



Software Engineering Department  
Braude College

## **User Guide**

# **EEG Classification Using Text Compression**

**25-1-R-2**

By:- Mohammad khateeb  
Jad taha

Advisors:- Dr. Samah Idrees Ghazawi  
Dr. Anat Dahan

Link to github:-  
<https://github.com/mhmdkh1905/EEG-recordings.git>

- **Purpose**

This user guide provides clear operational instructions for running the EEG classification pipeline using Normalized Compression Distance (NCD) with ZB2 compression. The guide walks the user through the key steps required to prepare data, execute the classification workflow, and analyze the output in Google Colab.

- **System Overview**

The system performs classification of EEG signals to distinguish between ADHD and control participants. It relies on text-based similarity measures using compression techniques rather than machine learning. EEG recordings are preprocessed, segmented, transformed into symbolic form, compared using NCD, and then classified based on similarity statistics.

- **Requirements**

- Google Colab (recommended: Pro with high RAM)
- Python 3.x (via Colab)
- Internet connection
- Google Drive access with the project folders and files
- Required Python packages (installed within the notebook):
  - numpy, pandas, matplotlib, scipy, zlib, os, glob, tqdm, xlswriter

- **Folder Location**

Your project should be located in:  
[/content/drive/MyDrive/finalProject/](#)

Ensure the following subdirectories exist:

- [/adhdcsv/](#) - Original EEG CSV files for ADHD group
- [/controlcsv/](#) - Original EEG CSV files for Control group
- [/filteredadhdcsv/](#) - Filtered EEG files (1–40 Hz) for ADHD group
- [/filteredcontrolcsv/](#) - Filtered EEG files (1–40 Hz) for Control group
- [/brainwave\\_sequence\\_2s/](#), [/brainwave\\_sequence\\_5s/](#), [/brainwave\\_sequence\\_6s/](#), [/brainwave\\_sequence\\_8s/](#), [/brainwave\\_sequence\\_10s/](#) – Brainwave region sequences extracted using different time window sizes
- [/parts\\_ncd\\_2s/](#), [/parts\\_ncd\\_5s/](#), [/parts\\_ncd\\_6s/](#), [/parts\\_ncd\\_8s/](#), [/parts\\_ncd\\_10s/](#) – Contains the computed NCD values for each part-to-part comparison between participants
- [/min/](#), [/avg/](#), [/median/](#) – Excel files summarizing the minimum, average, and median NCD values per participant pair
- [/classification\\_based\\_on\\_median/](#) – Contains per-participant classification results for each method (min, avg, median) using median-based group assignment
- [/high\\_scores\\_summary/](#) – Includes Excel files and graphs showing all classification scores  $\geq 60\%$  for each participant and method

- **Files Provided**

- `finalProject2.ipynb` – The main notebook for signal processing, part-to-part NCD computation, and classification analysis.
- `/adhdcsv/` – Contains the raw EEG CSV files for ADHD group participants.
- `/controlcsv/` – Contains the raw EEG CSV files for Control group participants.
- ZB2 compression module – integrated directly within the notebook for computing NCD values.

- **Operating Instructions**

1. **Open Google Colab** and upload `finalProject2.ipynb`.

2. **Mount Google Drive and Install Dependencies**

- Go to the cell titled:

- `#Install libraries, Import modules, Mount Google Drive`

- Run this cell to set up the environment and link to your Drive.

3. **Verify Folder Structure**

- Ensure the following directories exist in your Google Drive under `/finalProject/`:

- `/adhdcsv/`, `/controlcsv/`
- `/filteredadhdcsv/`, `/filteredcontrolcsv/`
- `/brainwave_sequence_2s/`, `/brainwave_sequence_5s/`,  
`/brainwave_sequence_6s/`, `/brainwave_sequence_8s/`,  
`/brainwave_sequence_10s/`
- `/parts_ncd_2s/`, ..., `/parts_ncd_10s/`
- `/classification_based_on_median/`, `/high_scores_summary/`

4. **Set Main Project Path**

- Run the cell titled:

- `#Main project folder.`

- Confirm the path is correct and accessible (e.g., `/content/drive/MyDrive/finalProject/`).

