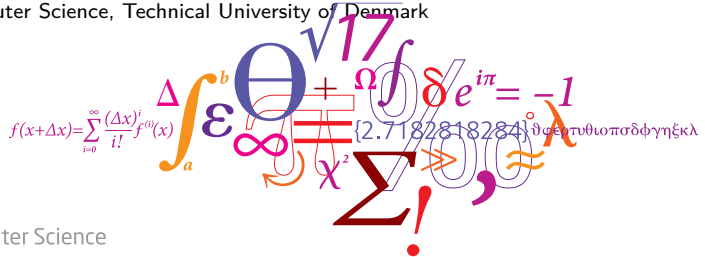# Feasibility of Few-shot Learning for the Automatic Transcriptions of Clinical-child Conversation in Danish

Speech Processing for clinical conversations

Sneha Das

Department of Applied Mathematics and Computer Science, Technical University of Denmark
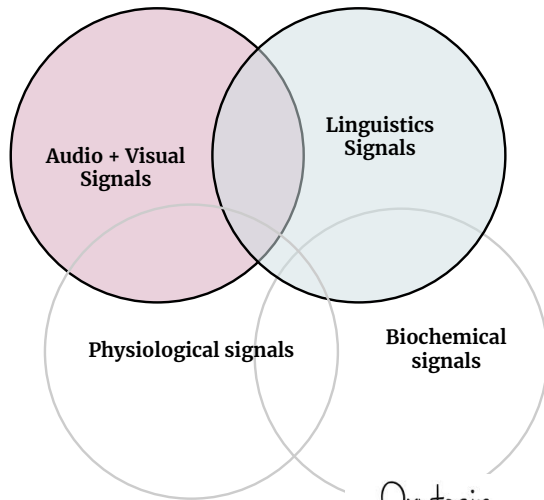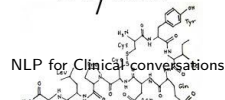
# WristAngel: Research for Intervention and Management of OCD

\* Progression and severity of disorder.

\* Improve efficiency in CIB (Coding Interactive Behavior)

¤ Identify and predict impending OCD events.

¤ Aid in delivering cognitive behavioral therapy to patients.

¤ Provide useful interventions for management.

SPEECH

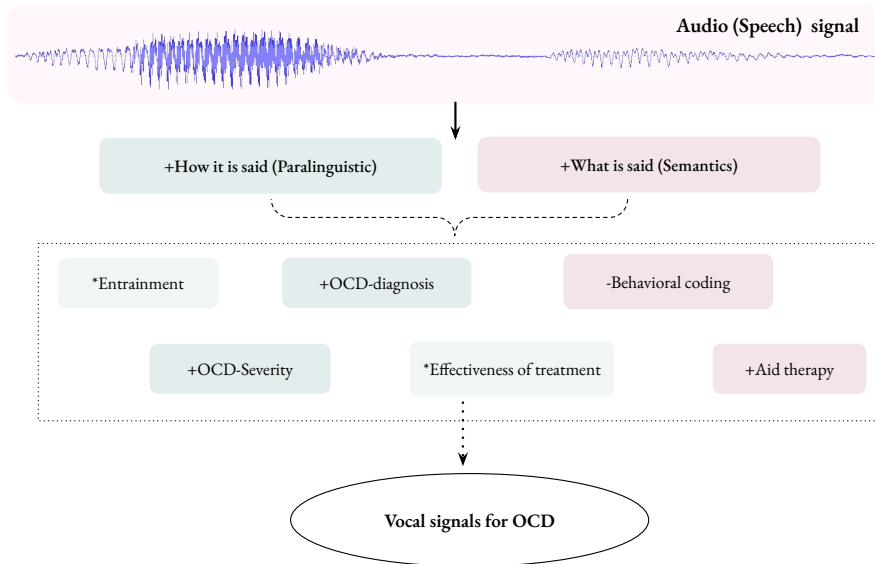## Audio (Speech) in OCD Management

Audio (Speech) signal

+How it is said (Paralinguistic)    +What is said (Semantics)

*Entrainment    +OCD-diagnosis    -Behavioral coding

+OCD-Severity    *Effectiveness of treatment    +Aid therapy

Vocal signals for OCD

## Speech preprocessing

❶ Pre-processing: Conversations → speech segments.

