# TENSORFLOW SPEECH RECOGNITION CHALLENGE
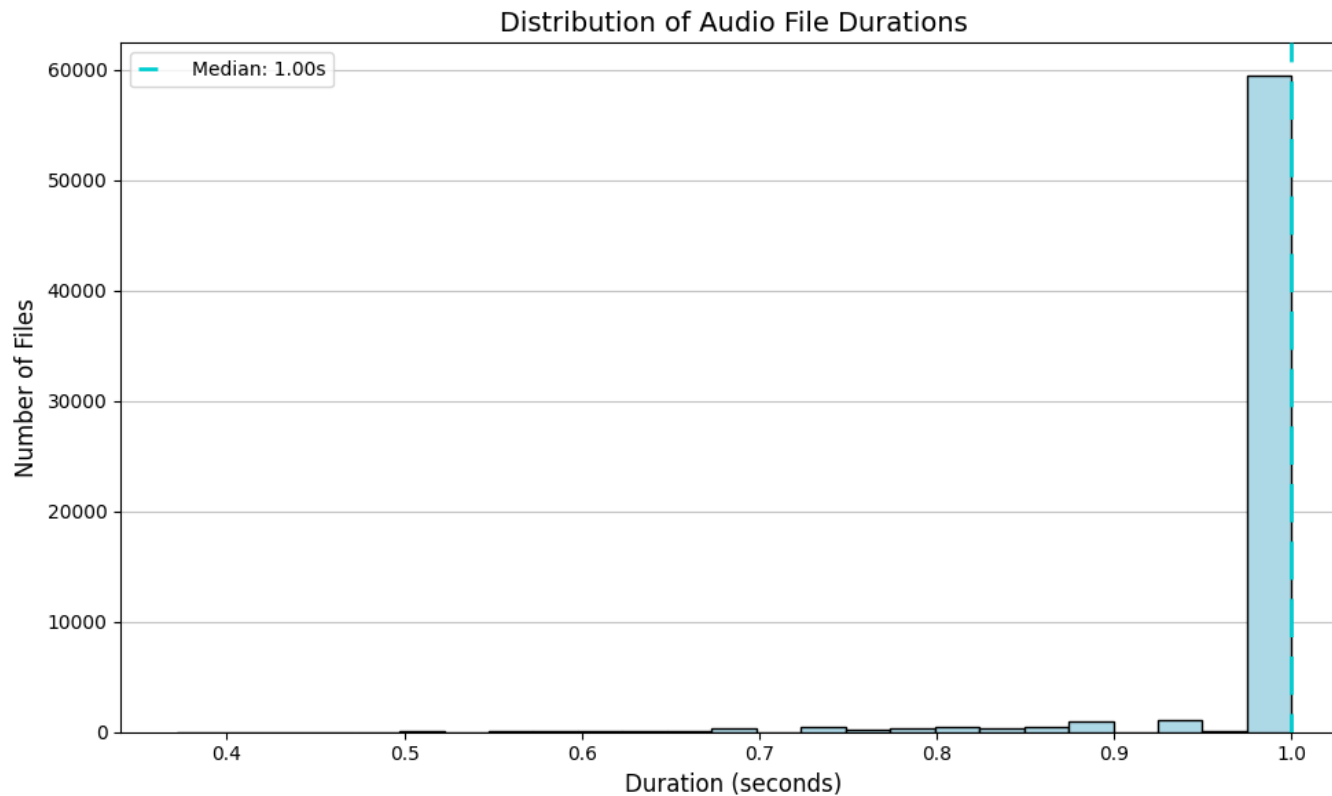
Paulina Kulczyk,

Jan Pogłód

| Class | Number of Files |
|---|---|
| _background_noise_ | 6 |
| bed | 1713 |
| bird | 1731 |
| cat | 1733 |
| dog | 1746 |
| down | 2359 |
| eight | 2352 |
| five | 2357 |
| four | 2372 |
| go | 2372 |
| happy | 1742 |
| house | 1750 |
| left | 2353 |
| marvin | 1746 |
| nine | 2364 |
| no | 2375 |
| off | 2357 |
| on | 2367 |
| one | 2370 |
| right | 2367 |
| seven | 2377 |
| sheila | 1734 |
| six | 2369 |
| stop | 2380 |
| three | 2356 |
| tree | 1733 |
| two | 2373 |
| up | 2375 |
| wow | 1745 |
| yes | 2377 |
| zero | 2376 |

# DATA OVERVIEW

- **Dataset**: TensorFlow Speech Recognition Challenge form Kaggle

- **Twenty core command words** were recorded, with most speakers saying each of them five times: "Yes", "No", "Up", "Down", "Left", "Right", "On", "Off", "Stop", "Go", "Zero", "One", "Two", "Three", "Four", "Five", "Six", "Seven", "Eight", and "Nine". **There are also ten auxiliary words**, which most speakers only said once: "Bed", "Bird", "Cat", "Dog", "Happy", "House", "Marvin", "Sheila", "Tree", and "Wow". And 6 files with noise (silence).

