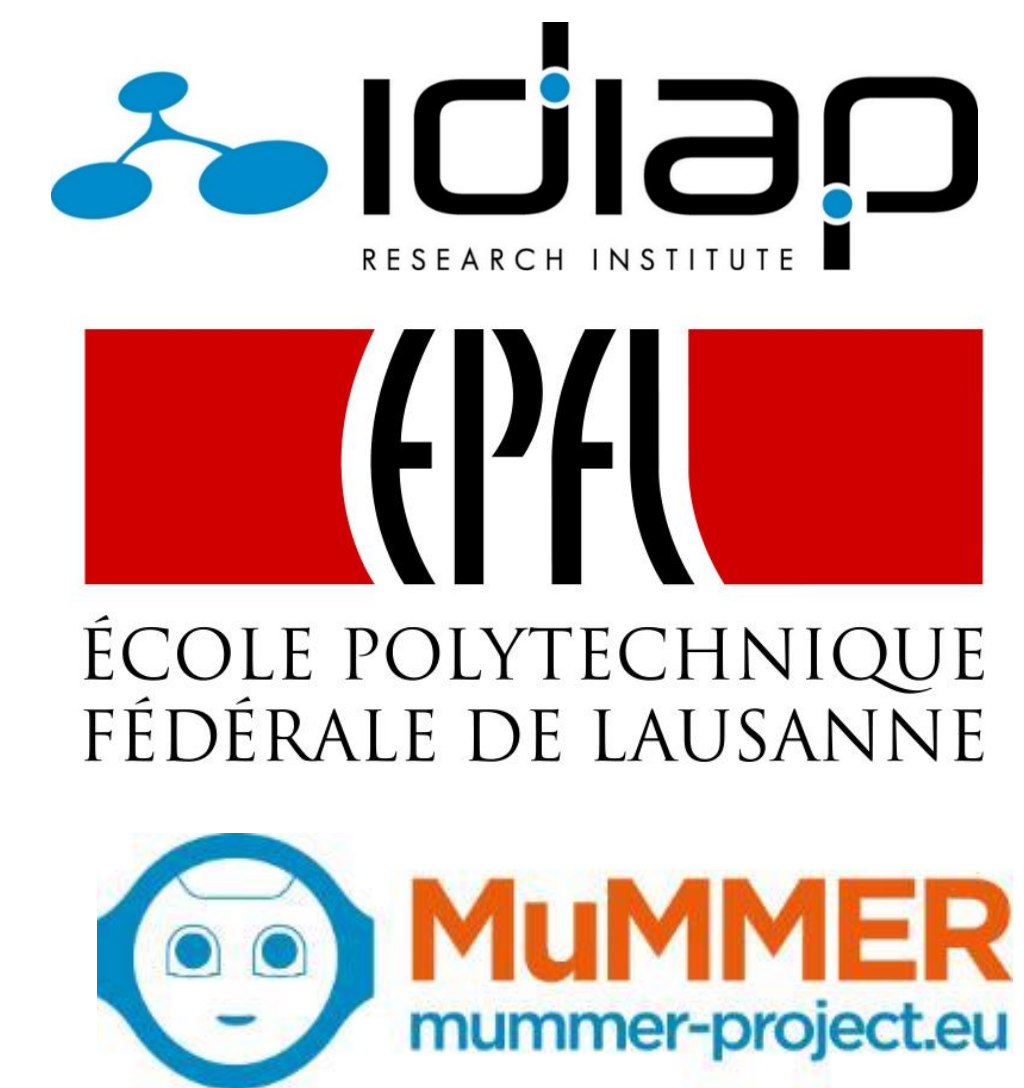# Deep Neural Networks for Multiple Speaker Detection and Localization