❷ Manual pre-processing: resource intensive

❸ Approx. 13 minutes /per minute of annotation → 260 individual hours for annotating 10 minute long audio conversation for 120 audio samples.

❹ Popular approach: ML pre-trained models pre-processing.
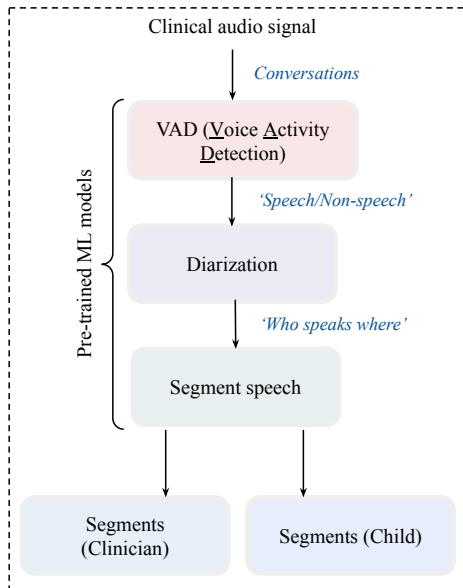


Figure: ML pre-processing pipeline.

# Speech pre-processing
## Challenges:

- Performance difference between clinicians and children.

- Errors (variance of error) higher for children in patient group.

- Correlation between error and OCD-severity score!!!



False negative rates over groups



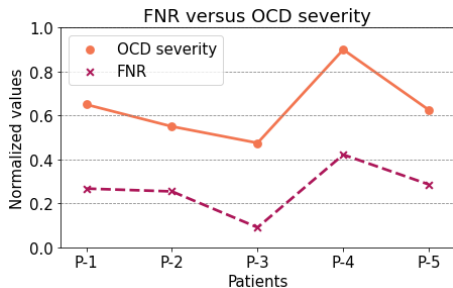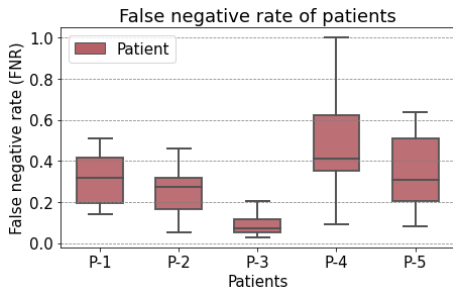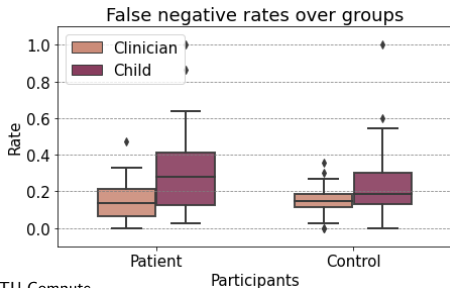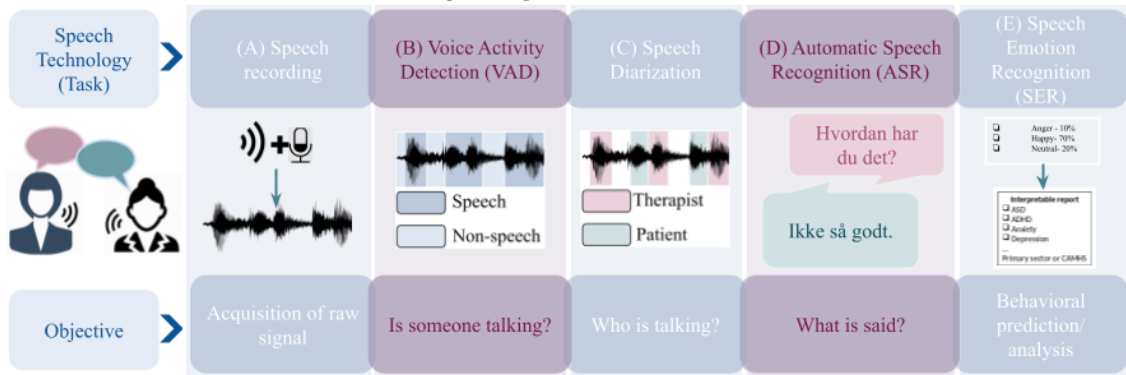False negative rate of patients