5. **Filter EEG Signals (1–40 Hz)**

- Section:

- `# **preprocessing functions**`

- Cell:

- `# Bandpass filter settings`

- Run all related cells to clean and save filtered data into:  
`/filteredadhdcsv/` and `/filteredcontrolcsv/`

6. **Extract Brainwave Regions (per time window)**

- Section:

- `**Extracting Dominant EEG Band Sequences from Filtered Signals Using a Sliding Window Approach**`

- Run all 5 cells that set `output_path` for each version (2s, 5s, 6s, 8s, 10s)

- This will create the `/brainwave_sequence_Xs/` folders

## 7. Compute NCD Between Parts

➤ Section:

# **\*\*NCD Funtcion\*\***

→ Run the NCD functions for part-to-part comparisons

➤ Then go to the section:

# **dividing to parts**

→ Run all cells to divide the brainwave sequences into 1000-character parts and apply NCD

## 8. Calculate min / avg / median NCD per participant

➤ Sections to run:

- # Classification using the parts and average
- # Classification using the parts and median
- # Classification using the parts and min

→ These cells generate participant-to-participant scores and save them to Excel

## 9. Final Classification Based on Median Method

➤ Section:

**\*\*Classification based on median for (avg, median, min)\*\***

→ Run all cells to generate classification decisions for each method

→ Output is saved in **/classification\_based\_on\_median/**

## 10. Generate Final Graphs and Visualizations

➤ Section:

**\*\*graphs\*\***

→ Run all cells to produce bar charts and summaries

→ Output is saved in **/high\_scores\_summary/**

### • Outputs

#### ○ Filter EEG Signals (1–40 Hz)

After running the bandpass filter step, the system generates filtered EEG signals for each participant. These signals are stored as **.csv** files, each named in the format:

**filtered\_vXp.csv** where **X** represents the participant number (e.g., **filtered\_v121p.csv** for participant 121).

These filtered files retain all 19 EEG channels and contain only the frequency range between 1 Hz and 40 Hz, effectively removing noise and irrelevant components from the original signal.



brainwave activity (e.g., 'D' for Delta, 'T' for Theta, 'B' for Beta), capturing the EEG structure in a simplified, compression-friendly form.

○ **Compute NCD for Each Part-to-Part**

When you run this step, the system calculates the Normalized Compression Distance (NCD) between each pair of parts from two participants' signals for a specific channel. These part-to-part comparisons produce a CSV file that lists every possible match between segments, including the computed similarity value (NCD).

Each output file is named:

vXp\_vs\_vYp\_channelZ.csv

Where:

- X is the first participant number
- Y is the second participant number
- Z is the EEG channel index

Each row in the file corresponds to:

- Part #X and its symbolic sequence
- Part #Y and its symbolic sequence
- The NCD value between the two parts

