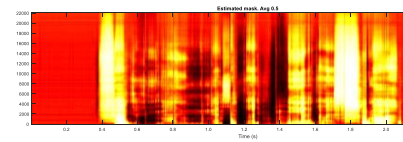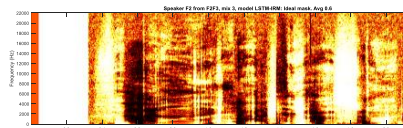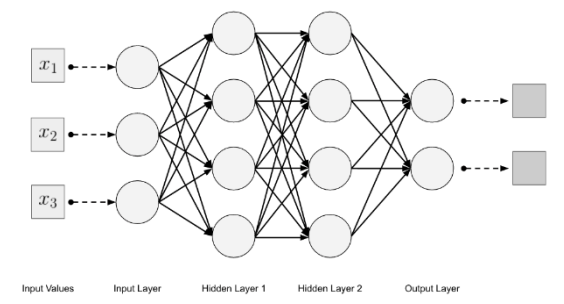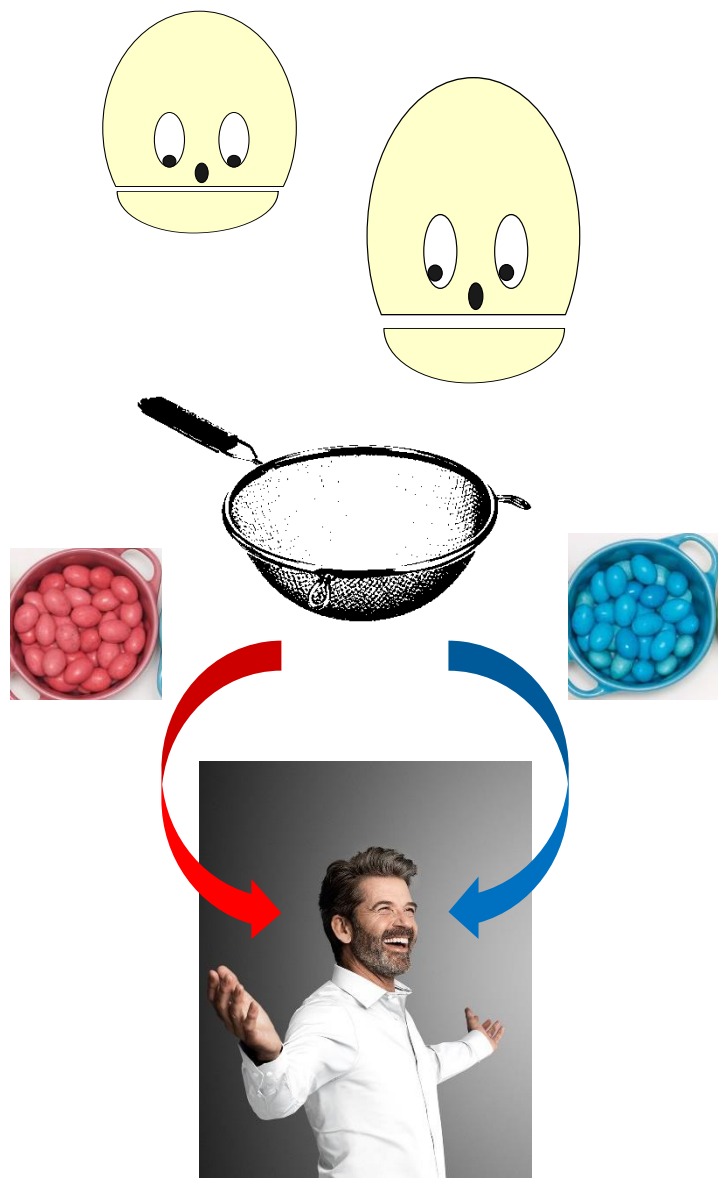# Overview