- **Task**: classify words to one of the following classes: **"Down", "Go", "Left", "No", "Off", "On", "Right", "Stop", "Up", "Yes", "Unknown"** (with other words from Dataset) and Silence (where none word was recognized)

# DURATION OF FILES
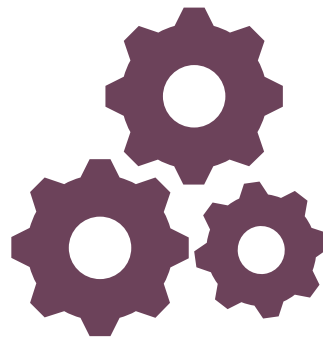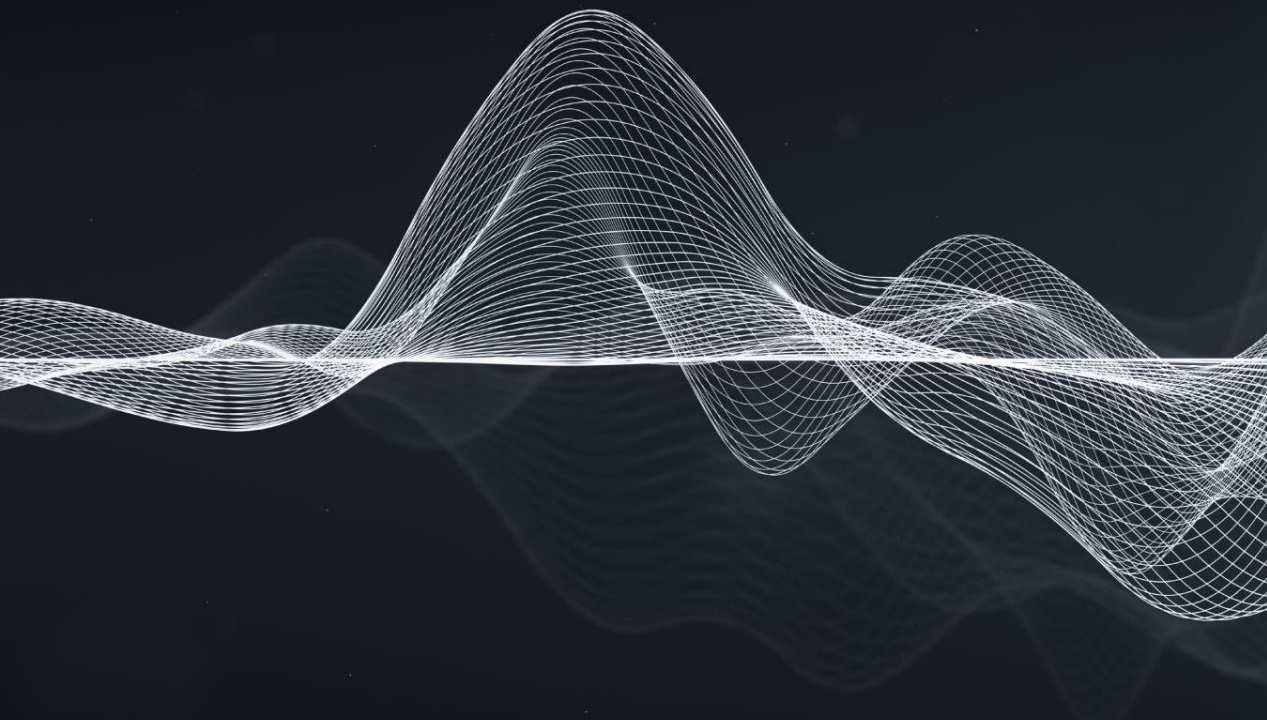


Distribution of Audio File Durations

- **Number of files:** 64721

- **Min:** 0.37s**, Max:** 1.00s

- **Mean duration:** 0.98 sekundy

- **Sampling rate**: 16 000 Hz

# PREPROCESSING

- Deletion of 308 files that last shorter then 0.6s

- Cutting 6 silence files to many shorter files which lasts about 1s (as our files from dataset)

- Mapping data to unknown and silence (randomly choosing 3000 files from uknown class with given seed)

- Bootstraping files per class and sampling from unknown and silence to have well balanced classes (3000 audio tracks per class)

```
📊 Number of files per class:
down: 3000 files
go: 3000 files
left: 3000 files
no: 3000 files
off: 3000 files
on: 3000 files
right: 3000 files
silence: 3000 files
stop: 3000 files
unknown: 3000 files
up: 3000 files
yes: 3000 files
```

# SCIENTIFIC BACKGROUND

- Voice can be displayed as an **acoustic wave**

- Main measurements of waves that has impact to sound are: **amplitude and frequency**

# SCIENTIFIC BACKGROUND - MEASUREMENTS

- **Mean RMS Amplitude** – is the square root of the average of the squares of a series of measurements of values of sign of wave (almost amplitude).

