

Learning Adolescent Mental Health from Speech Data

Yuqi Wang and Shurui Bai

iBHE, Tsinghua University, China

ABSTRACT

- We proposed a model in adolescent mental health identification on speech data.
- Based on LSTM/GRU, this model can learn contextual features among utterances while using the features of a single speech, so as to judge the mental health of adolescents.
- Through experiments, our model can achieve good classification results.

INTRODUCTION

- **Adolescence**: a critical period for physical and mental development.
- At present, there are **few** deep learning researches in adolescent mental health, and **very few** related studies using speech data.
- **Current** screening method: combination tests of mental health questions (low pertinence, easy to deceive).
- **Deep learning**: to extract audio data features and analyze (accurate, time-saving, meaningful).

METHOD

- **Data Preprocessing** (Figure 1, 2)
- **Feature Extraction of each Uttrance** (Figure 3)
- **Model**: Contextual-LSTM/GRU/BiLSTM/BiGRU (Figure 4)

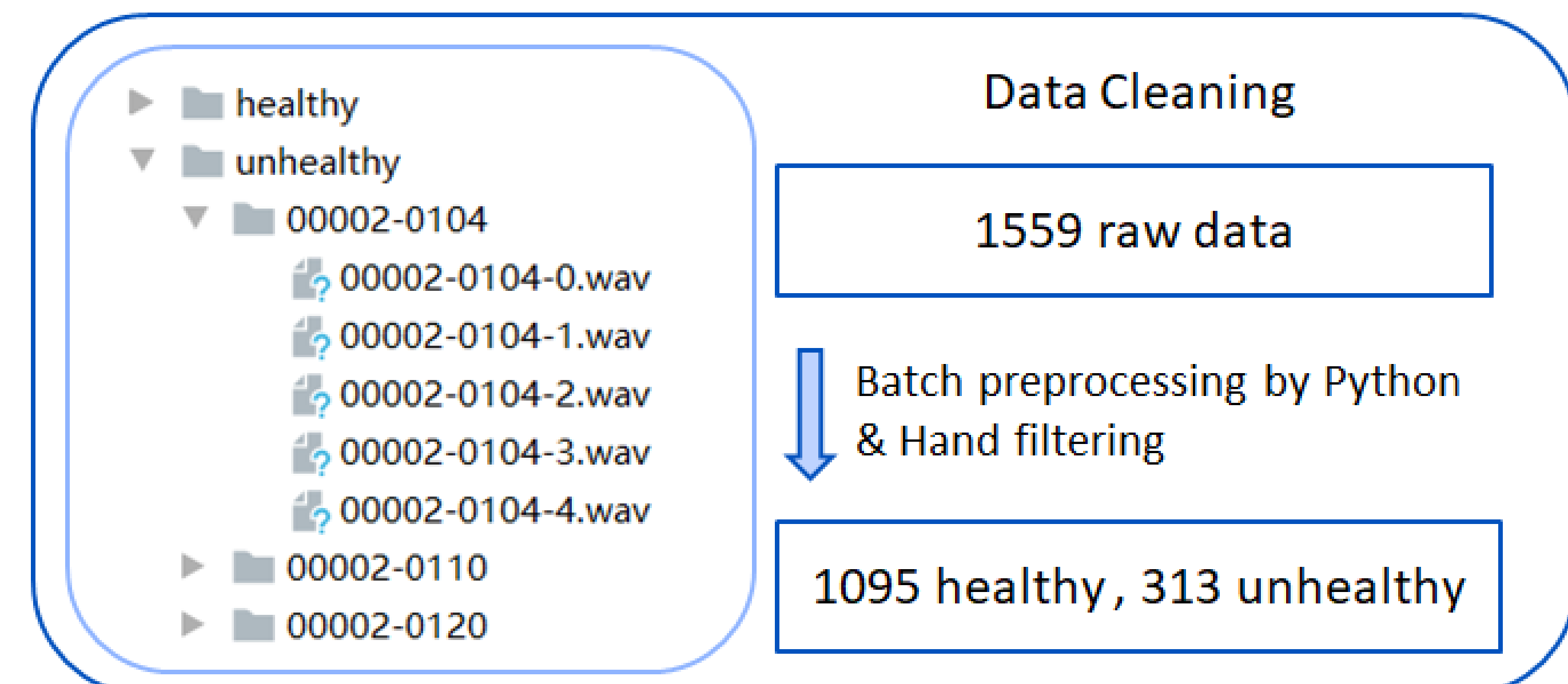


Figure 1. Dataset and Data Cleaning

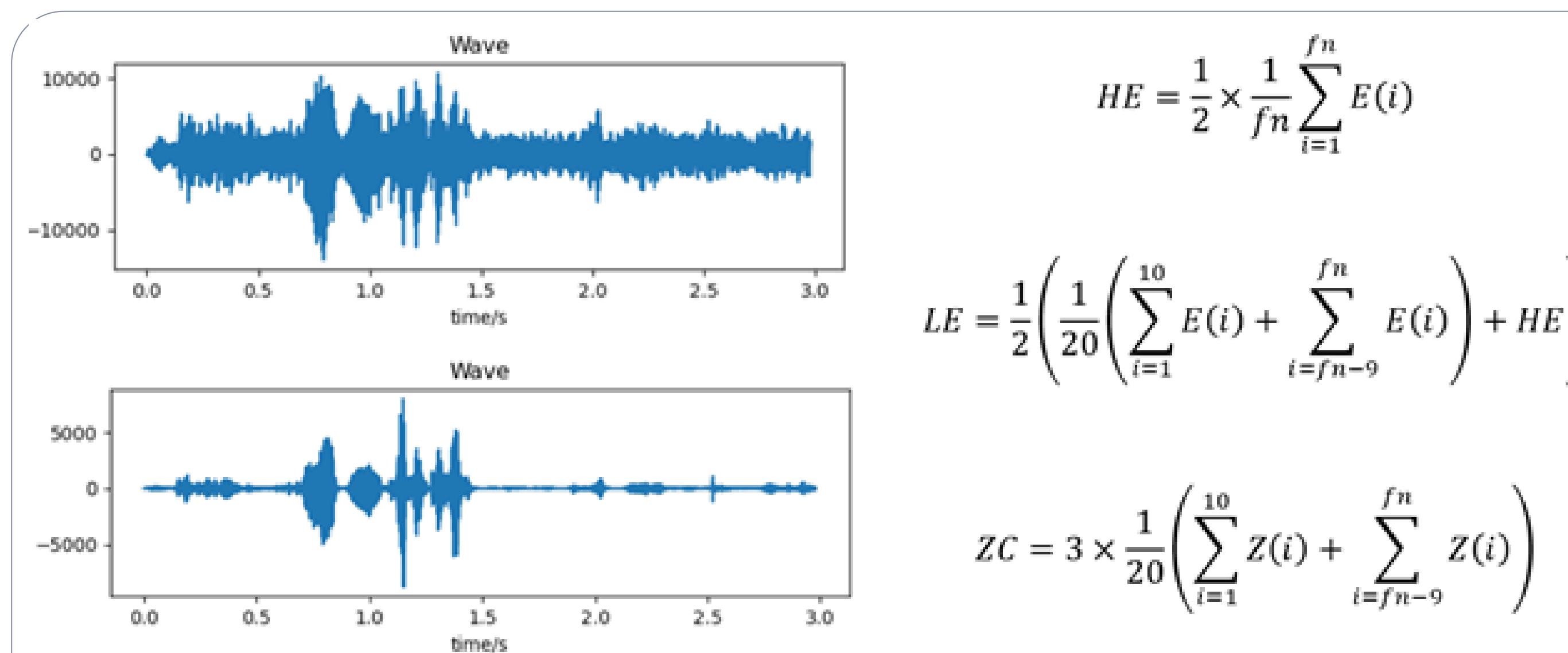


Figure 2. Noise Reduction & Endpoint Detection

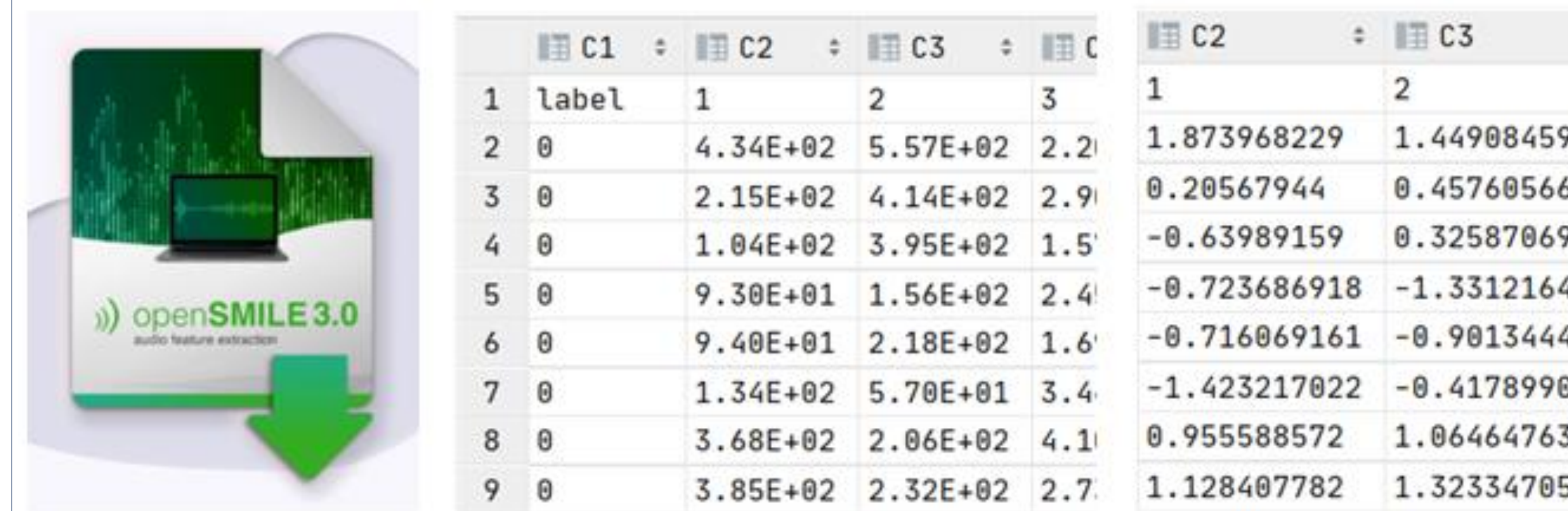


Figure 3. Feature Extraction of a Single Speech (a) Software tool: openSMILE3.0; b) Extracted IS10 speech features; c) Features after Z-score standardization)