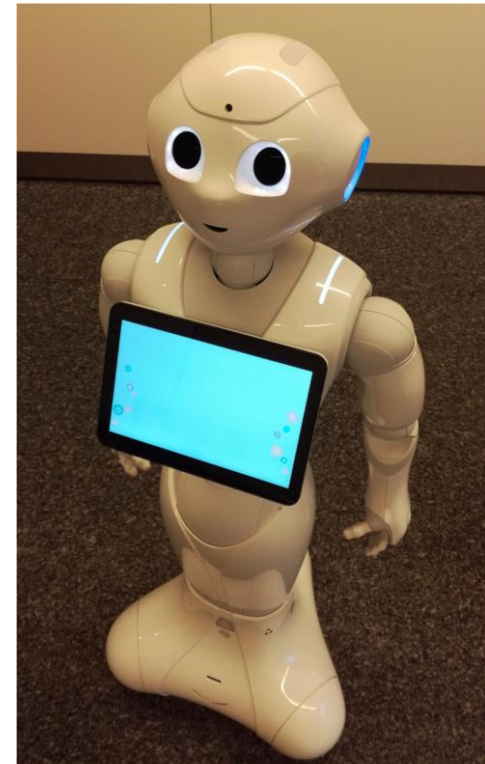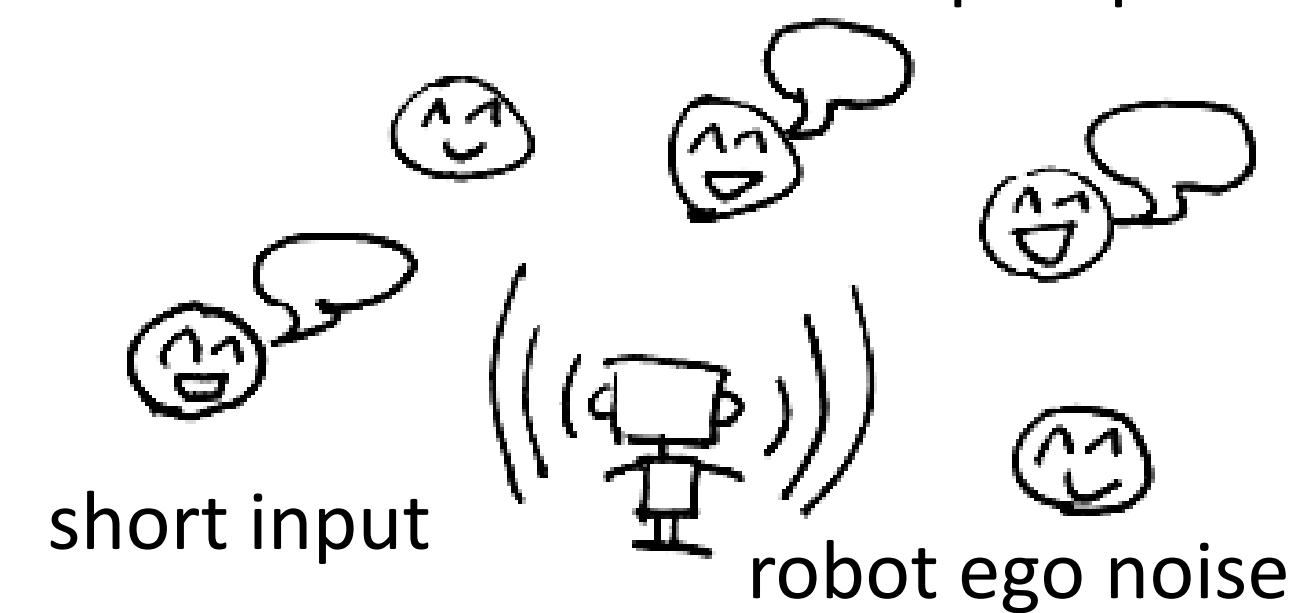## [1]Idiap Research Institute, Switzerland        [2]EPFL, Switzerland

Weipeng He[1,2], Petr Motlicek[1], Jean-Marc Odobez[1,2]

## Introduction

**Task:** Sound source localization in real HRI scenarios

unknown number of multiple speakers
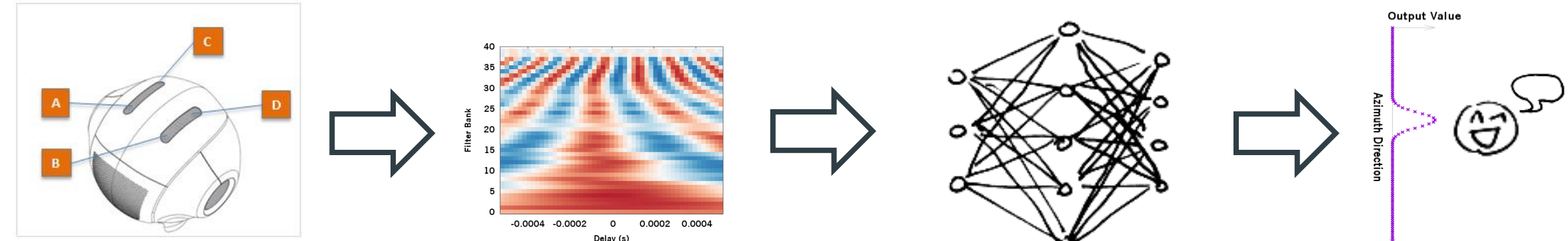
short input    robot ego noise