A	B	C	D	E
Part #120p	Sequence 120p	Part #121p	Sequence 121p	NCD Value
1	DDDDDDDDDDDDDDDDDD	1	DDDDDDDDDDDDDDDDDD	0.3582089552
1	DDDDDDDDDDDDDDDDDD	2	DDDDDDDDDDDDDDDDDD	0.3582089552
1	DDDDDDDDDDDDDDDDDD	3	TTTTTTTTTTTTTTTTTTTT	0.3432835821
1	DDDDDDDDDDDDDDDDDD	4	TTDDDDDDDDDDDDDDDTT	0.44
1	DDDDDDDDDDDDDDDDDD	5	DDDDDDDDDDDDDDDDDD	0.3529411765
1	DDDDDDDDDDDDDDDDDD	6	DDDDDDDDDDDDDDDDDD	0.3880597015
1	DDDDDDDDDDDDDDDDDD	7	DDDDDDDDDDDDDDDDDD	0.3432835821
1	DDDDDDDDDDDDDDDDDD	8	DDDDDDDDDDDDDDDDDD	0.4
1	DDDDDDDDDDDDDDDDDD	9	DTTTTTTTTDDDDDTTTTTTT	0.328358209
1	DDDDDDDDDDDDDDDDDD	10	DDDDDDDDDDDDDDDDDD	0.2985074627
1	DDDDDDDDDDDDDDDDDD	11	DDTTTTTTTDDDDDDDDDDDD	0.328358209
1	DDDDDDDDDDDDDDDDDD	12	DDDDDDDDDDDDDDDDDD	0.3880597015
2	DDDDDDDDDDDDDDDDDD	1	DDDDDDDDDDDDDDDDDD	0.4242424242
2	DDDDDDDDDDDDDDDDDD	2	DDDDDDDDDDDDDDDDDD	0.3939393939
2	DDDDDDDDDDDDDDDDDD	3	TTTTTTTTTTTTTTTTTTTT	0.3939393939
2	DDDDDDDDDDDDDDDDDD	4	TTDDDDDDDDDDDDDDDTT	0.4933333333
2	DDDDDDDDDDDDDDDDDD	5	DDDDDDDDDDDDDDDDDD	0.4558823529
2	DDDDDDDDDDDDDDDDDD	6	DDDDDDDDDDDDDDDDDD	0.4393939394
2	DDDDDDDDDDDDDDDDDD	7	DDDDDDDDDDDDDDDDDD	0.3939393939
2	DDDDDDDDDDDDDDDDDD	8	DDDDDDDDDDDDDDDDDD	0.4571428571
2	DDDDDDDDDDDDDDDDDD	9	DTTTTTTTTDDDDDTTTTTTT	0.3939393939
2	DDDDDDDDDDDDDDDDDD	10	DDDDDDDDDDDDDDDDDD	0.3939393939
2	DDDDDDDDDDDDDDDDDD	11	DDTTTTTTTDDDDDDDDDDDD	0.4090909091
2	DDDDDDDDDDDDDDDDDD	12	DDDDDDDDDDDDDDDDDD	0.4393939394
3	DDDDDDDDDDDDDDDDDD	1	DDDDDDDDDDDDDDDDDD	0.3846153846
3	DDDDDDDDDDDDDDDDDD	2	DDDDDDDDDDDDDDDDDD	0.4
3	DDDDDDDDDDDDDDDDDD	3	TTTTTTTTTTTTTTTTTTTT	0.4
3	DDDDDDDDDDDDDDDDDD	4	TTDDDDDDDDDDDDDDDTT	0.52
3	DDDDDDDDDDDDDDDDDD	5	DDDDDDDDDDDDDDDDDD	0.4264705882
3	DDDDDDDDDDDDDDDDDD	6	DDDDDDDDDDDDDDDDDD	0.3692307692

The image shows an output CSV with symbolic sequences from Participant 120 and Participant 121 for Channel 18 (2-second version). The final column contains the NCD similarity values. These values quantify how structurally similar two sequences are based on compression—lower values imply higher similarity.

- **Calculate min / avg / median NCD per participant**

When you run this step, the system summarizes the NCD values between a selected participant and all other participants for a specific channel, based on the chosen method: minimum, average, or median. Each output is a .txt file containing:

- The filename of the original NCD comparison file
- The computed NCD score using the selected method

Each file is named and stored in the following structure: /<method>/<Xs>/channelY/participantZ.txt  
Where:

- method is one of avg, min, or median
- X is the window size (2s, 5s, 6s, 8s, 10s)
- Y is the channel number
- Z is the participant number

```
Participant 121 - Channel 18 (Min-Based)
v76p_vs_v121p_channel18.xlsx: 0.02326
v77p_vs_v121p_channel18.xlsx: 0.02326
v78p_vs_v121p_channel18.xlsx: 0.02326
v79p_vs_v121p_channel18.xlsx: 0.02326
v80p_vs_v121p_channel18.xlsx: 0.02326
v81p_vs_v121p_channel18.xlsx: 0.02326
v82p_vs_v121p_channel18.xlsx: 0.02326
v83p_vs_v121p_channel18.xlsx: 0.02326
v84p_vs_v121p_channel18.xlsx: 0.02326
v85p_vs_v121p_channel18.xlsx: 0.02326
v86p_vs_v121p_channel18.xlsx: 0.02326
v87p_vs_v121p_channel18.xlsx: 0.02326
v88p_vs_v121p_channel18.xlsx: 0.02326
v89p_vs_v121p_channel18.xlsx: 0.02326
v90p_vs_v121p_channel18.xlsx: 0.02326
v91p_vs_v121p_channel18.xlsx: 0.02326
v92p_vs_v121p_channel18.xlsx: 0.02326
v93p_vs_v121p_channel18.xlsx: 0.02326
v94p_vs_v121p_channel18.xlsx: 0.02326
v95p_vs_v121p_channel18.xlsx: 0.02326
v96p_vs_v121p_channel18.xlsx: 0.02326
v97p_vs_v121p_channel18.xlsx: 0.02326
v98p_vs_v121p_channel18.xlsx: 0.02326
v99p_vs_v121p_channel18.xlsx: 0.02326
v100p_vs_v121p_channel18.xlsx: 0.02326
v101p_vs_v121p_channel18.xlsx: 0.02326
v102p_vs_v121p_channel18.xlsx: 0.02326
v103p_vs_v121p_channel18.xlsx: 0.02326
v104p_vs_v121p_channel18.xlsx: 0.02326
v105p_vs_v121p_channel18.xlsx: 0.02326
v106p_vs_v121p_channel18.xlsx: 0.02326
v107p_vs_v121p_channel18.xlsx: 0.02326
v108p_vs_v121p_channel18.xlsx: 0.02326
v109p_vs_v121p_channel18.xlsx: 0.02326
v111p_vs_v121p_channel18.xlsx: 0.02326
```



