MSc in Artificial Intelligence
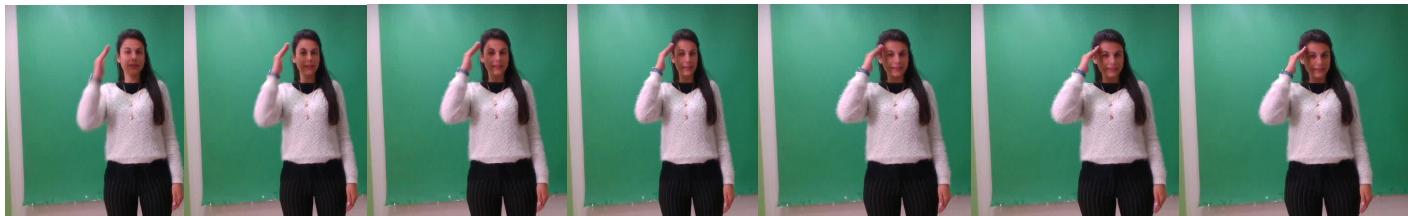
DEMOKRITOS

# "Multimodal machine learning techniques for sign language recognition"

Zidianakis Matthaios - MTN2008

Papadopoulos Georgios - MTN2025

# Overview

- Sign language recognition refers to:
  - Classification problem of a sequence of video frames
  - Subcategory of gesture recognition in motion



- Dataset: *GSL* (Greek Sign Language)
  - Non multilabel (Isolated words)
  - 7 signers
  - 5 repetitions of each signers
  - 5 scenarios
  - 310 individuals words
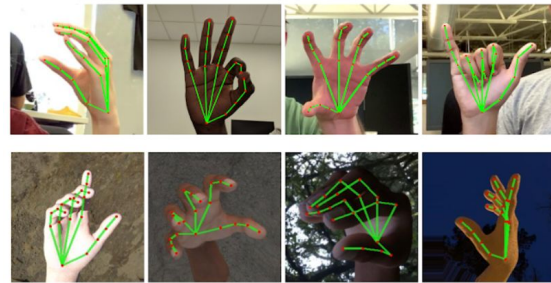  - 30 FPS - RGB (+Depth) - 848×480

*Not used here*

# Tools



- Keypoints extraction   MediaPipe

- Deep learning methods   PyTorch

- Classical Machine learning (SVM)   scikit learn

# Dataset distribution

- Predefined *train*, *validation*, *test* sets → **not mixed** signers

- 4 flavours of the same dataset:

  1. Keeping the first **10** unique words.

  2. Keeping the first **12** unique words which have more than 500 samples (undersampled **exactly to 500**).

  3. Keeping the first **10** unique words containing **more than 500** samples for each words.
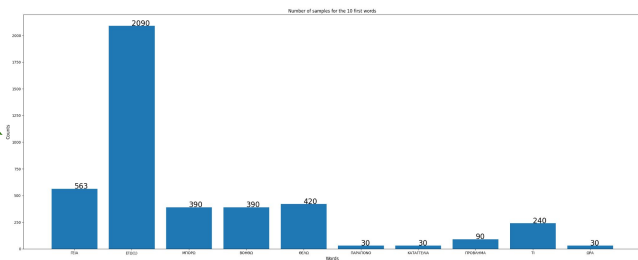
  4. Keeping all the glosses (**310**).

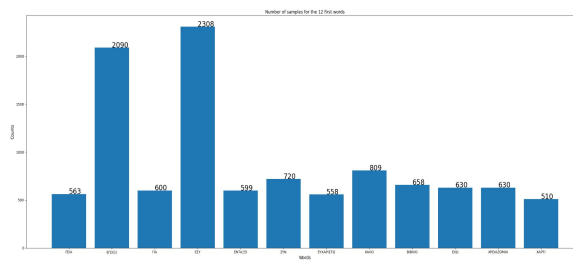

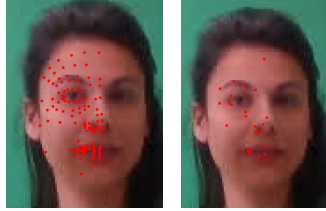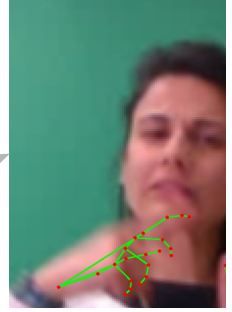Figure 1: Word distribution plot for dataset 1 of train set.

Figure 2: Word distribution plot for dataset 2 of train set before undersampling.

