# *voxplorer*

# –

# The Voice Explorer

## Alessandro De Luca
## LiRI – CL UZH

Universität Zürich UZH

# Schedule

1. Motivation
2. What is *voxplorer* and where to find it
3. How to install *voxplorer*
4. DEMO 1: Visualizer mode (Saarbrücken data)
5. DEMO 2: Speaker embeddings mode (CANDOR)
6. Feature extraction mode: parameter settings
7. DEMO 3: Feature extraction mode (CANDOR)
8. DEMO 4: Filtering mode
9. Future plans, goals, and Q&A

# Motivation

- With the recent and fast improvements in computational power and storage capacity, driven by the need of more and better descriptive data that covers a large study population, speech corpora have started to become increasingly large and complex.

- E.g. CANDOR:
  - Multi-modal data: archived (.zip) = 1TB
  - Only full-length WAV files $\cong$ 440 GB
  - Processed (VAD + filtering and snipping) 28480 files [942 speakers]

# Motivation

**Better access to computational resources for all!**

# What is *voxplorer*?

- *voxplorer* started as an interactive visualization tool, BUT…

- An interactive tool for exploration and subsetting of large corpora

- A powerful automatic feature extraction tool

Where to find *voxplorer*:

https://github.com/delale/voxplorer

# Installation demo…

- Clone the <u>repository</u>
- Install <u>miniconda</u> or <u>Anaconda</u>
- Create the *voxplorer* conda environment
- Activate the environment
- Have fun!

```
// clone the repository
 ~ [ base]
→  git clone https://github.com/delale/voxplorer.git


 ~ [ base]
→  cd voxplorer


// create conda environment
 ~ [ base]
→  conda env create -f voxplorer_env_[platform].yml


// activate conda environment
 ~ [ base]
→  conda activate voxplorer


// run voxplorer
 ~ [ voxplorer]
→  python voxplorer.py
```

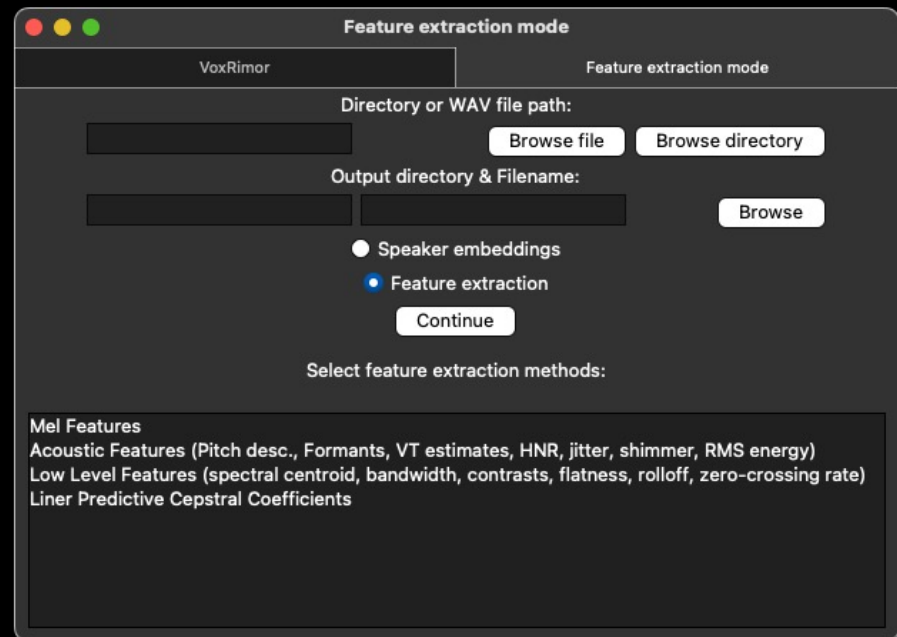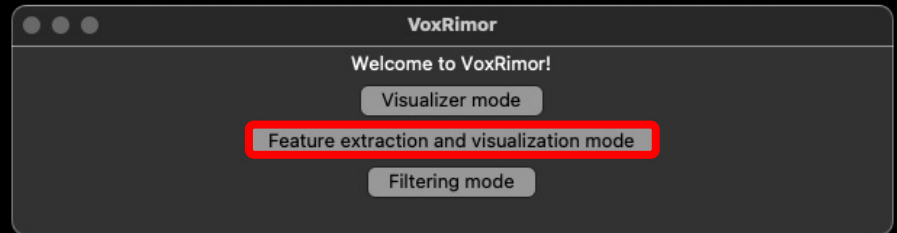# DEMO 1: Visualizer mode

# DEMO 2: Speaker embeddings mode

# Feature extraction mode

- Mel-features

- Acoustic features

- Low-level features

- Linear predictive cepstral coefficients

# Feature extraction mode: Mel-features

- MFCCs: coefficients representing the short-term power spectrum of an audio signal computed using Mel-frequency scaling and often used in automatic speech processing tasks.

