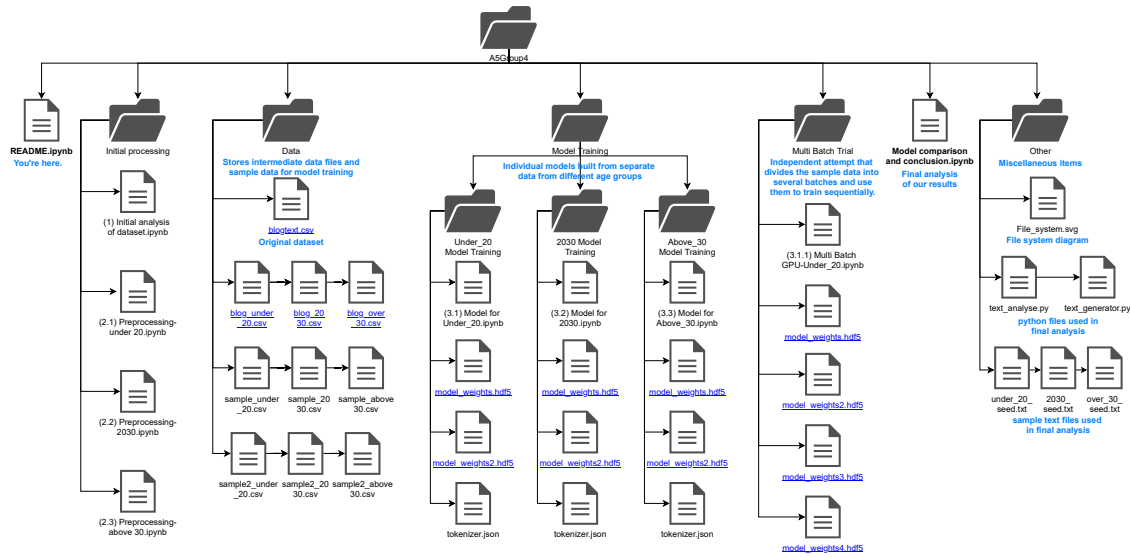
## 2 Structure of folder & files

Below is a diagram of the file system of this project; a corresponding Google Drive link is provided if it is not included in the submission file due to its size. (All the .pdf files are not shown here.)

### 2.1 Note:

In case the display is too small for reading in this notebook, please open **File_system.svg** inside the `Other` folder manually.

```
from IPython.core.display import SVG
SVG(filename='Other/File_system.svg')
```

A5Group4

README.ipynb
**You're here.**

Initial processing

Data
**Stores intermediate data files and sample data for model training**

Model Training
**Individual models built from separate data from different age groups**

Multi Batch Trial
**Independent attempt that divides the sample data into several batches and use them to train sequentially.**

Model comparison and conclusion.ipynb
**Final analysis of our results**

Other
**Miscellaneous items**

(1) Initial analysis of dataset.ipynb

(2.1) Preprocessing-under 20.ipynb

(2.2) Preprocessing-2030.ipynb

(2.3) Preprocessing-above 30.ipynb

blogtext.csv
**Original dataset**

blog_under _20.csv    blog_20 30.csv    blog_over _30.csv

sample_under _20.csv    sample_20 30.csv    sample_above 30.csv

sample2_under _20.csv    sample2_20 30.csv    sample2_above 30.csv

Under_20 Model Training

2030 Model Training

Above_30 Model Training

(3.1) Model for Under_20.ipynb

(3.2) Model for 2030.ipynb

(3.3) Model for Above_30.ipynb

model_weights.hdf5    model_weights.hdf5    model_weights.hdf5

model_weights2.hdf5    model_weights2.hdf5    model_weights2.hdf5

tokenizer.json    tokenizer.json    tokenizer.json

(3.1.1) Multi Batch GPU-Under_20.ipynb

model_weights.hdf5

model_weights2.hdf5

model_weights3.hdf5

model_weights4.hdf5

File_system.svg
**File system diagram**

text_analyse.py    text_generator.py
**python files used in final analysis**

under_20_ seed.txt    2030_ seed.txt    over_30_ seed.txt
**sample text files used in final analysis**

# 3 Recommended reading order

1. Subsequent to this notebook, one can start by looking at the `Initial processing` folder, which contains notebooks for analysis and preprocessing of the dataset:

   - **(1) Initial analysis of dataset.ipynb** provides an overview of the data we use and some data splitting.
   - The other three notebooks starting with (2.x) prefix are further data preprocessing for each age group.

**Note: As the three (2.x) notebooks have almost the same structure, only** *(2.1) Preprocessing-under20.ipynb* **will contain detailed explnations of the steps involved.**

2. The `Data` folder is simply output from the notebooks of the `Initial processing` folder and therefore is not of importance here.

3. The `Model Training` folder contains training noteboks and some related output files for each age group's data.

**Note: As the three training notebooks have almost the same structure, only** *(3.1) Model for Under_20.ipynb* **will contain detailed explnations of the steps involved.**

4. The `Multi Batch Trial` is intended to record one of our attempts that divided a large sample data file into different groups and used them to train the model sequentially. It's worth noting that this attempt was only applied on the sample data of the under-20 age group (please refer to the corresponding notebook for more detailed reasons). Moreover, the attempt is independent from the ones inside the `Model Training` and can be ignored if one is only interested in the main part of the project.

5. Finally, the **Model comparison and conclusion.ipynb** loads the previously trained models from the .hdf5 files within the folder `Model Training` and shows predictions of each model, from which some analysis and conclusions will be drawn.