FNR versus OCD severity

Figure 1: Speech tasks

# Automatic Speech Recognition and Transcriptions

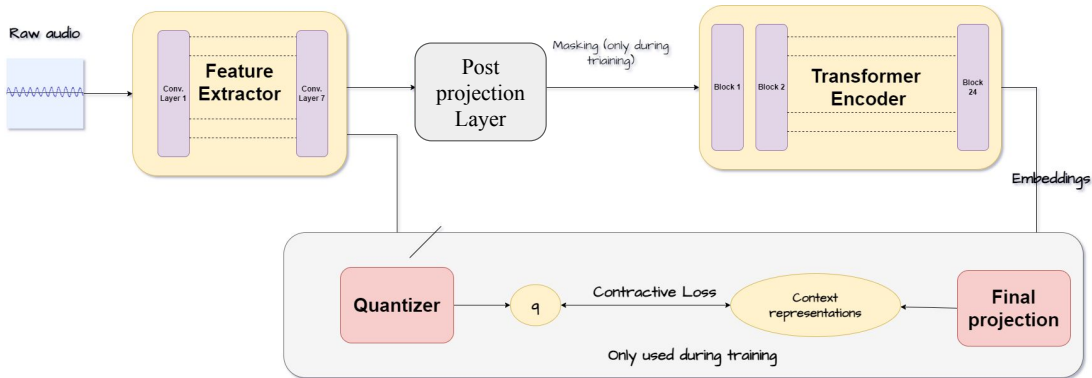- Clinical documentation

- Screening, diagnosis, management.

## Automatic Speech Recognition and Transcriptions

**1** State-of-the-art Models $\rightarrow$ English + Adults

**2** State-of-the-model for Danish $\rightarrow$ Alvenir

**3** Challenges:

- Transcribe speech from children in Danish
- Clinical conversations between clinician and child.
- Do we have data?