- Delta features:
  - delta: first derivative of the MFCCs
  - delta-delta: second derivative of the MFCCs

# Feature extraction mode: Mel-features

- n_mfcc: number of coefficients
- n_mels: number of Mel bands
- win_length: length of analysis frame (ms)
- overlap: length of overlap between successive frames (ms)
- fmin: lowest frequency (Hz)
- fmax: highest frequency (Hz)
- premphasis: pre-emphasis coefficient
- lifter: cepstral filtering coefficient
- deltas: compute also delta and delta-delta features?
- summarise: summarise ($\mu$ and $\sigma$) of each feature at the utterance (file) level

# Feature extraction mode: Acoustic features

- Pitch: mean, median, minimum, maximum, std. dev.
- Formants: F1, F2, F3, F4
- VT estimates: formant dispersion, avg. formant, formants' geometric mean, [Fitch VTL](), [VTL Δf]()
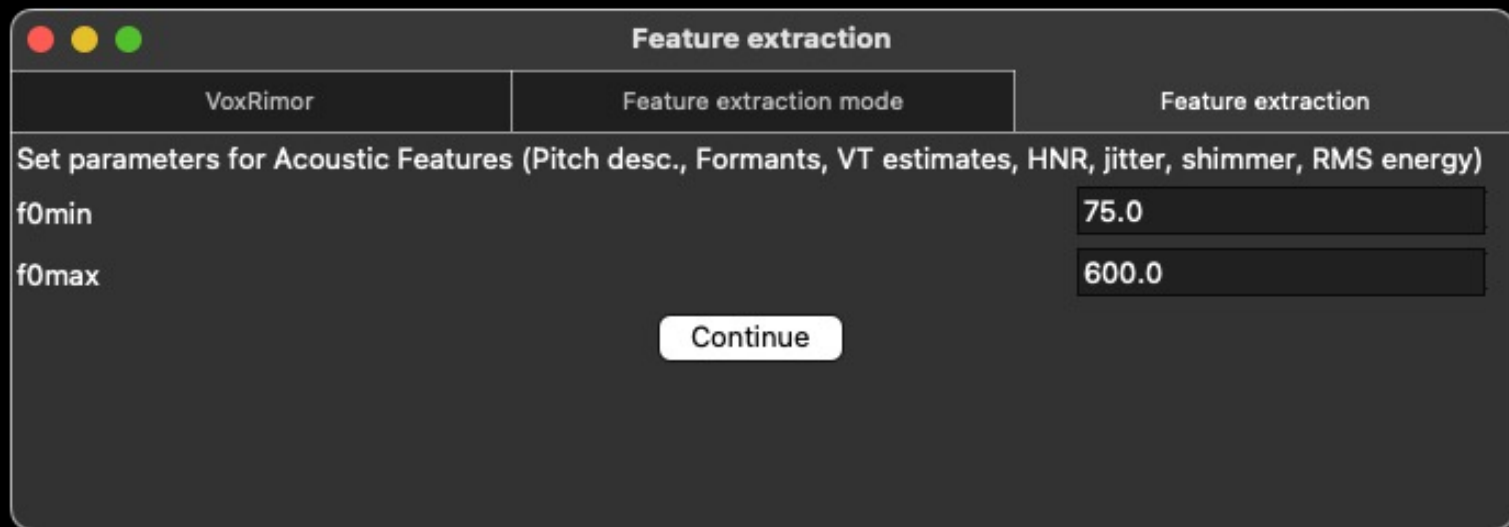- HNR
- Jitter & Shimmer
- RMS energy

# Feature extraction mode: Acoustic features

- f0min: pitch floor (Hz)

- f0max: pitch ceiling (Hz)

Uses Praat (Sound: To Pitch (cc)); Soon more parameters to come…

# Feature extraction mode: Low-level features

- spectral centroid: "centrer of mass" of avg. frequency in frame

- spectral bandwidth

- spectral contrast: mean energy contrast between peaks and valleys in each frequency sub-band

- spectral flatness: how "noise-like" vs. "tone-like" is the sound?