**Contributions:**
- Novel deep learning-based multiple sound source localization method.
- Likelihood-based output encoding handles an arbitrary number of sources.
- Investigation of three network architectures based on different motivations.
- Study sub-band cross correlation information as an input feature for better localization cues in speech mixtures.
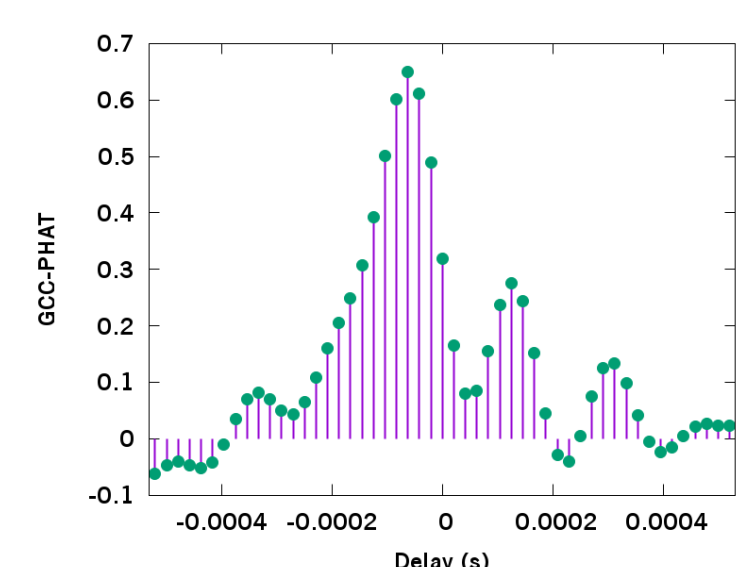- Collected and released a benchmark dataset of real recordings.

## Approach

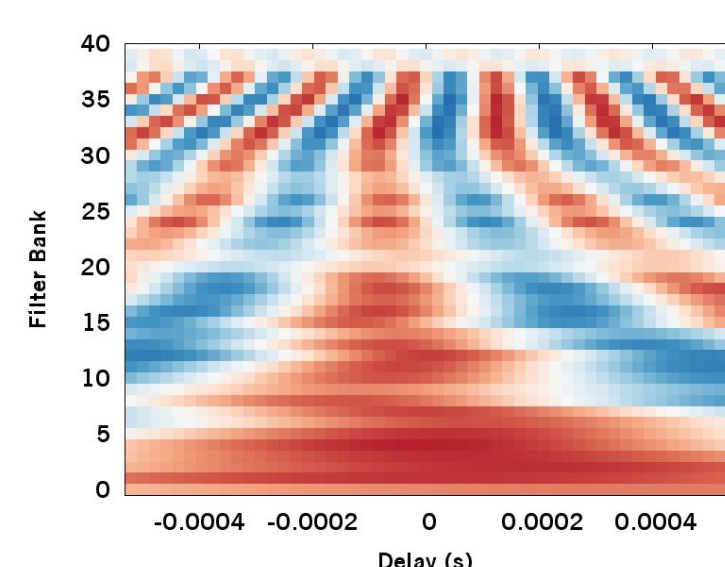4-channel audio → Features → Neural Network → Localize

**Features: Cross-correlation between pairs of microphones**
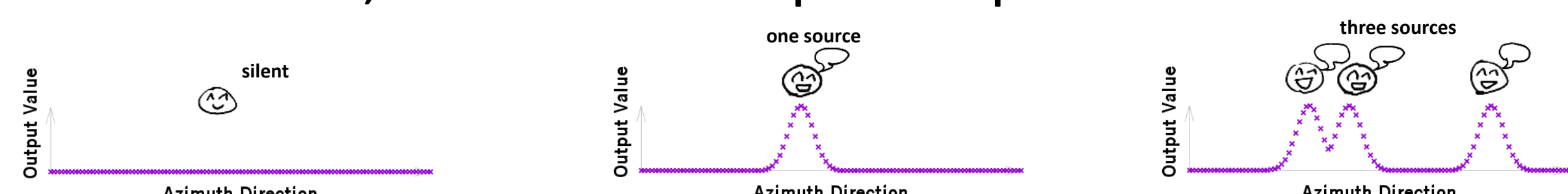
GCC-PHAT coefficients    ×6 pairs

