

## Dataset

### Step one: Combine subject and assessment and plan for original long dialogue

- **Dataset:** Data/ACIBENCH\_train\_processed, ACIBENCH\_validation\_processed and dialogue-processed

## Step Two: Data cleaning and data split into train/val/test

- **Dataset:** Data/Long\_dialogue/train\_long.csv, test\_long.csv, val\_long.csv
- **Dataset:** Data/Short\_dialogue/train\_short.csv, test\_short.csv, val\_short.csv
- **Notebook:** Notebook/Data\_preprocessing and analysis/Data\_Split
- **Note:** These are the dataset to run the **baseline model**

### Step Three: BioBert/Key Bert/Clinical NER extraction

- **Dataset:** all the data in folder **long\_dialogue\_NER\_extraction** and **short\_dialogue\_NER\_extraction**
- **Note:** Each NER model has it's own test/train/val dataset,
- **Note:** output for NER model are tokens, further merge into word and phrase per IBO tagging

### Step Four: BioBert/Key Bert/Clinical NER cleaning

- **Dataset:** all the data in folder **long\_dialogue\_NER\_cleaning** and **short\_dialogue\_NER\_cleaning**
- **Notebook:** Notebook/NER/ALL\_NER\_Cleaning
- NER\_cleaning
  - some NER appears several times and decided to just keep **unique** entities and keep **all entities**
  - keep the NER with it's **label** and **without it's label**
  - Below is a example for the for BioBert

[illegible]

**Bio\_ner\_label:** this is with all entities and their label (list of dictionary)

**Bio\_ner\_unqiuqe\_label:** this is the unique and their label (list of dictionary)

**Bio ner no label:** This is all entities without label (list of string)

**Bio\_ner\_unique\_no\_label:** This is unique entities without label (list of string)

- **Note:** These datasets are used for improved model