- Speaker separation: two competing voices

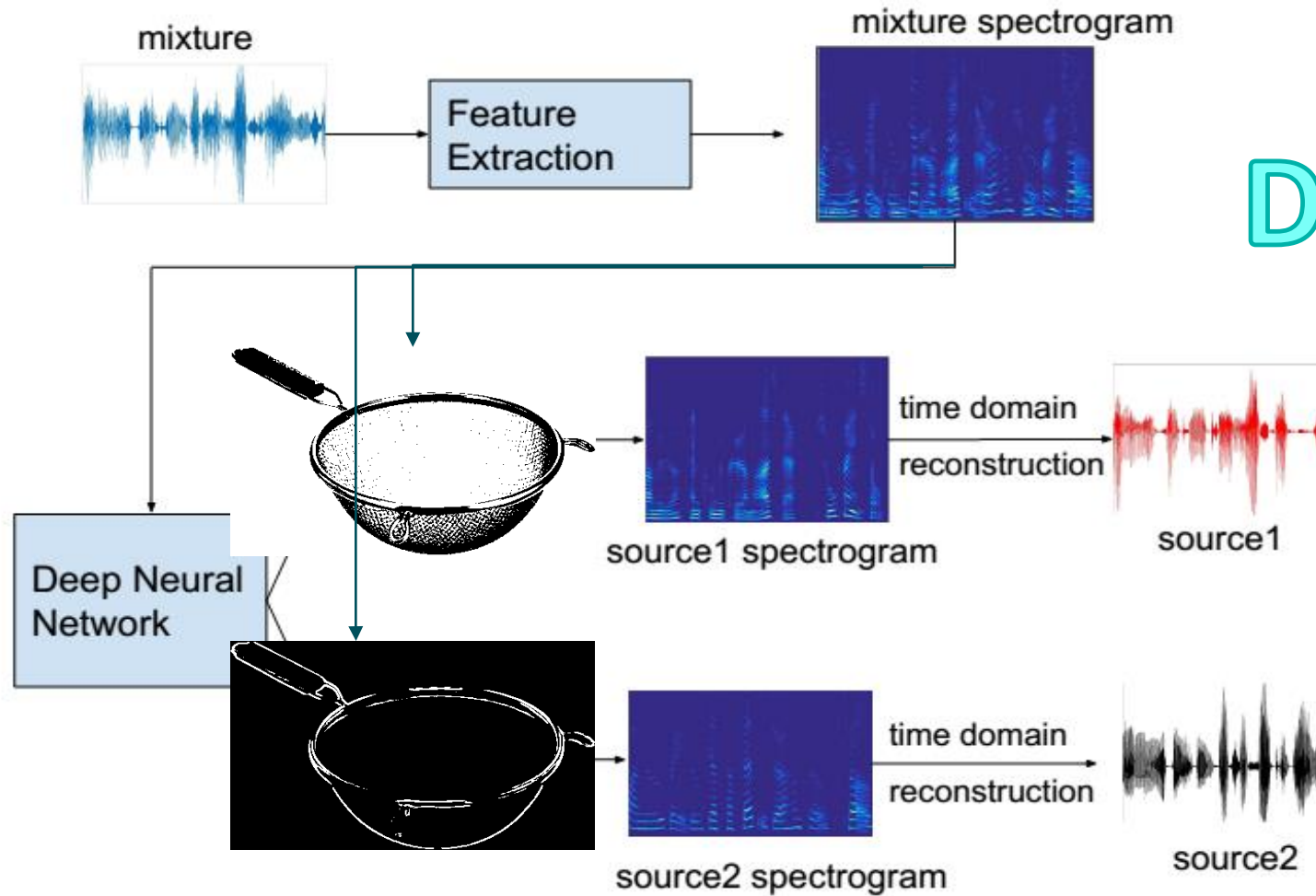- Noise reduction: voice-in-noise

- Summary

# Speaker separation (voice-on-voice)

# Separating two voices with low delay



DELAY 8 ms

# Training the DNN with the truth

# Ideal and estimated ratio mask



Speaker F2 from F2F3, mix 3, model LSTM-IRM: Ideal mask. Avg 0.6

Estimated mask. Avg 0.5

# Competing voices separated

## Derhjemme ser vi ikke ked

## Statuen har ikke noget hoved

- Example: pairs of sentences from the Danish Hearing In Noise Test
- Voices known in training
- Lots of glimpsing possible