GCC-PHAT on filter bank    ×6 pairs

**Output: Likelihood of sound source being in each direction**
- **Encoding:** Gaussian functions around true sources
- **Decoding:** Finding peaks
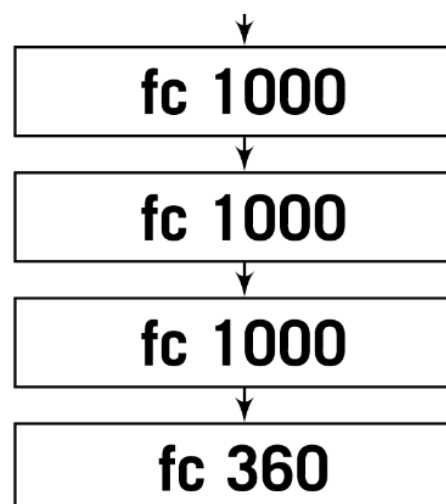- No softmax; Resembles a spatial spectrum

silent    one source    three sources

**Different network architectures:**
- **Multi-layer Perceptron** (MLP-GCC)
  - Basic structure (baseline)
- **Convolutional neural network** (CNN-GCCFB)
  - Convolution to reduce number of parameters
- **Two-stage network** (TSNN-GCCFB):
  - Considers the sparsity of speech signal in time-frequency points
  - First predict on sub-bands
  - Then aggregate early predictions across all frequencies.
  - Training also in 2 steps: (1) Pretrain first subnet (2) End-to-end
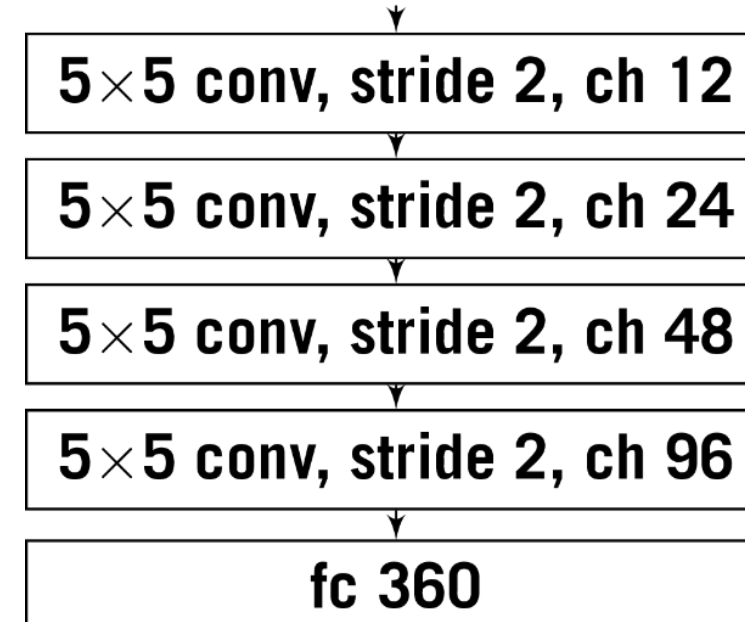
### MLP-GCC
GCC-PHAT (51×6)
| fc 1000 |
| fc 1000 |
| fc 1000 |
| fc 360 |
DOA Likelihood (360)

### CNN-GCCFB
GCC-FB (51×40×6)
| 5×5 conv, stride 2, ch 12 |
| 5×5 conv, stride 2, ch 24 |
| 5×5 conv, stride 2, ch 48 |
| 5×5 conv, stride 2, ch 96 |
| fc 360 |
DOA Likelihood (360)

