

# Data Intake Report

Name: en-fr-translation-dataset

Report date: 12/06/2023

Internship Batch: LISUM21

Version: 1.0

Data intake by: Bogdan-Remus Pintilie Data

intake reviewer:

Data storage location: <https://www.kaggle.com/datasets/dhruvildave/en-fr-translation-dataset>

## Tabular data details:

|                                     |          |
|-------------------------------------|----------|
| <b>Total number of observations</b> | 22520376 |
| <b>Total number of files</b>        | 1        |
| <b>Total number of features</b>     | 2        |
| <b>Base format of the file</b>      | .csv     |
| <b>Size of the data</b>             | 8.2 GB   |

## Proposed Approach:

- Set up an environment
- Upload csv file
- Read the file using different methods (pandas, dask, modin)
- Use a smaller subset of the data, reading the entire file at once can be memory-intensive
- Measure the computational efficiency of each method and document findings
- Perform basic validation on data columns (remove special characters and white spaces)
- Create a YAML file defining the separator and column names
- Read the YAML file and validate the number of columns and column names
- Write the file in pipe-separated (“,”) text format in gzipped format

The screenshot shows the Visual Studio Code interface. The Explorer sidebar on the left displays a file tree for a 'NEW FOLDER' containing files: 'en-fr.csv', 'how\_long\_does\_it\_tak...', 'magicfile.yaml', 'number\_of\_rows.py' (selected), 'result.txt.gz', and 'week6\_code.py'. The main editor window shows the code for 'number\_of\_rows.py':

```
1 import pandas as pd
2
3 file_path = 'C:\\Users\\Bogdan\\Desktop\\New folder\\en-fr.csv'
4 chunk_size = 10000
5
6 total_rows = 0
7 for chunk in pd.read_csv(file_path, chunksize=chunk_size):
8     total_rows += len(chunk)
9
10 print(f"Total number of rows in the DataFrame:", total_rows)
```

Below the editor, the TERMINAL panel shows the command executed and its output:

```
PS C:\Users\Bogdan\Desktop\New folder> & C:/Users/Bogdan/AppData/Local/Programs/Python/Python311/python.exe "c:/Users/Bogdan/Desktop/new folder/number_of_rows.py"
Total number of rows in the DataFrame: 22520376
PS C:\Users\Bogdan\Desktop\New folder>
```

The status bar at the bottom indicates the cursor is at Line 10, Column 60, with 4 spaces, UTF-8 encoding, CRLF line endings, Python 3.11.4 64-bit, and the date 04-Aug-23.

The screenshot shows the Visual Studio Code interface. The Explorer sidebar on the left displays a file tree for a '2GB' folder containing files: 'en-fr.csv' and 'samafut.py' (selected). The main editor window shows the code for 'samafut.py':

```
1 import pandas as pd
2
3 df = pd.read_csv("D:\\New folder\\week 6\\2GB\\en-fr.csv", nrows=5)
4
5 print(df.head())
6
```