- **Zero-crossing rate (ZCR)** - normally it measures how often the signal changes sign — meaning it crosses the horizontal axis (from positive to negative or negative to positive) — within a frame. But in our study to better capture fast waveform fluctuations, we modified the ZCR calculation.
Instead of counting strict signal sign changes, we measure how often the waveform touches the zero axis, averaged over 0.25-second frames (as most files last 1s).
This approach approximates the traditional ZCR while improving sensitivity to rapid oscillations.

# SCIENTIFIC BACKGROUND – REFLECTED IN REAL LIFE EXAMPLES

- The **loudness of a sound** is directly proportional to the quadrature of the **amplitude of the sound wave**



I See Fire - Ed Sheeran
Mean RMS: 0.05803



Famous Ay-Oh! - Freddie Mercury
Mean RMS Amplitude: 0.12956

# SCIENTIFIC BACKGROUND - REFLECTED IN REAL LIFE EXAMPLES

- Indirectly, **how fast we speak** or **how we stretch syllables has an effect on the frequency of the acoustic wave**

# SCIENTIFIC BACKGROUND – OUR DATA

# BINARY CLASSIFICATION PROBLEM FOR RECOGNITION "YES" AND "NO"



Figure 7: Learning process of the basic CNN architecture for binary problem



Figure 9: Learning process of the simple transformer architecture for binary problem



Figure 8: Learning process of the LSTM architecture for binary problem



Figure 10: Learning process of the advanced architecture (CNN + transformer) for binary problem

| Architecture | Valid accuracy (mean) |
|---|---|
| CNN | 97.5% |
| LSTM | 98,1% |
| Transformer | 97.9% |
| CNN + Transformer | 98.7% |

# CNN COMPARISON OF PARAMETERS – BATCH SIZE

# TRANSFORMER COMPARISON OF PARAMETERS – BATCH SIZE
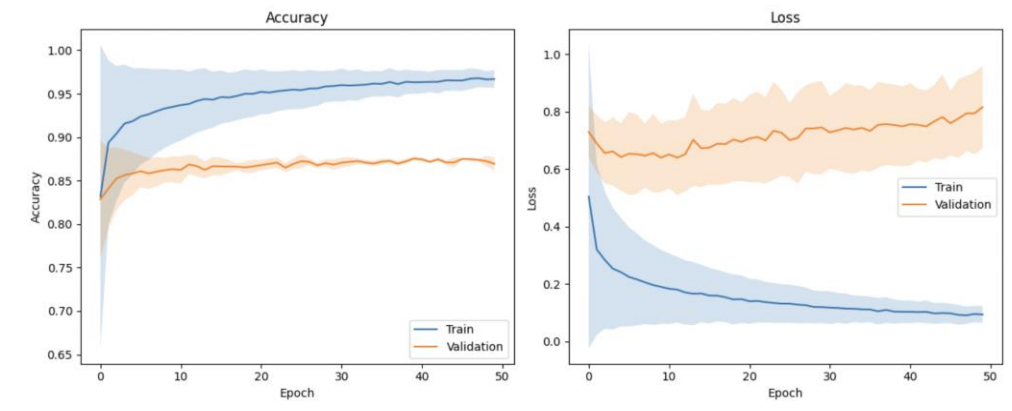
# TRANSFORMER COMPARISON OF PARAMETERS – LEARNING RATE

