UNIVERSITÀ
DEGLI STUDI
DI PADOVA

DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

Simone Milani
Room 216
DEI A

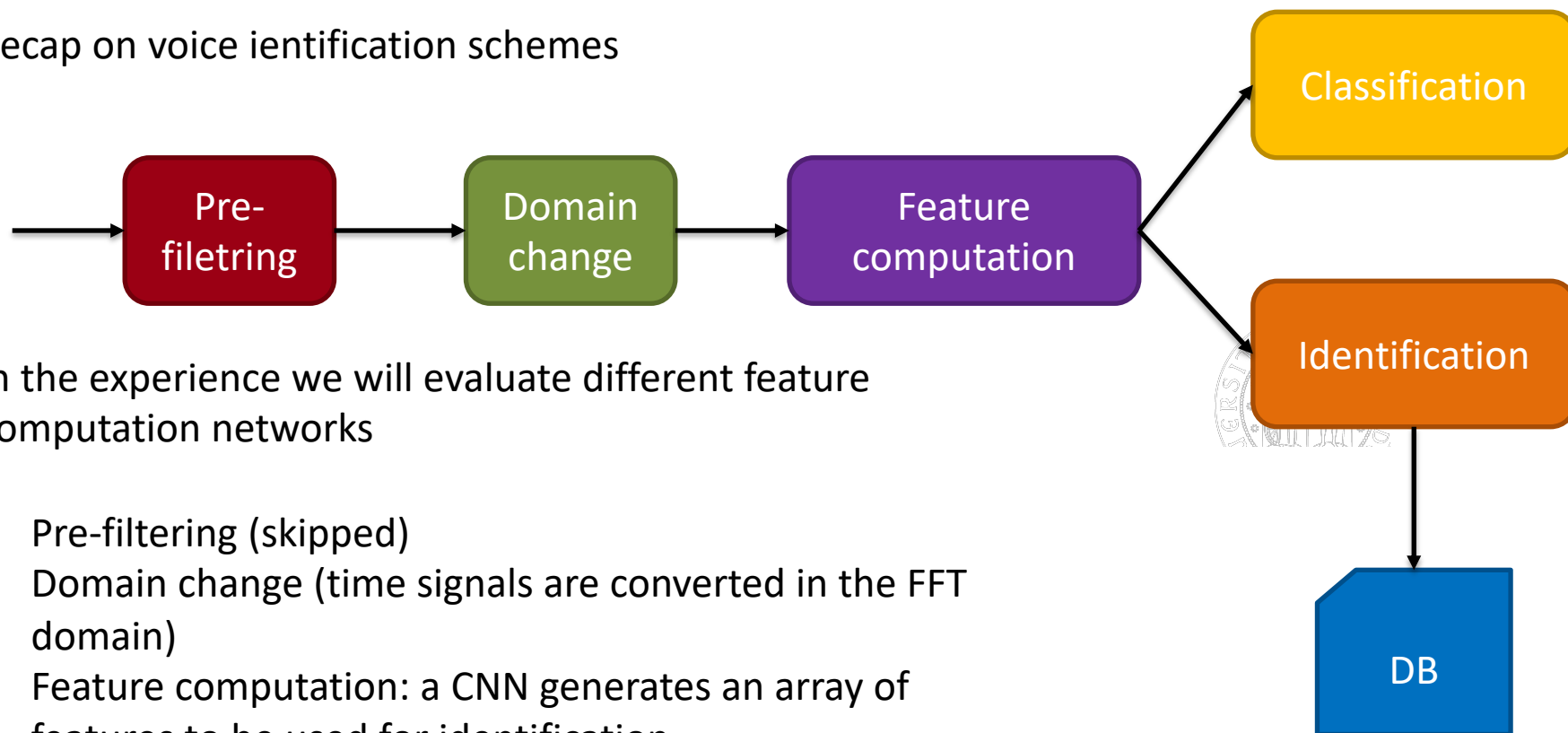Phone: 049 827 7641
E-mail:
simone.milani@dei.unipd.it