M1 + M2                    ~M1                    ~M2

# Danish HINT material

- Overall: 13 lists of 20 sentences each
- Talkers: 3 male, 3 female (originally 1 male)
  - Combined in male-male, female-female and male-female pairs.
  - The DNN is speaker-specific, trained per speaker pair
- DNN validation: 1 list
- Training material: 4 lists
- Listening test: 8 lists
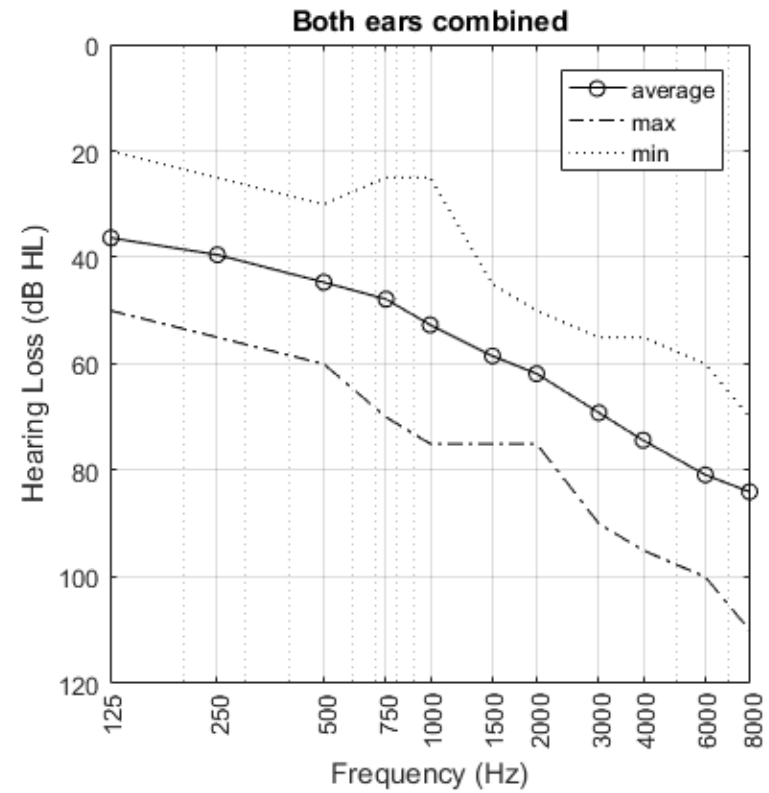
# Speech intelligibility benefit: Competing voices test

Statuen har ikke noget til hoved
Derhjemme spiser vi ikke kød



- Pairs of sentences from the Danish HINT (Hearing In Noise Test).
- Cueing by
  - First word = Single target
  - Last word = Dual targets: Competing Voices (CVT)

- Hearing loss compensated individually (NAL-R)

# 15 hearing-impaired listeners

# Processing

1. Sum (unprocessed)

2. Separate (ideal)

3. Feed Forward DNN (FDNN)

4. Long-Term Short-Term Memory Neural Net (LSTM)

5. Convolutional Recurrent Neural Net (CRNN)

Roughly 3.5 mio weights

Naithani et al, CHAT 2017, Stockholm
Naithani et al, WASPAA 2017, Mohonk

# Speech segregation results (competing voices)

# Speech separation results (single target)

# Different benefit for different people

**Noise reduction**

# Voice in noise

- Known and unknown voices in known noise
  - more common scenario
  - evaluate generalization ability
- New DNN+mask candidates

# Named DNN conditions

1. Sum (= input)
2. FDNN known voice
3. LSTM known voice
4. LSTM unknown voice
5. LSTM unknown voice + multi resolution mask
6. LSTM unknown voice + phase sensitive mask
7. Ideal ratio mask

Maximum 20 dB attenuation (except 7.)

Roughly 3.5 mio weights

# Test stimuli

- Danish HINT sentences
  - M1-M6, F1-F6 (12 talkers)
  - 200 – 260 sentences ~ 6 min

- Target talkers:
  - M1, M2, F1, F3
  - Speaker dependent: train on these (test other sentences)
  - Speaker independent: do <span style="color:red">not</span> train on these (test all sentences)

- Noise from the 'ICRA natural sound library'
  - P1: Party noise
    - train at -3..+3 dB SNR, test at +0 dB
  - S1: Shopping center noise
    - train at -3..+3 dB, test at +0 dB.