- spectral roll-off: roll-off frequency (freq. below which 85% of the energy is contained)

- zero-crossing rate

# Feature extraction mode: Low-level features

- win_length: length of analysis frame (ms)

- overlap: length of overlap between successive frames (ms)

- premphasis: pre-emphasis coefficient

- n_bands_contrasts: number of frequency bands by which to divide the spectrum to calculate contrasts (num. contrasts = n_bands+1)

- use_mean_contrasts: calculates the average contrast

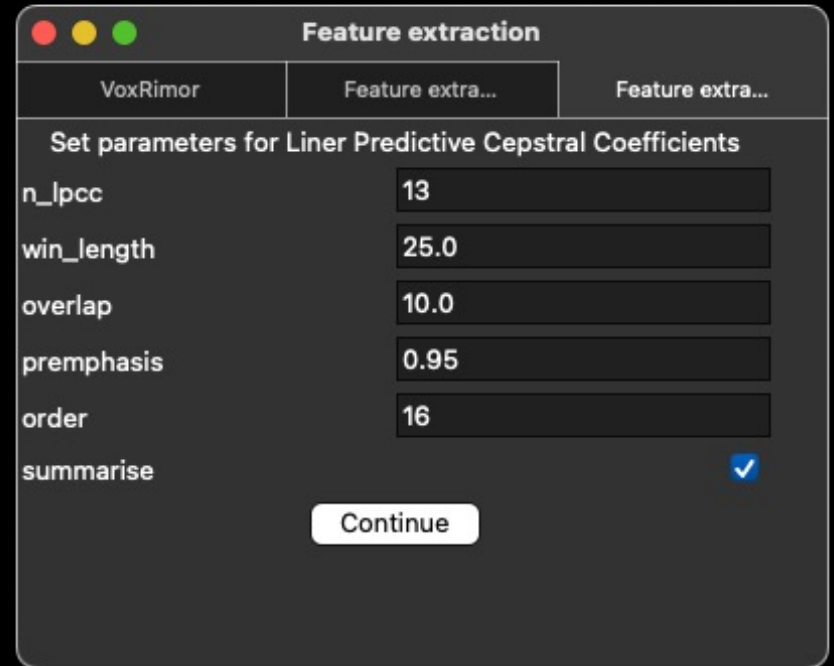- summarise: summarise ($\mu$ and $\sigma$) of each feature at the utterance (file) level

# Feature extraction mode: LPC features

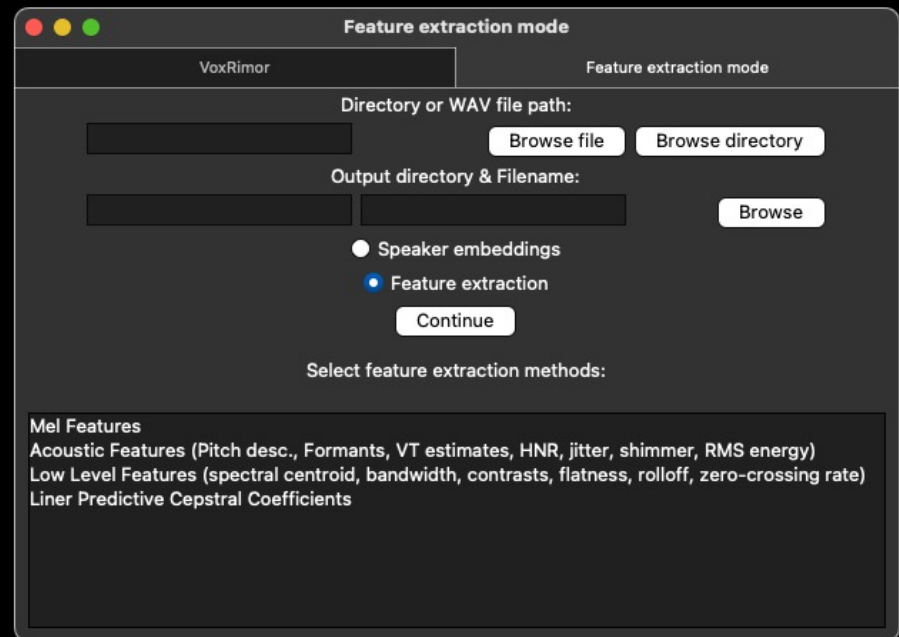- LPC: linear prediction coefficients (Burg); a representation of the spectral envelope after applying a linear filter