The image displays the contents of [participant\\_121.txt](#) located under [/min/2s/channel18/](#). It lists NCD scores between participant 121 and other participants using the minimum-based method for each comparison. All values are consistent here (0.02326), highlighting similarity patterns across comparisons.

#### ○ **Final classification Based on median method**

When you run this step, the system generates an Excel file for each participant that summarizes their classification scores across all EEG channels using three different methods: average (avg), median, and minimum (min). These scores reflect how well the participant's EEG signals align with others from the same or opposite group, based on the chosen metric.

Each output file follows the format:

participant\_[X](#).xlsx , Where [X](#) is the participant number.

These files are stored in the folder:

[/classification\\_based\\_on\\_median/version/channelY/](#)

A	B	C	D	E	F	G	H	I	J	K
	Channel 0	Channel 1	Channel 2	Channel 3	Channel 4	Channel 5	Channel 6	Channel 7	Channel 8	Channel 9
avg	0.49167	0.55	0.55833	0.55833	0.625	0.575	0.60833	0.525	0.575	0.525
median	0.50833	0.50833	0.53333	0.525	0.53333	0.60833	0.55	0.48333	0.60833	0.55
min	0.49167	0.5	0.49167	0.50833	0.49167	0.48333	0.5	0.5	0.51667	0.50833

The image shows a classification summary table from [participant\\_121.xlsx](#), containing scores across channels 0–18. The rows correspond to the three methods ([avg](#), [median](#), [min](#)), and the columns represent EEG channels. Higher values typically indicate stronger classification confidence.

#### ○ **High Score Summary & Graphs**

When you run this step, the system generates an Excel file summarizing all classification results where the score is  $\geq 60\%$ . The file includes:

- The method used ([avg](#), [median](#), or [min](#))
- The segmentation window size (e.g., [2s](#), [5s](#), [10s](#))
- The participant number
- The EEG channel on which the score was calculated
- The classification score

The goal is to track strong classification results across all conditions and help analyze patterns of performance.

These results are saved in the folder: [/high\\_scores\\_summary/](#)

In addition to the Excel file, this step generates bar chart graphs visualizing:

- Number of high scores per participant
- Method-wise comparison
- Performance across segmentation versions



A	B	C	D	E	F
score number	version	score	participant	channel	method
1290	8s	0.68333	110	14	avg
32	2s	0.675	63	14	median
58	2s	0.675	65	14	median
301	2s	0.675	89	14	median
429	2s	0.675	106	14	avg
506	2s	0.675	115	8	avg
551	2s	0.675	120	14	avg
1138	8s	0.675	84	14	avg
1166	8s	0.675	90	14	avg
1320	8s	0.675	115	14	avg
1438	10s	0.675	90	8	avg
13	2s	0.65833	26	4	min
53	2s	0.65833	65	14	avg
63	2s	0.65833	66	14	avg
77	2s	0.65833	67	14	median
122	2s	0.65833	72	14	avg
156	2s	0.65833	75	14	avg
247	2s	0.65833	84	14	avg
262	2s	0.65833	85	14	median
273	2s	0.65833	86	14	median
284	2s	0.65833	87	14	median
380	2s	0.65833	100	14	median
387	2s	0.65833	101	14	avg
500	2s	0.65833	114	14	avg
510	2s	0.65833	115	8	median
525	2s	0.65833	117	8	median

This image shows a portion of the [high\\_scores\\_summary.xlsx](#) file. Each row represents a high-scoring classification instance. You can see the segmentation version, the score, the participant, the EEG channel, and the method used.