# Data preprocessing



Sampling every 5 face points

Blurring, 400% rescale and interpolation



Magnitude: $M_t^N = \sqrt{\left(D_{N_{x(t+1)}} - D_{N_{x(t)}}\right)^2 + \left(D_{N_{y(t+1)}} - D_{N_{y(t)}}\right)^2}$
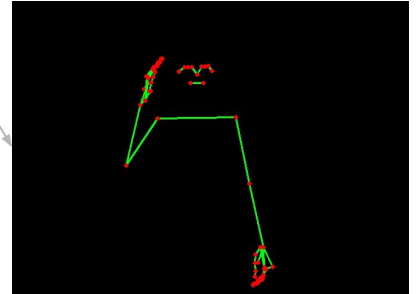
Preprocessing involves points selection:

Angles: $A_t^N = tanh^{-1}\left(\dfrac{D_{N_{y(t+1)}} - D_{N_{x(t)}}}{D_{N_{x(t+1)}} - D_{N_{x(t)}}}\right)$
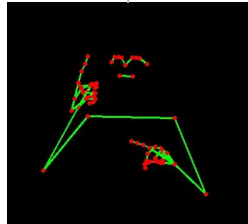


Upper body Pose points

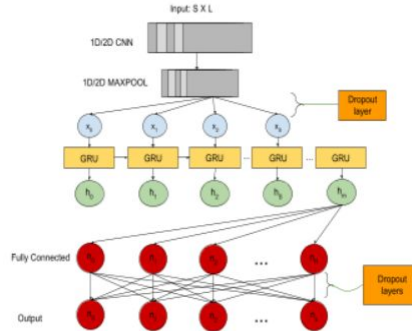Pose skeleton in blank image

# Methods: Deep learning models

# Results – 1D CNN-GRU *dataset 1* (sampling face points)



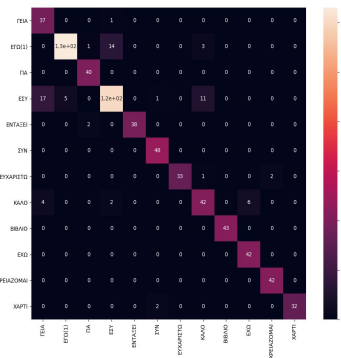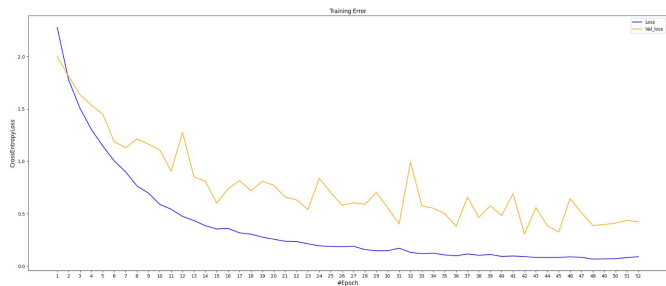| F1-macro score | | |
|---|---|---|
| | Validation set | Test set |
| With class weights | 53% | 59% |
| Without class weights | 61% | 67% |

Totally lose 2 classes!

$$weight_i = 1/\left(\ |samples_i|\ /\ \sum_{i=1}^{n}\ |samples_i|\ \right)$$

$$normalized\ weight_i = weight_i / \sum_{i=1}^{n} weight_i$$

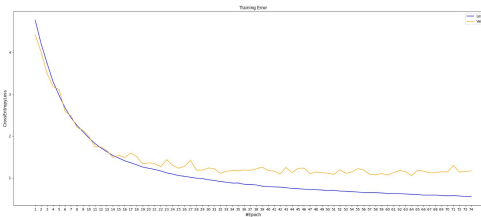# Results – 2D/1D CNN-GRU (sampling face points)