# TRANSFORMER COMPARISON – NUMBER OF ATTENTION HEADS

# CNN      VS    TRANSFORMER
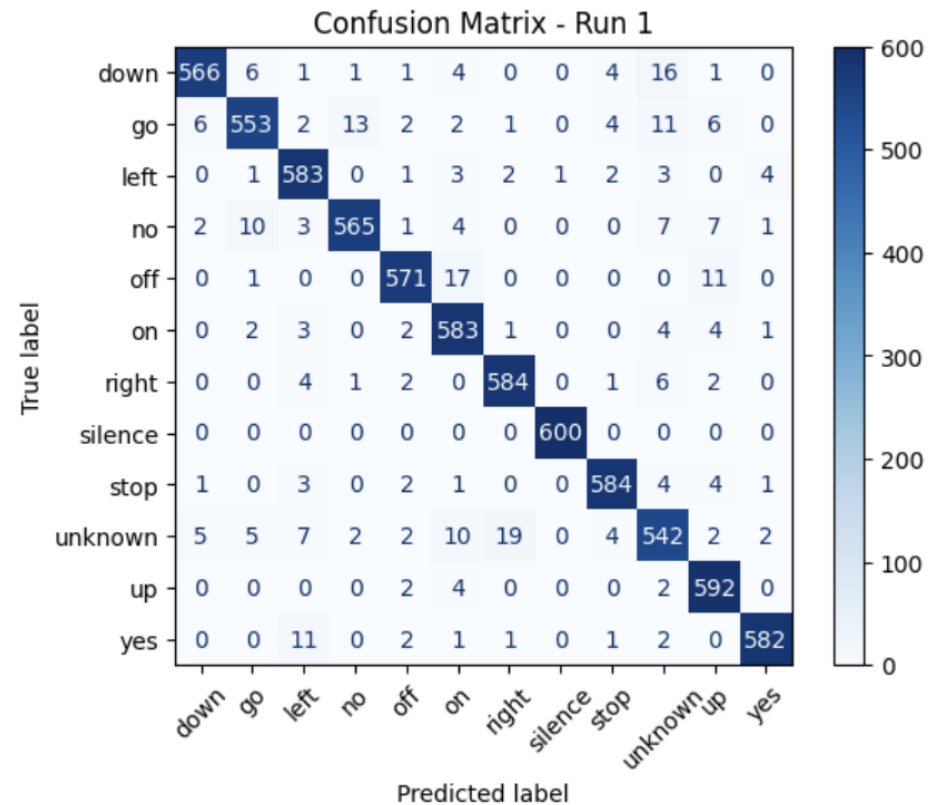