https://icra-audiology.org/

# Voice-on-noise test

## Statuen har ikke noget hoved

- Sentences from the Danish HINT
- 0 dB SNR

M1 + P1                 ~M1

M1 + S1                 ~M1

# Results: post hoc



HINT Test: Speech Reception Scores

noise
- P1
- S1

Processing
1. Sum (= input)
2. FDNN known voice
3. LSTM known voice
4. LSTM unknown voice
5. LSTM unknown voice + multi resolution mask
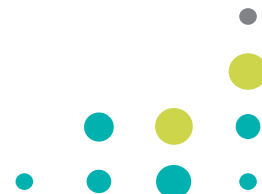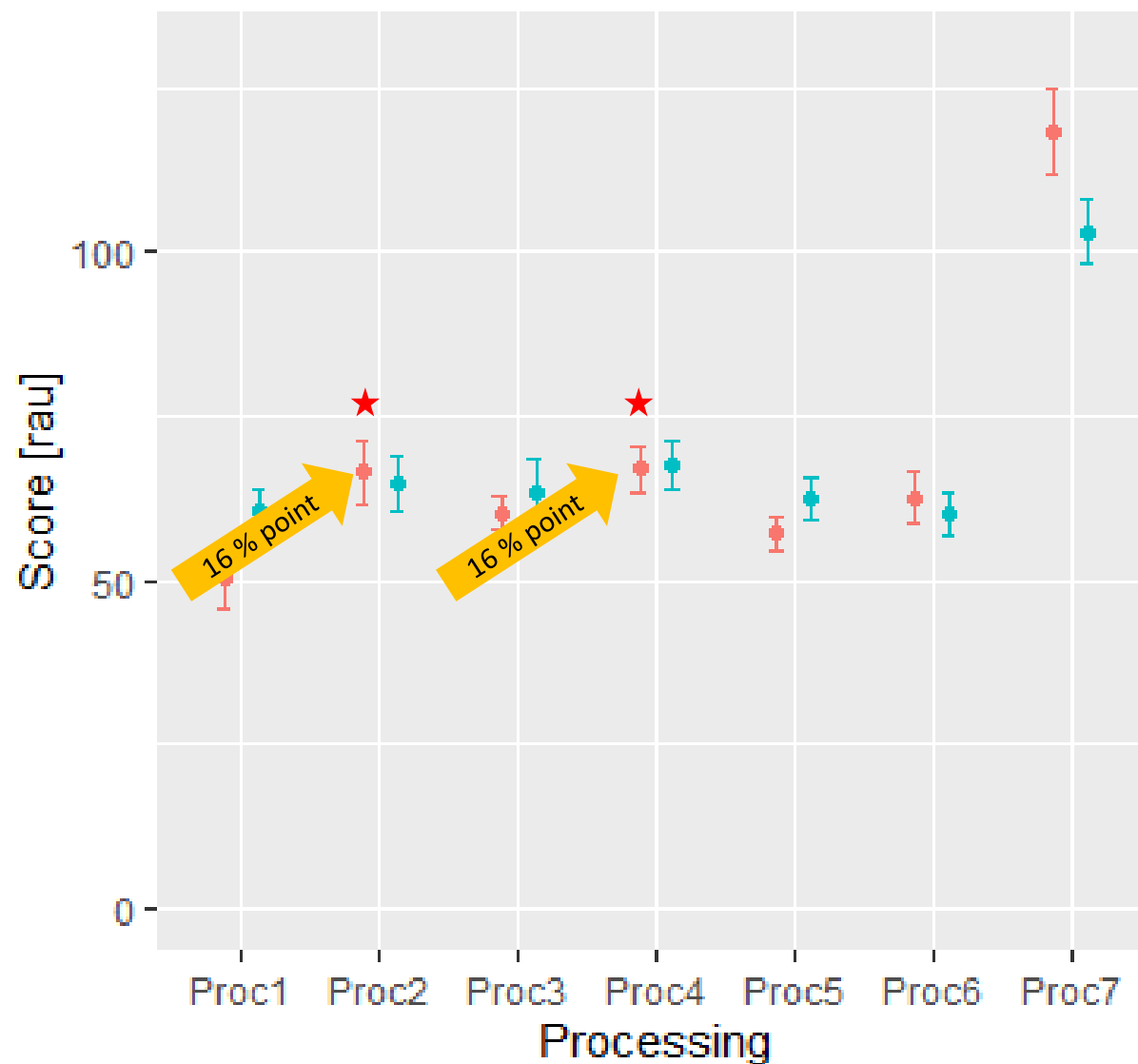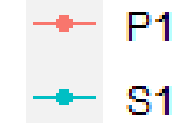6. LSTM unknown voice + phase sensitive mask
7. Ideal ratio mask

# Conclusion, speaker separation

- Competing voices test: Relevant, significant effect = 13% point. The user has all the information!

- Target-masker test: Large effect = 37% point The user must chose!

- All DNN modes (topologies) give the same improvement.

Bramsløw, L., Naithani, G., Hafez, A., Barker, T., Pontoppidan, N. H., and Virtanen, T. (2018). "Improving competing voices segregation for hearing impaired listeners using a low-latency deep neural network algorithm," J. Acoust. Soc. Am., 144, 172–185. doi:10.1121/1.5045322

# Conclusion, noise reduction

- Party noise: ~16 %-point (1.5 dB)
  - known voice FDNN
  - unknown voice LSTM

- Shopping centre: no benefit
  - Less modulated = less glimpses

- Unknown noise is a challenge!

**Thank you!**