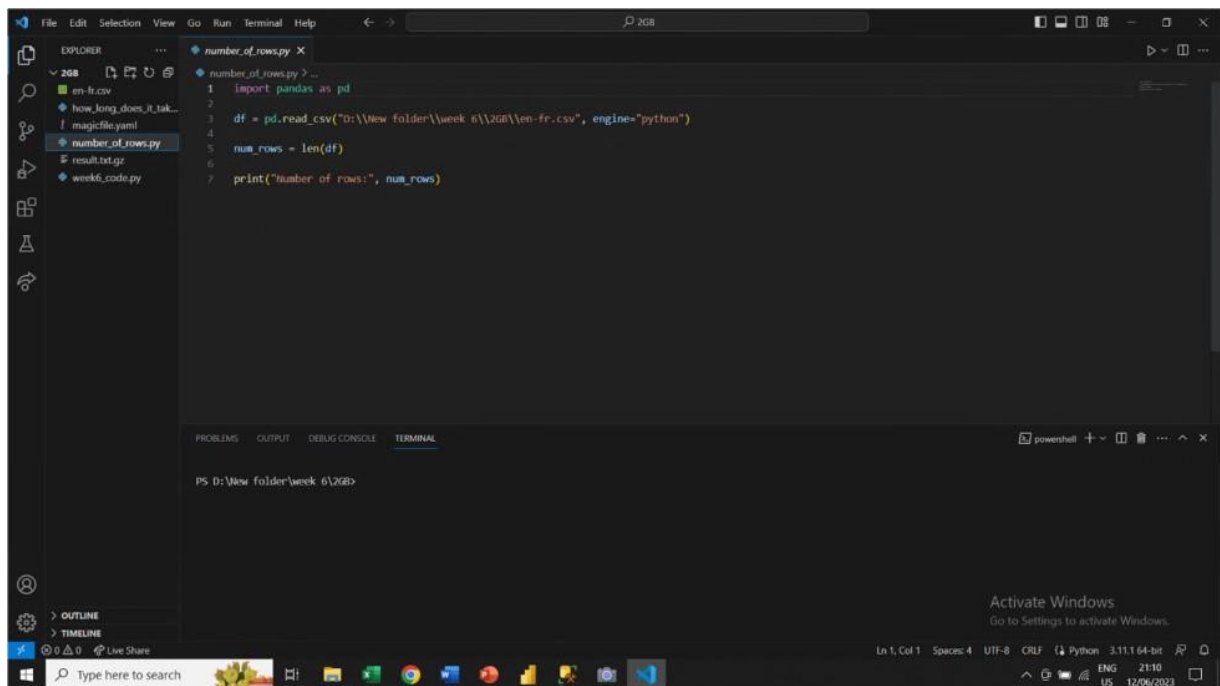
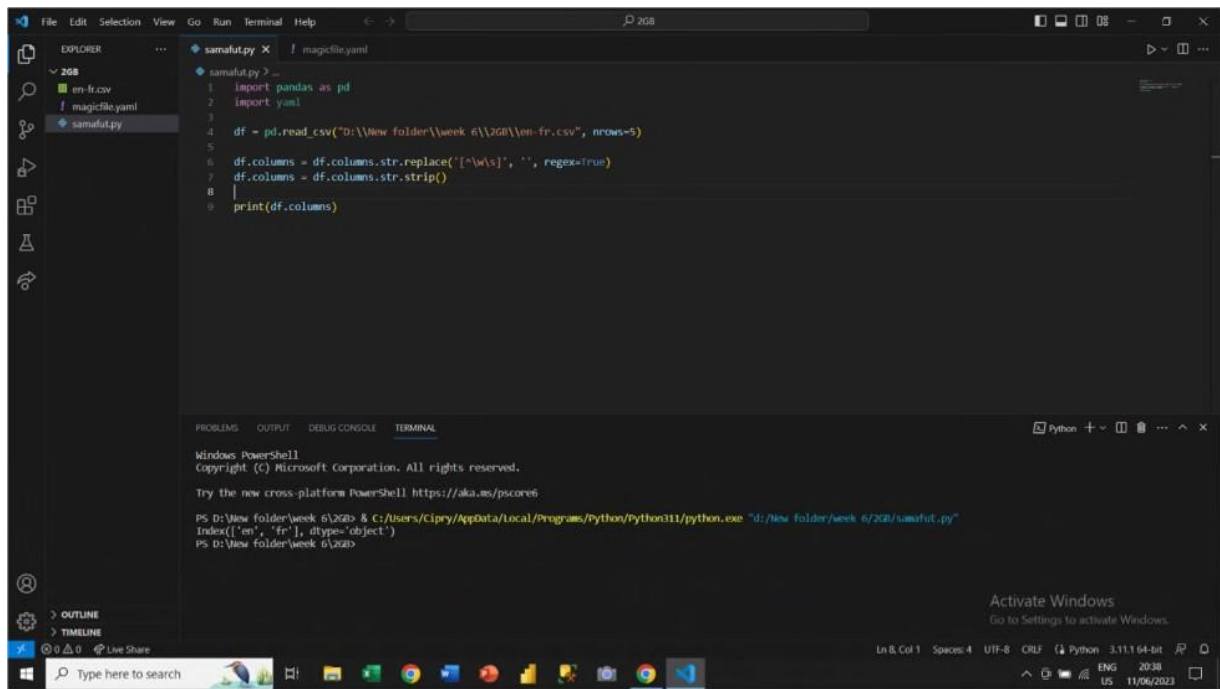
Below the editor, the TERMINAL panel shows the command executed and its output:

```
Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Try the new cross-platform PowerShell https://aka.ms/pscore6

PS D:\New folder\week 6\2GB> & C:/Users/Cipry/AppData/Local/Programs/Python/Python311/python.exe "d:/New folder/week 6/2GB/samafut.py"
en
fr
0  Changing Lives | changing Society | How It Wor...  Il a transformé notre vie | Il a transformé la...
1  Site map  Site map  Plan du site
2  Feedback  Feedback  Rétroaction
3  Credits  Credits  Crédits
4  Français  Français  Français
PS D:\New folder\week 6\2GB>
```