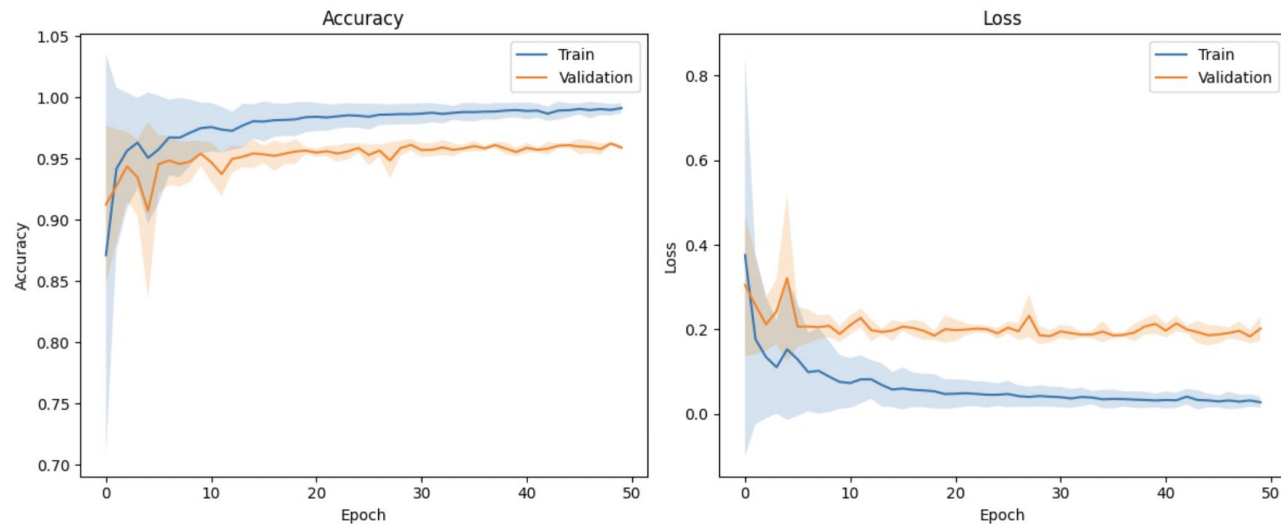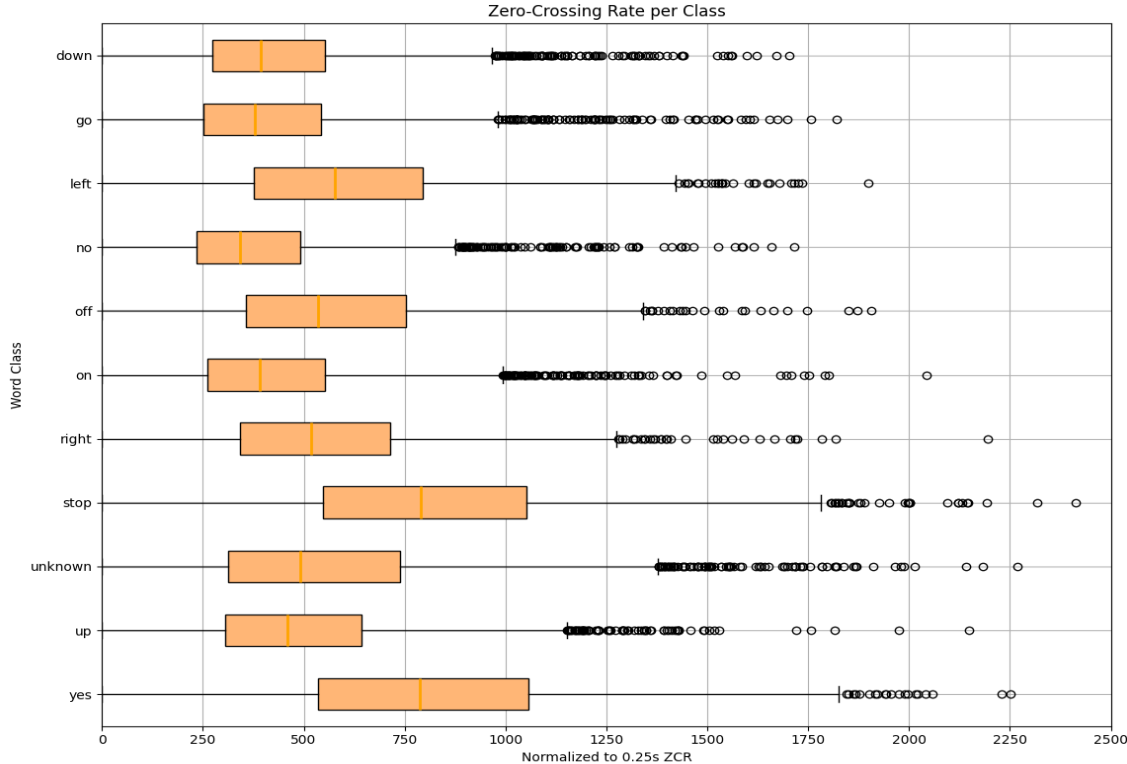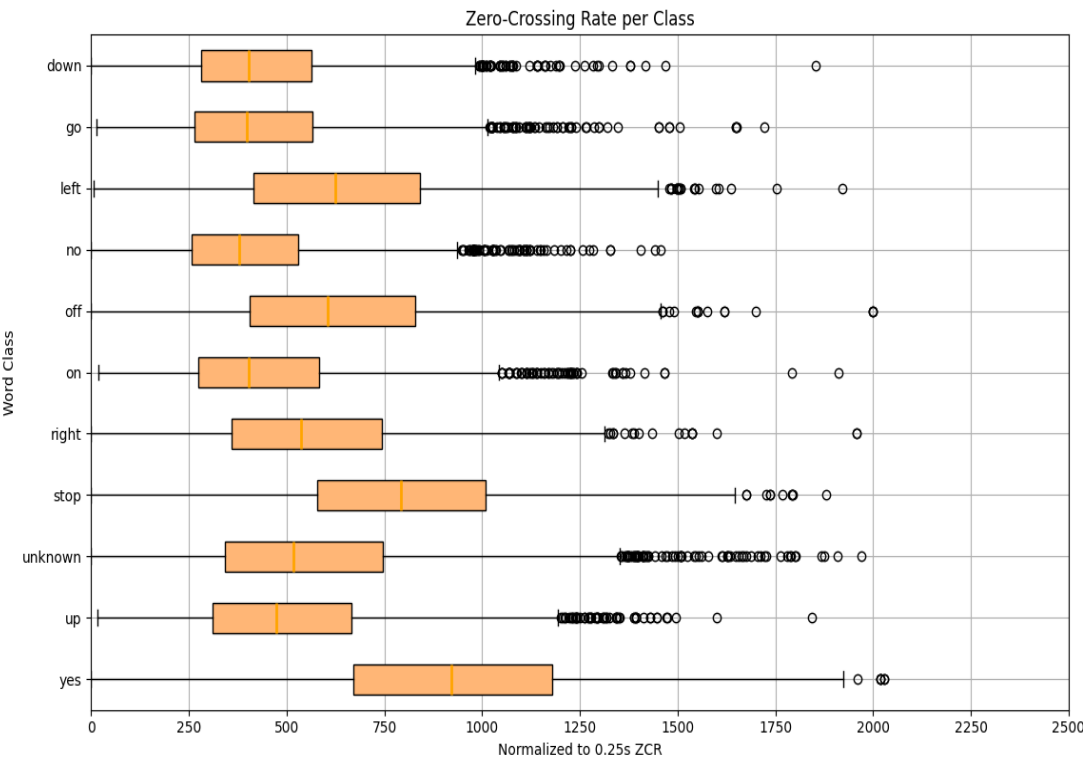
# TRANSFORMER + CNN ARCHITECTURE