## 1D CNN-GRU *dataset 2*



| F1-macro score | |
|---|---|
| Validation set | Test set |
| 91.77% | 93.89% |

## 2D CNN-GRU *dataset 2 (skeletal)*

| F1-macro score | |
|---|---|
| Validation set | Test set |
| 42.1% | 43.1% |

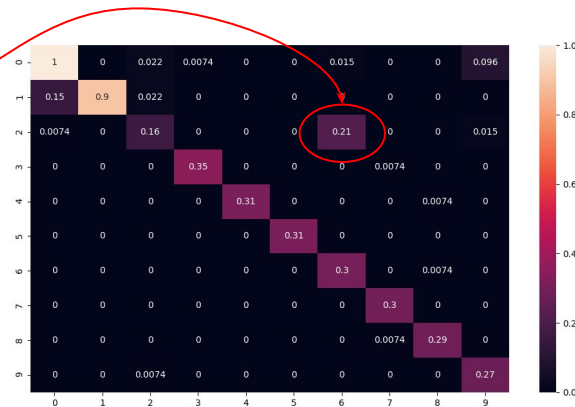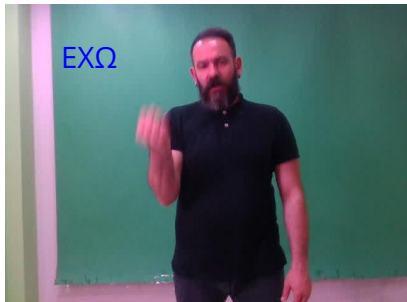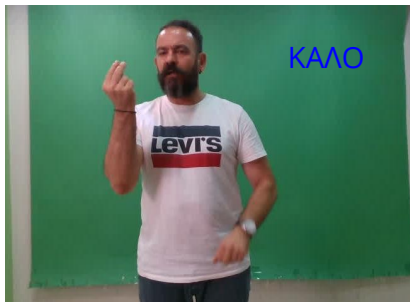## 1D CNN-GRU *dataset 4 (all words)*



| F1-macro score | |
|---|---|
| Validation set | Test set |
| 56.81% | 60.5% |

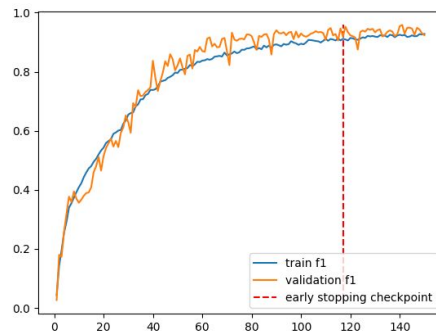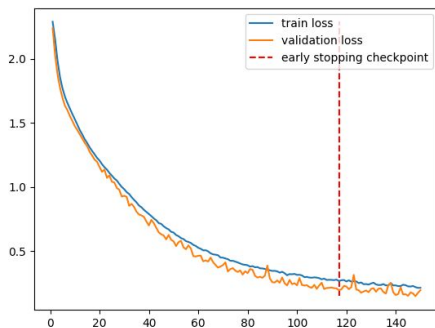# Results - Early experiments

- *LSTM* architecture
  - full head keypoints → noisy features
  - misclassified two similar words

# Results - LSTM



*Training/ Validation*
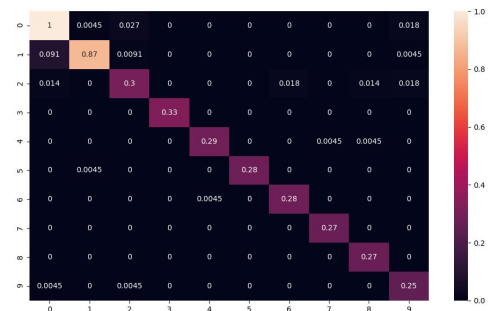
94%

- *F1 macro score*
- *L: upper body points*

*Test*

95%

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| ΕΣΥ | ΕΓΩ (1) | ΚΑΛΟ | ΣΥΝ | ΒΙΒΛΙΟ | ΧΡΕΙΑΖΟΜΑΙ | ΕΧΩ | ΓΙΑ | ΕΝΤΑΞΕΙ | ΓΕΙΑ |

# Results – Miscellaneous

| SVM F1-macro scores (sampling face points) | | |
|---|---|---|
| | Validation set | Test set |
| Dataset 1 | 78% | 88% |
| Dataset 2 | 81% | 82% |
| Dataset 4 | 39% | 42% |

| Angles - Magnitude F1-macro score (*dataset 2*) | | |
|---|---|---|
| | Validation set | Test set |
| 1D CNN-GRU | 70% | 65% |
| SVM | 45% | 49% |

| Filters applied -  F1-macro score (*dataset 2*) | | |
|---|---|---|
| | Validation set | Test set |
| 1D CNN-GRU | 88.25% | 91.86 |

False points recognition!

# Conclusions



- *Mediapipe* → suitable for keypoints feature extraction
    - Especially for different deep learning methods
    - Adequate classical ML (SVM) performance

- Improve with **only** the representative points

- Applied filters were weaker

- Preprocessing proved to be significant!

- Future work:
    - Keypoints position relative to image → independent to captured image angle
        - Division by shoulder points coordinates
        - Video frames data augmentation
    - Examine *GANs* architecture → denoise video frames from motion blur

# THANK YOU