The status bar at the bottom indicates the cursor is at Line 5, Column 17, with 17 spaces, UTF-8 encoding, CRLF line endings, Python 3.11.1 64-bit, and the date 11/06/2023. An 'Activate Windows' watermark is visible in the bottom right corner.



The screenshot shows the Visual Studio Code editor interface. The Explorer sidebar on the left displays a file tree with the following items: `en-fr.csv`, `how_long_does_it_take.py` (selected), `magicfile.yaml`, `number_of_rows.py`, `result.txt.gz`, and `week6_code.py`. The main editor area displays the content of `how_long_does_it_take.py`, which is a Python script for benchmarking data loading performance using pandas, dask, and modin. The script includes imports for `pandas`, `dask.dataframe`, `modin.pandas`, and `time`. It defines a `file_path` and measures the execution time for reading a CSV file using each library. The terminal window at the bottom shows the current directory as `PS D:\New folder\week 6\2GB>`. The status bar at the bottom indicates the file is at Line 1, Column 1, with 4 spaces, using UTF-8 encoding and CRLF line endings. The system tray shows the date and time as 12/06/2023, 21:11.

```
1 import pandas as pd
2 import dask.dataframe as dd
3 import modin.pandas as mpd
4 import time
5
6 file_path = 'D:\\New folder\\week 6\\2GB\\en-fr.csv'
7
8 start_time = time.time()
9 df_pandas = pd.read_csv(file_path)
10 end_time = time.time()
11 execution_time_pandas = end_time - start_time
12
13 start_time = time.time()
14 df_dask = dd.read_csv(file_path, nrows=100)
15 end_time = time.time()
16 execution_time_dask = end_time - start_time
17
18 start_time = time.time()
19 df_modin = mpd.read_csv(file_path)
20 end_time = time.time()
21 execution_time_modin = end_time - start_time
22
23 print(f"Pandas Execution Time: {execution_time_pandas} seconds")
24 print(f"Dask Execution Time: {execution_time_dask} seconds")
25 print(f"Modin Execution Time: {execution_time_modin} seconds")
```

PS D:\New folder\week 6\2GB>

The screenshot shows the Visual Studio Code editor interface with the `magicfile.yaml` file selected in the Explorer sidebar. The main editor area displays the content of `magicfile.yaml`, which is a simple YAML configuration file. The terminal window at the bottom shows the current directory as `PS D:\New folder\week 6\2GB>`. The status bar at the bottom indicates the file is at Line 5, Column 1, with 2 spaces, using UTF-8 encoding and CRLF line endings. The system tray shows the date and time as 12/06/2023, 21:11.

```
1 separator: ','
2 columns:
3   - en
4   - fr
5
```

PS D:\New folder\week 6\2GB>