# DATA AUGMENTATION -> BETTER GENERALIZATION

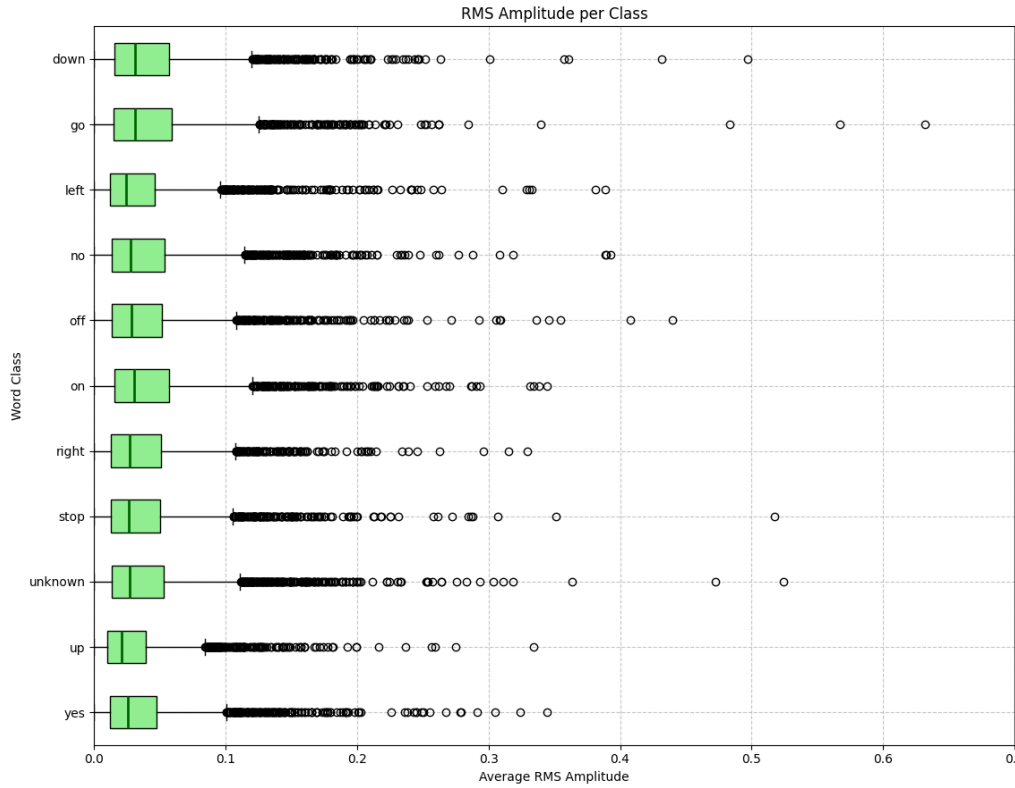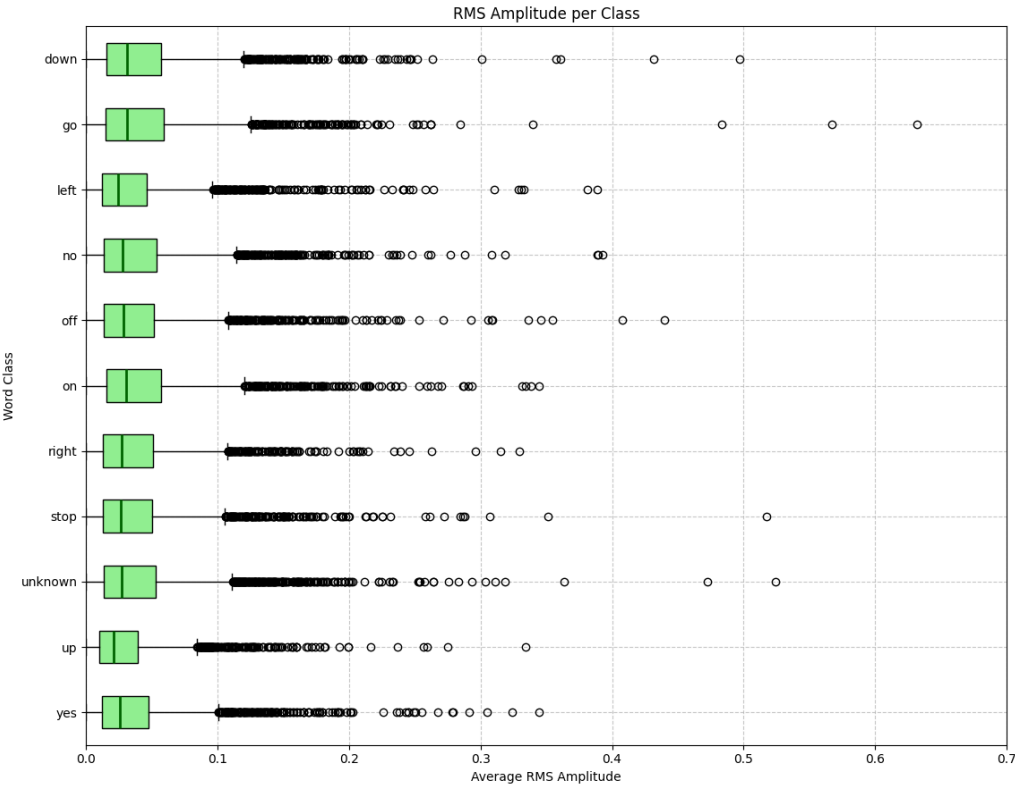- We have large dataset but… we can still enlarge it with also changing a little bit our data to reflects on more ways how people accent words