- LPCCs: linear predictive cepstral coefficients; coefficients representing the cepstrum (log-power spectrum) of the linear prediction coefficients

# Feature extraction mode: LPC features

- n_lpcc: number of coefficients

- win_length: length of analysis frame (ms)

- overlap: length of overlap between successive frames (ms)

- premphasis: pre-emphasis coefficient

- order: order of the linear filter

- summarise: summarise ($\mu$ and $\sigma$) of each feature at the utterance (file) level
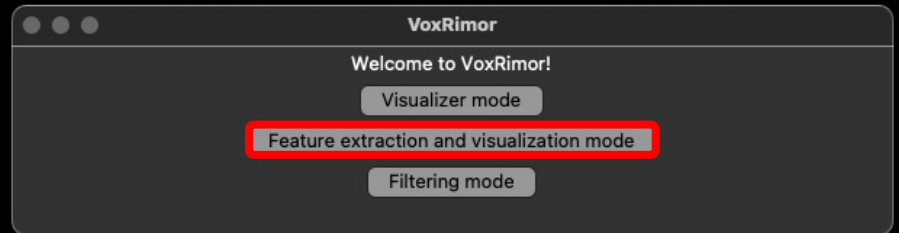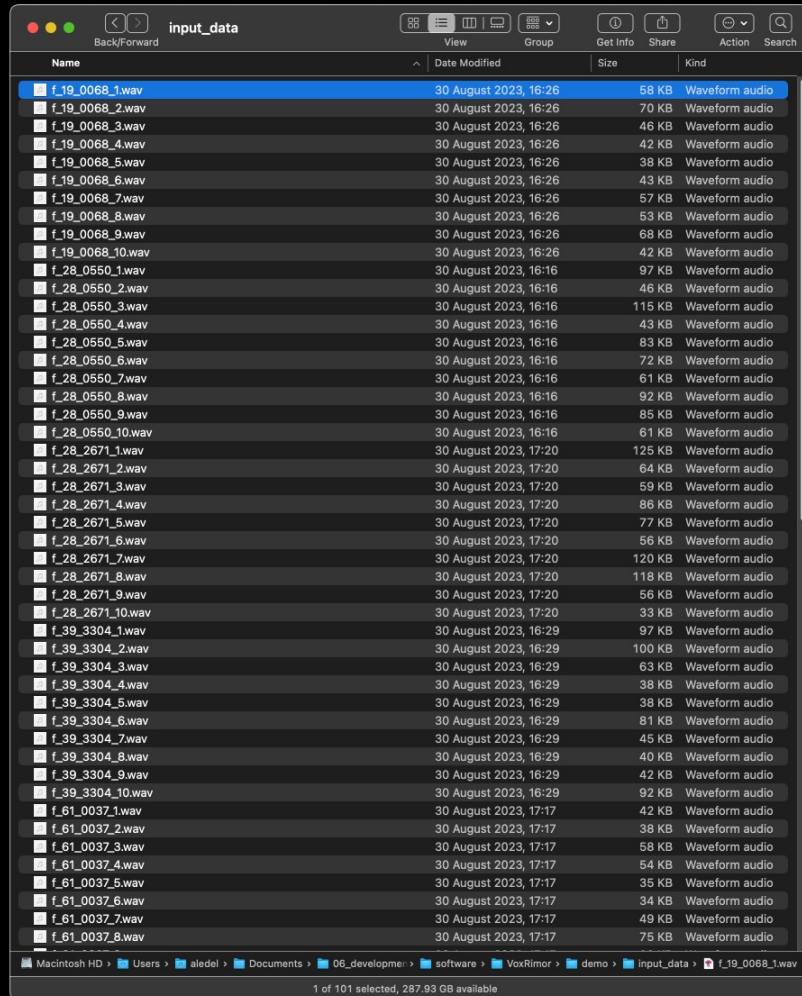