### TSNN-GCCFB
GCC-FB ×6
Filter bank (40)    Delay (51)    Subnet 1    in: 51×5×6    out: 360
Latent Feature
Filter bank (36)    DOA (360)
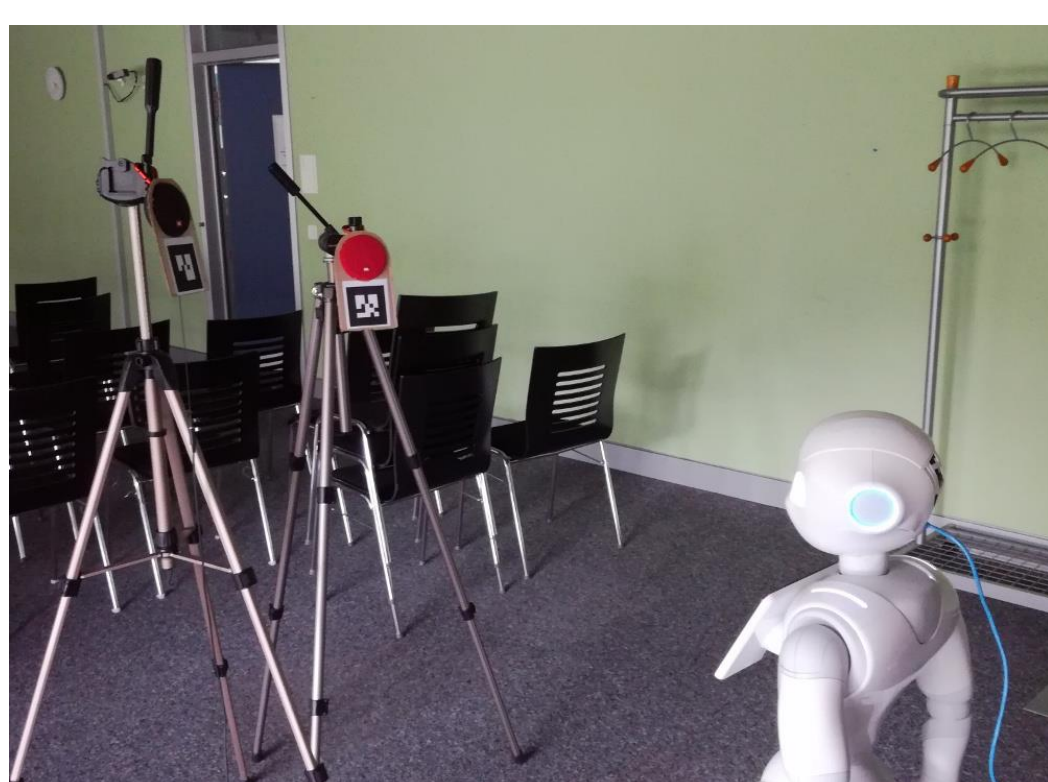in: 11×36    Subnet 2    out: 1
DOA Likelihood (360)

## Experiments

**Data:**
- 24 hours of real recordings of Pepper.
- Up to two simultaneous speakers.

Loudspeakers
(16h train / 8h test)

Human talkers
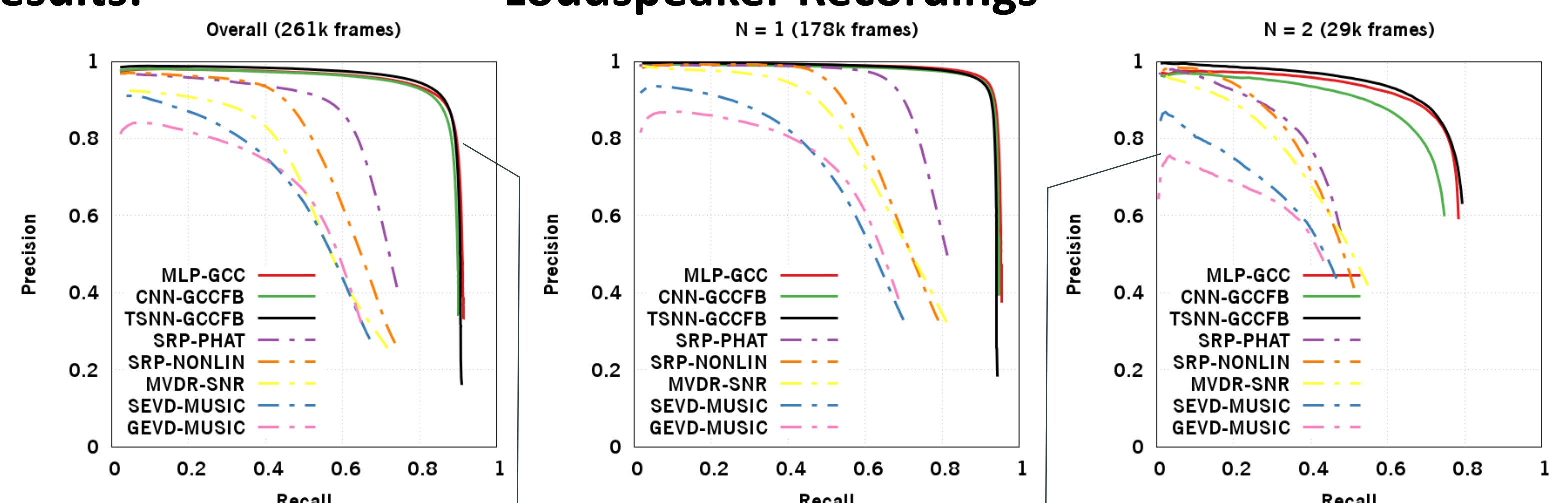(4 min test)