- People can say word slower or faster, can accent something louder or ending word really quietly, or in general speaking quietly, can also speak in a noisy space

- We want to reflect all these situation by applying the following augmentations: stretching speech, changing volume, adding noise
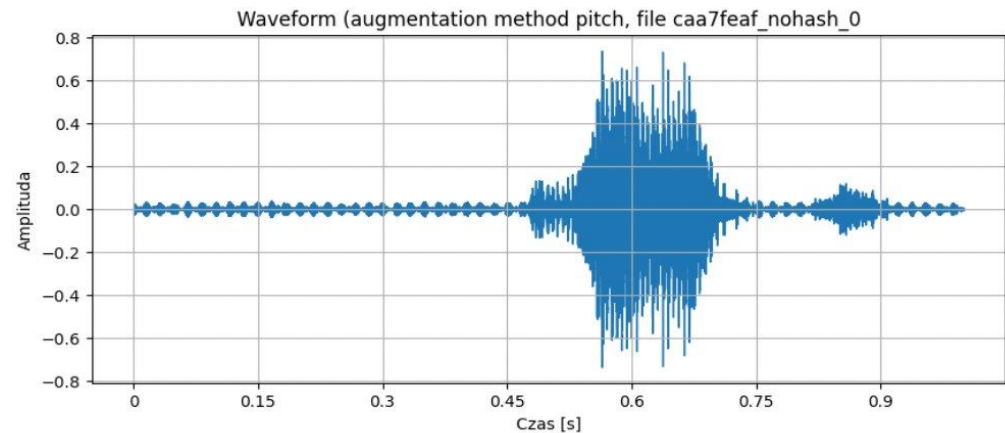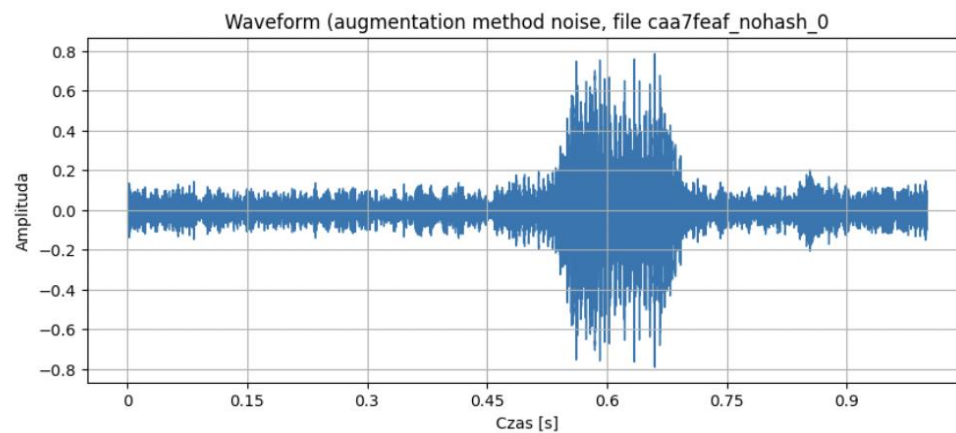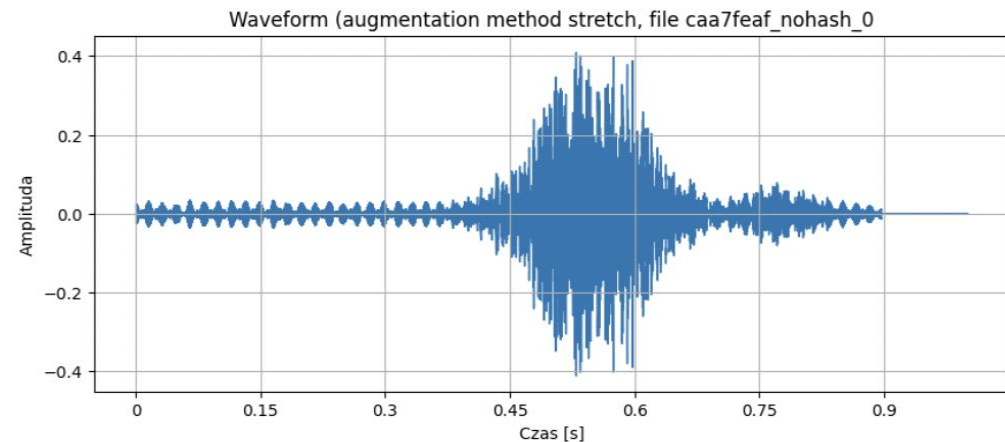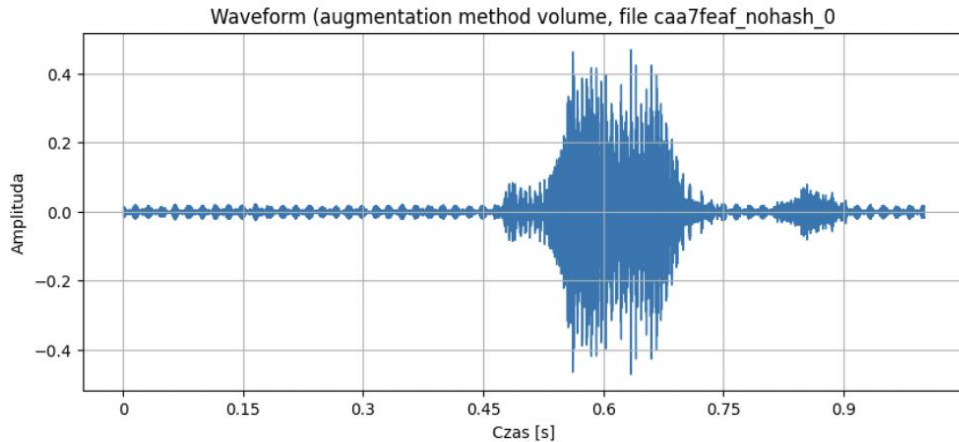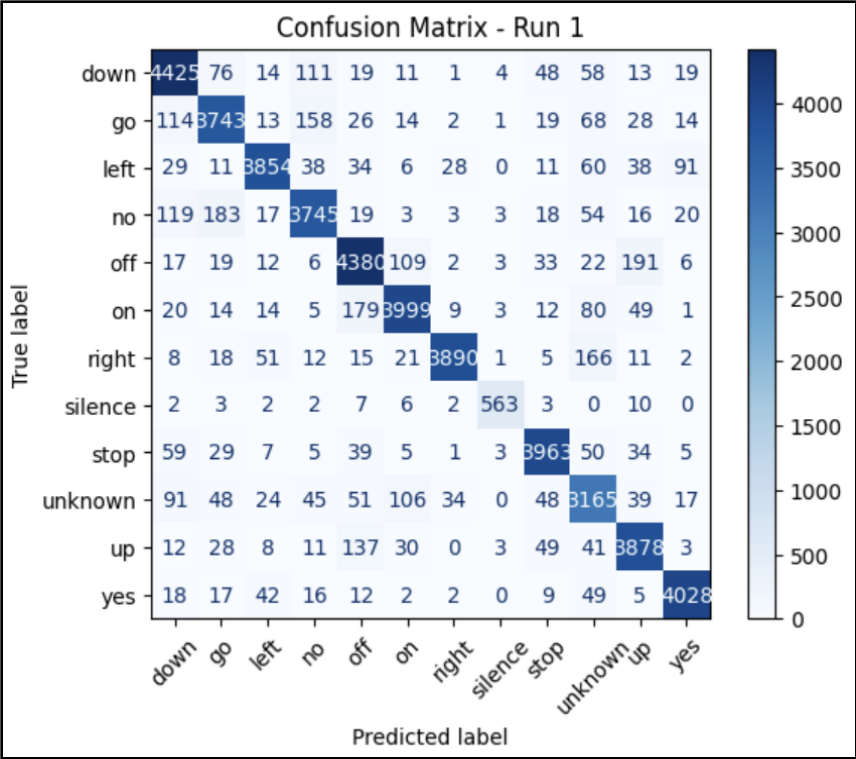
# STRETCHING

# VOLUME CHANGING

# DATA AUGMENTATION TECHNIQUES

# CNN FOR AUGMENTED DATA



| AUGMENTED DATA | ORIGINAL DATA |
|:---:|:---:|
| 91,57% | 89,65% |

# TRANSFORMER+CNN FOR AUGMENTED DATA



Accuracy / Loss / Confusion Matrix - Run 3

| AUGMENTED DATA | ORIGINAL DATA |
|---|---|
| 96,17% | 95,90% |

# FUTURE GOALS

- Explore more augmentation techniques
- Analyse different ensemble models (focus on classes like on/off)
- Explore different architectures of Transformer+CNN
- Test more hyperparameters
- Extend the range of possible labels

# THANK YOU FOR YOUR ATTENTION!