# DEMO 3: Feature extraction mode
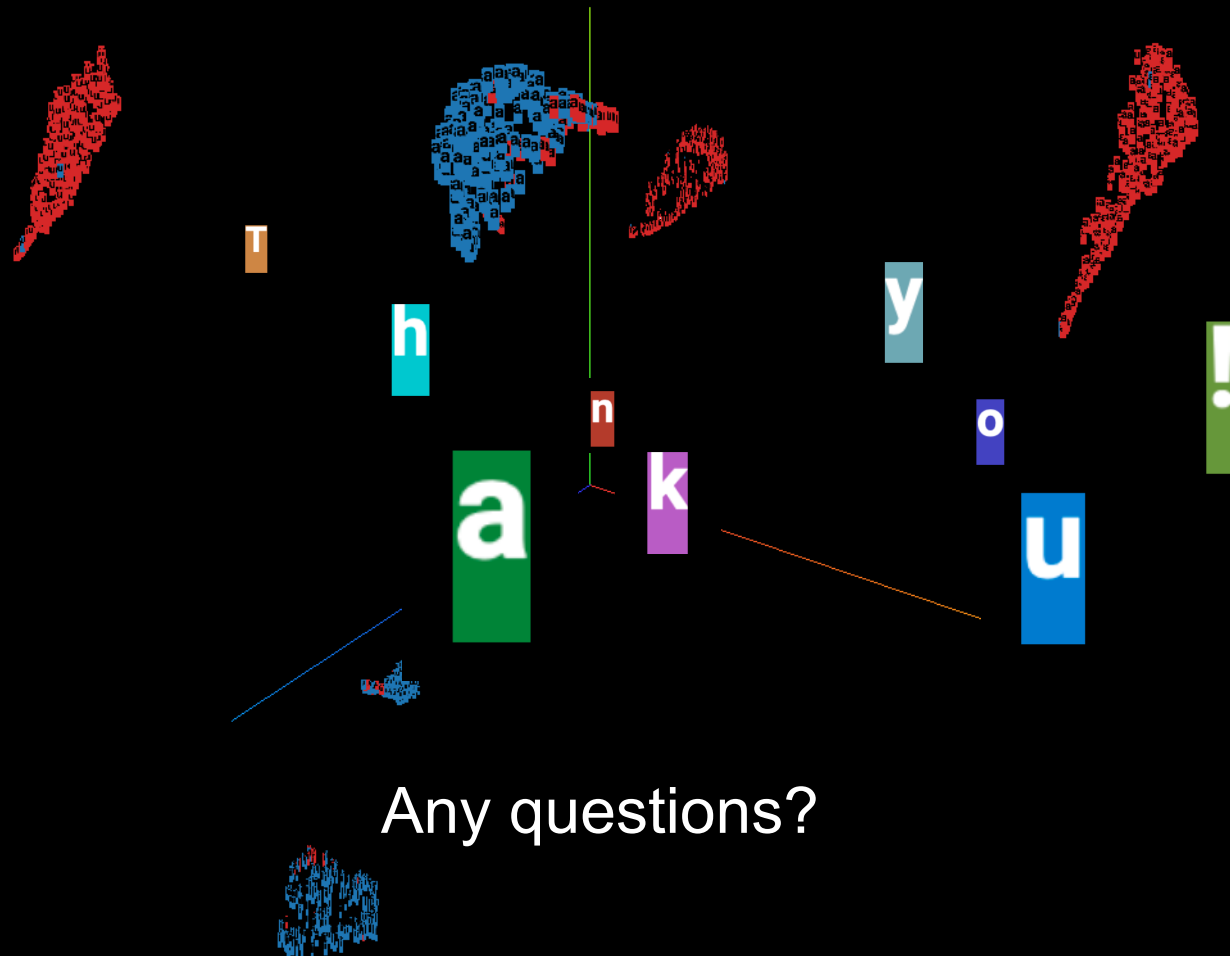
# Future plans, goals, and Q&A

- Praat TextGrid support.
- More parameter settings → better feature extraction.
- More features: e.g. VoiceSauce
- Integrated speaker verification system and classification tools (ECAPA-TDNN VoxCeleb2).

Fully-integrated web application and UI:
- *voxplorer* as an open-source online LiRI service.
- Direct selection & downloading of data.
- Downloadable reduced dimensions versions of data.
- Statistical tools: distances, box-plots, distributions.

*voxplorer* Thanks You!

Any questions?