# BIOMETRICS

A.A. 2020/2021

## LECTURE 23

## LAB SESSION 2 – VOICE RECOGNITION

Recap on voice ientification schemes



In the experience we will evaluate different feature computation networks

- Pre-filtering (skipped)
- Domain change (time signals are converted in the FFT domain)
- Feature computation: a CNN generates an array of features to be used for identification
- Classification: identifies the speaker in a pool of 5 persons

Dataset for each speaker is divided into three sets:

- Training set: used to optimize parameters
- Validation set: used to monitor the training process
- Test set: used to verify the final accuracy

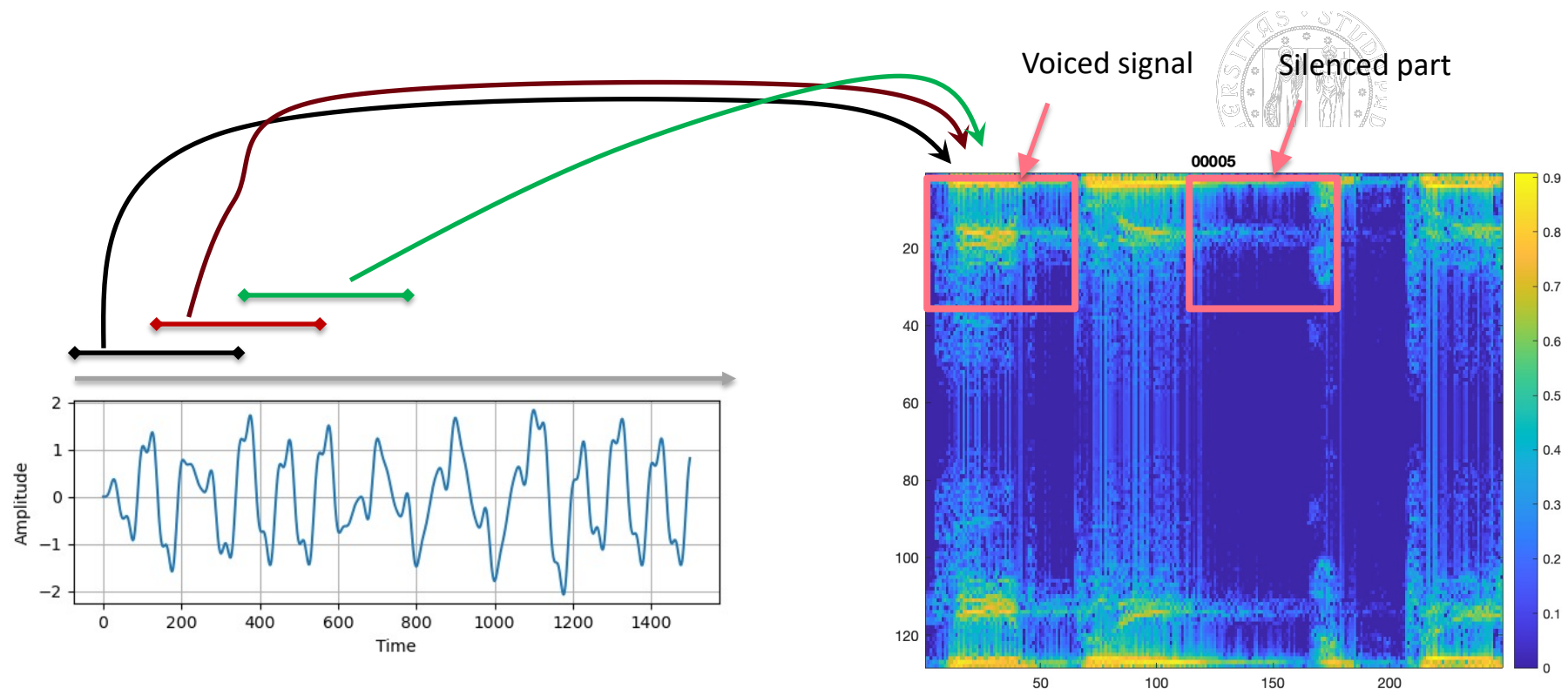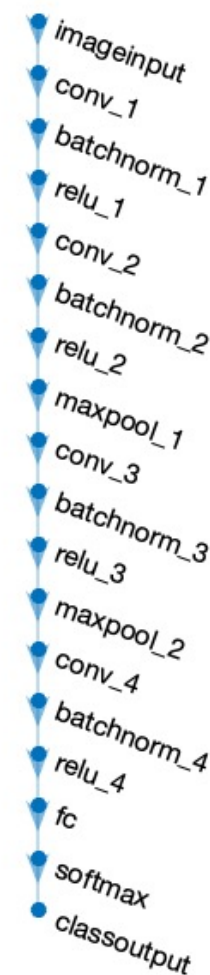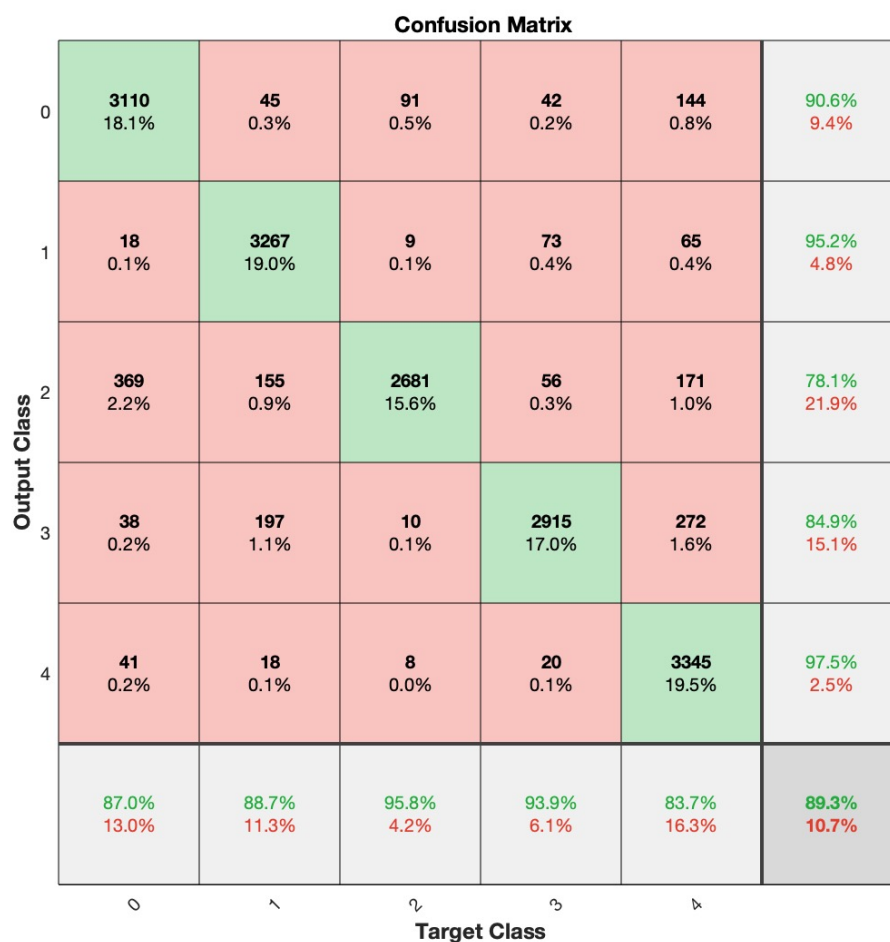Each set includes different files.

With different cardinalities.

- Spectrogram is created by windowing the input audio signal.
- Overlapping window are seletced
- FFT is computed on each window
- Results are included in the columns of the spectrogram image
- 32x32 relevant windows are selected (with voiced signal)



Voiced signal          Silenced part

- Features and classification is performed using a CNN
- Progressive set of convolutional, batch normalization, activation layers
- Final layer is a fully connected
- Final accuracy visualized with confusion matrix



Confusion Matrix

Use Cepstral values instead of spectrum values.

What is the difference in spectrograms?

How does the performance change?

Optimize the network in order to maximize the performance.

Which is the top accuracy you can get?