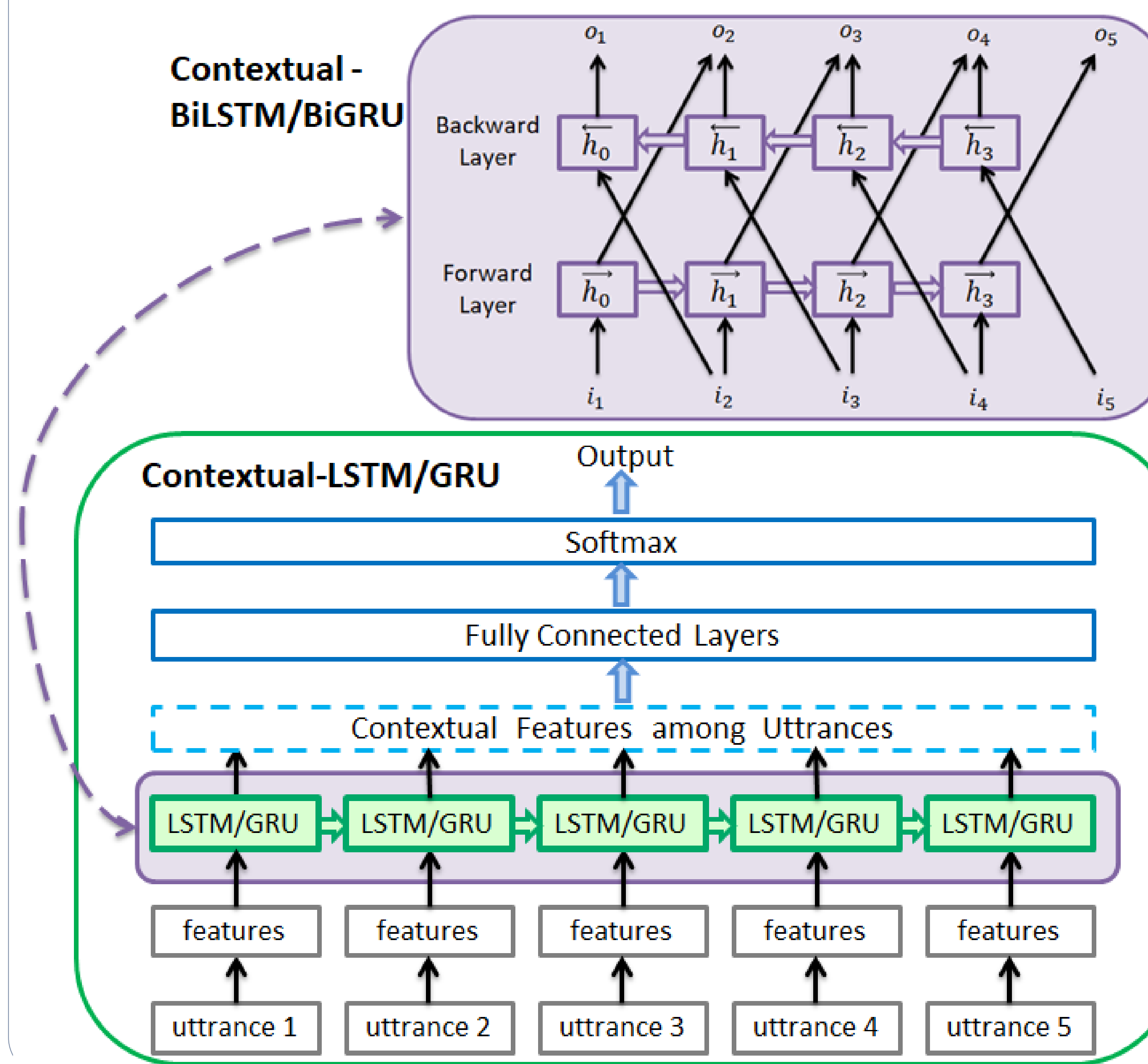


Figure 4. Model: Contextual-LSTM/GRU/BiLSTM/BiGRU

RESULTS

Table 1. Verification of the necessity of data preprocessing and standardization

| Data Processing | Standardization | Acc |
|---|-----------------|---------------|
| Raw Data (No Processing) | N | 0.6104 |
| | Y | 0.6427 |
| Noise Reduction | N | 0.6526 |
| | Y | 0.7866 |
| Noise Reduction + Endpoint Detection and Separation | N | 0.6948 |
| | Y | 0.8313 |

Table 2. Experimental results of different models

| Model | Feature | Acc | F1_score | Precision | Recall |
|----------|---------|---------------|----------|-----------|--------|
| C-LSTM | IS09 | 0.8313 | 0.8384 | 0.8249 | 0.8524 |
| | IS10 | 0.8511 | 0.8497 | 0.8682 | 0.8319 |
| C-GRU | IS09 | 0.8412 | 0.8498 | 0.8388 | 0.8619 |
| | IS10 | 0.8437 | 0.8431 | 0.8532 | 0.8333 |
| C-BiLSTM | IS09 | 0.8462 | 0.8537 | 0.8570 | 0.8504 |
| | IS10 | 0.8660 | 0.8672 | 0.8758 | 0.8590 |
| C-BiGRU | IS09 | 0.8437 | 0.8558 | 0.8498 | 0.8619 |
| | IS10 | 0.8561 | 0.8547 | 0.8814 | 0.8295 |

CONCLUSIONS

- We classify the **speech data** of teenagers, judge their **mental health status**, assist intelligent medical treatment.
- In this task, **bidirectional** contextual models and **LSTM** based models perform slightly better.
- It is a desirable way to make use of **contextual features** among utterances while using the **acoustic features** of each utterance.
- **Future**: **multi-class** classification (depression, anxiety, etc.); **multimodal** data fusion (with video, text, etc.).

REFERENCES

- [1] Poria, Soujanya, et al. "Context-dependent sentiment analysis in user-generated videos." *Proceedings of the 55th annual meeting of the association for computational linguistics (volume 1: Long papers)*. 2017.
- [2] Asokan, Ashish Ramayee, et al. "Interpretability for multimodal emotion recognition using concept activation vectors." *2022 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2022.