# Project II - Transformers
## Deep Learning 2024

Kinga Frańczak, 313335
Grzegorz Zakrzewski, 313555

# Contents

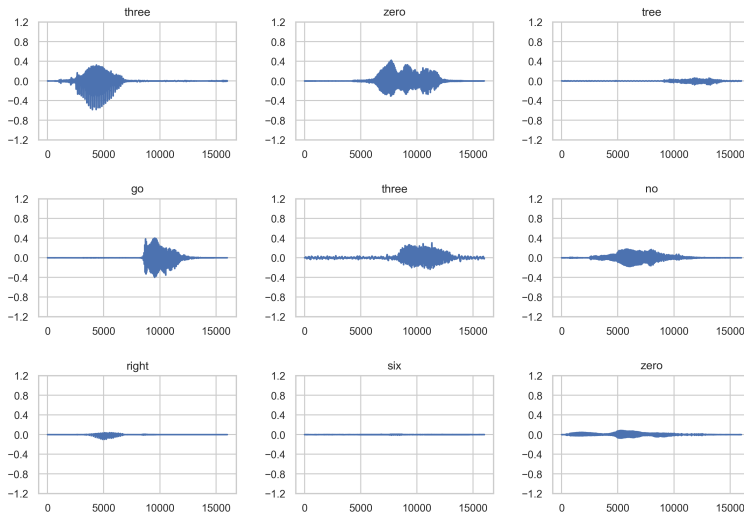# Description of the research problem



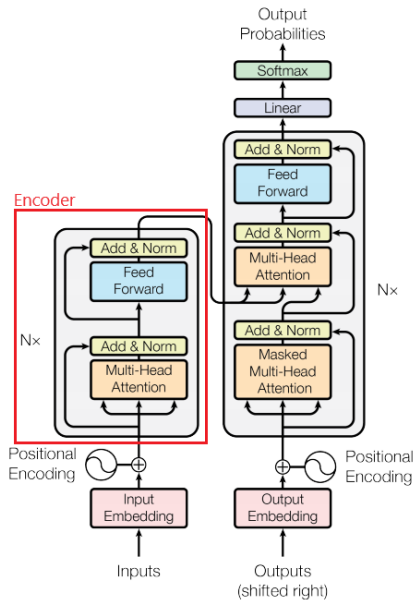Figure: Sample audio clips from Speech Commands dataset.

# Transformer - encoder



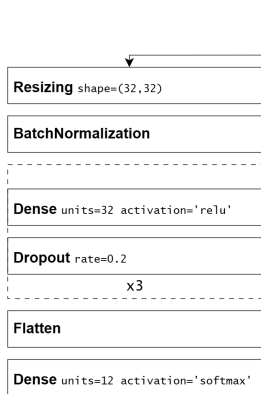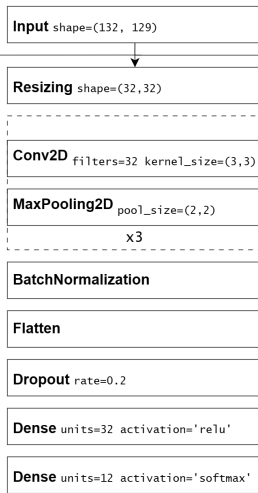Figure: The original "Attention is All You Need" Transformer diagram.

# Experiments - network architectures



Figure: Three neural network architectures used in the experiments.

# Experiments - details

| Experiment | Objective | Values |
|---|---|---|
| 1 | Architecture | Simple feed-forward |
| | | CNN |
| | | Transformer |
| 2.1 | Number of attention heads | 2 |
| | | 4 |
| | | 6 |
| | | 8 |
| | | 10 |
| 2.2 | Number of Encoder sub-layers | 2 |
| | | 4 |
| | | 6 |
| | | 8 |
| 3 | Handling *silence* and *unknown* classes | one network for all classes |
| | | separate network for special cases |

Table: Details of the experiments.