**Baseline methods:**
- SRP-PHAT, MVDR, MUSIC

**Evaluation:**
- Number of sources is unknown => detection problem
- Prediction is correct if error < 5°
- Compute precision vs recall

**Results:**

### Loudspeaker Recordings

Overall (261k frames)    N = 1 (178k frames)    N = 2 (29k frames)

MLP-GCC
CNN-GCCFB
TSNN-GCCFB
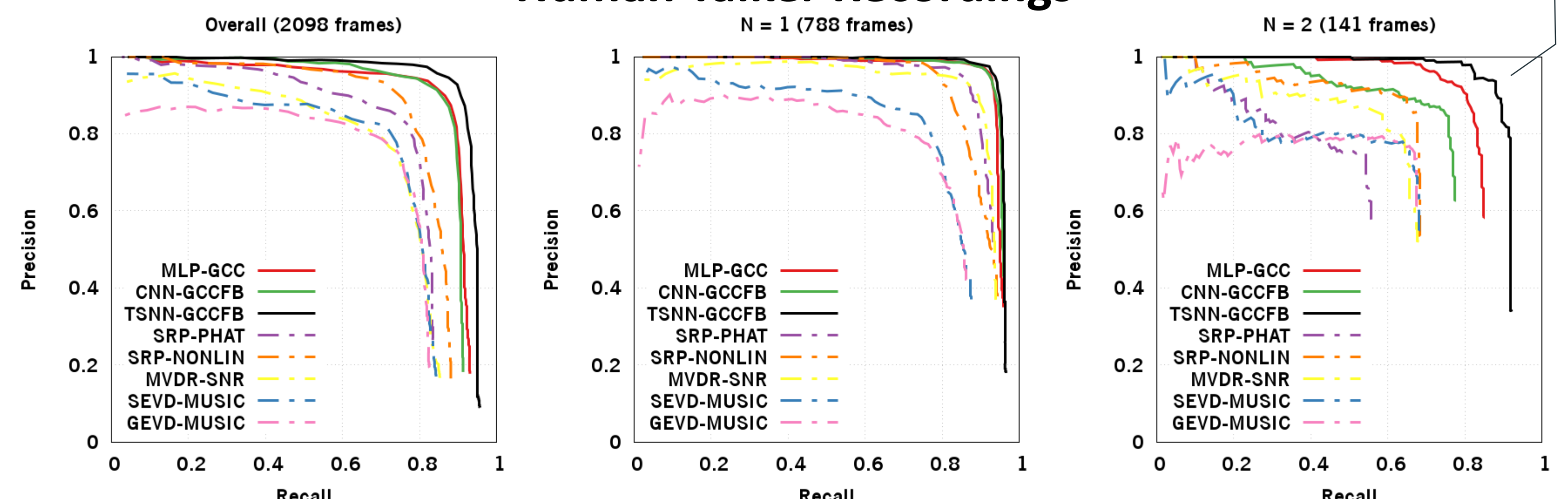SRP-PHAT
SRP-NONLIN
MVDR-SNR
SEVD-MUSIC
GEVD-MUSIC

1. Proposed methods have better overall performance

2. More significant with overlapping sources

3. Two-stage network performs the best

### Human Talker Recordings

Overall (2098 frames)    N = 1 (788 frames)    N = 2 (141 frames)

MLP-GCC
CNN-GCCFB
TSNN-GCCFB
SRP-PHAT
SRP-NONLIN
MVDR-SNR
SEVD-MUSIC
GEVD-MUSIC

## Conclusion

- >90% recall and precision.
- Significantly better than popular spatial spectrum methods.

## Resources

- Database:    https://www.idiap.ch/dataset/sslr
- Video:    https://youtu.be/_4EwuVlE_pU