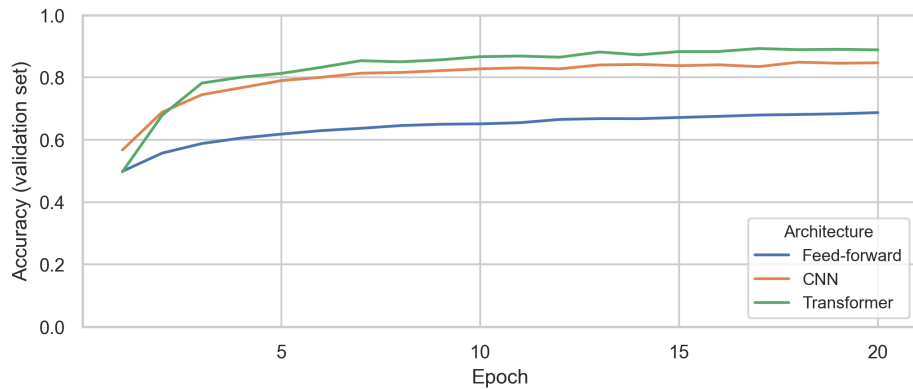# Experiment 1 - architecture



Figure: Accuracy computed on the validation subset for every epoch.

# Experiments 2.1 & 2.2 - hyper-parameters

| Number of attention heads | Accuracy | Validation accuracy |
|---|---|---|
| 2 | 0.930 (0.004) | 0.901 (0.003) |
| 4 | 0.914 (0.008) | 0.888 (0.012) |
| 6 | 0.866 (0.036) | 0.845 (0.025) |
| 8 | 0.794 (0.021) | 0.791 (0.016) |
| 10 | 0.803 (0.023) | 0.798 (0.012) |

(a) Experiment 2.1 - number of attention heads.

| Number of Encoder sub-layers | Accuracy | Validation accuracy |
|---|---|---|
| 2 | 0.941 (0.003) | 0.875 (0.004) |
| 4 | 0.930 (0.004) | 0.897 (0.006) |
| 6 | 0.905 (0.015) | 0.884 (0.009) |
| 8 | 0.468 (0.132) | 0.424 (0.174) |

(b) Experiment 2.2 - number of Encoder sub-layers.

Table: The mean (and standard deviation) of the best values of the accuracy achieved by models.
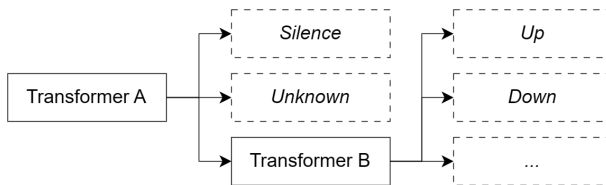
# Experiment 3 - handling *silence* and *unknown* classes.



Figure: Experiment 3 - diagram

|               | Accuracy      | Validation accuracy |
|---------------|---------------|---------------------|
| Transformer A | 0.947 (0.015) | 0.934 (0.011)       |
| Transformer B | 0.940 (0.003) | 0.913 (0.005)       |
| Combination   | 0.944 (0.014) | 0.917 (0.011)       |

Table: Experiment 3 - results.

# Confusion matrix


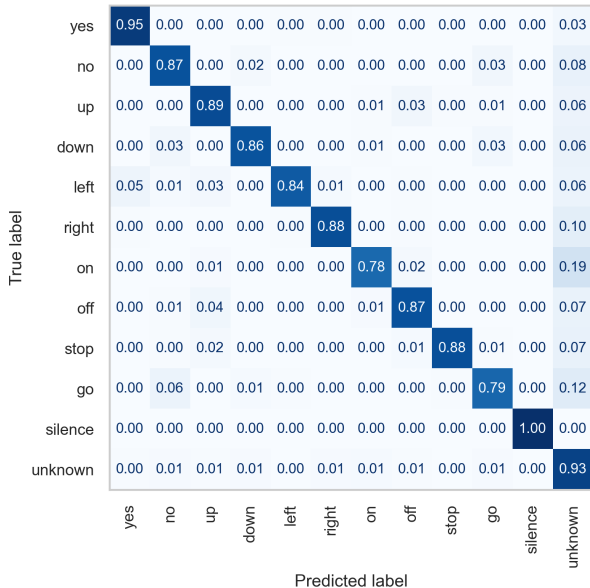
Figure: Confusion matrix prepared on the validation dataset, normalized by true conditions (rows).

# Conclusions

- The project was prepared according to the instructions.
- The best parameter settings:
    - number of attention heads $= 2$;
    - number of Encoder sub-layers $= 4$;
    - separate network for *silence* and *unknown* classes;
- The accuracy achieved on the validation subset (around 0.91) seems very high.
- The accuracy achieved on Kaggle is 0.67.