**Baseline and Wav2vec Model**

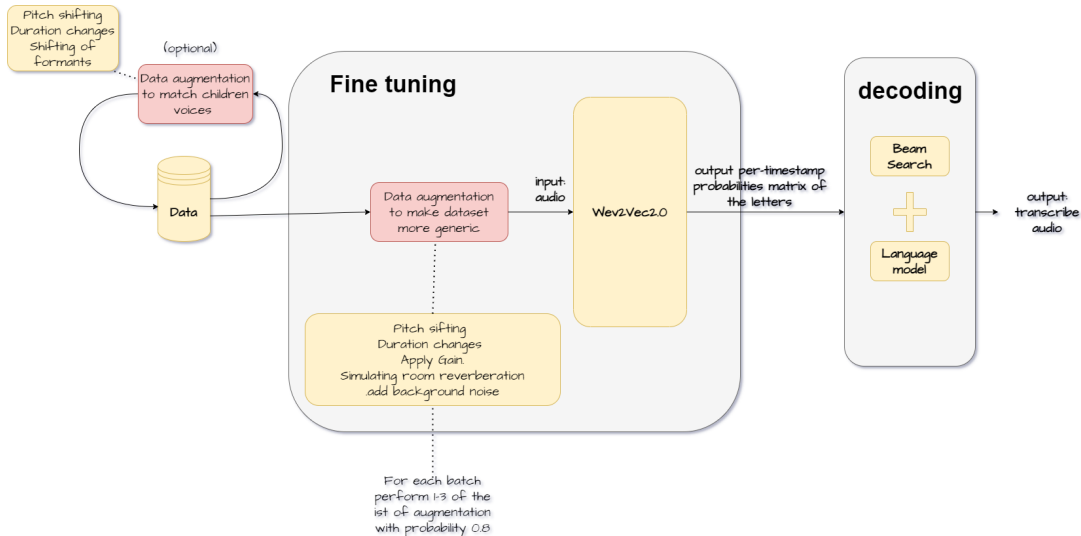**What to do when no data?**

## Data-augmentation
To aid in generalisation

- Gain change

- Reverberation

- Background noise

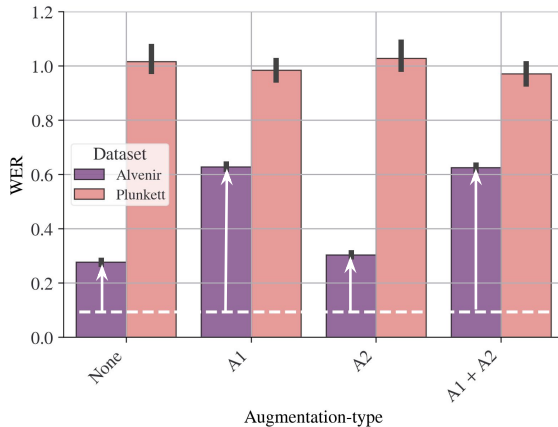- pitch and duration modification

To aid in transfer to children

- Formant-shift

- Pitch modification

- Duration modification

# Data augmentation



Pitch shifting
Duration changes
Shifting of formants

(optional)

Data augmentation to match children voices

**Fine tuning**

Data

Data augmentation to make dataset more generic

input: audio

Wev2Vec2.0

output per-timestamp probabilities matrix of the letters

Pitch sifting
Duration changes
Apply Gain
Simulating room reverberation
add background noise

For each batch perform 1-3 of the ist of augmentation with probability 0.8

**decoding**

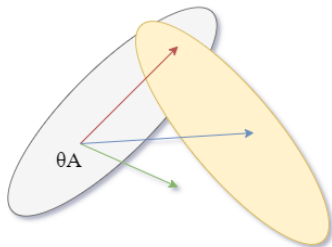Beam Search

+

Language model

output: transcribe audio

## Data augmentation



- Testing on Alvenir + Plunkett
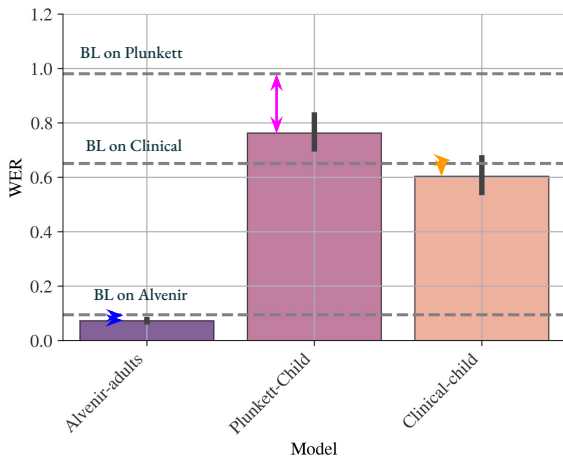- Catastrophic forgetting → Not acceptable (!)

## How to avoid Catastrophic forgetting?

- Weight freezing
  - Acoustic variability
  - Pronunciation variability
- Elastic weight consolidation: $L(\theta) = L_B(\theta) + \sum_i \frac{\lambda}{2} F_i (\theta_i - \theta_{A,i}^*)^2$



| Low error at task B | EWC |
| Low error at task A | L2 |
| | no penalty |

θA

## Results

Performance of the best model[1]

[1] Garofalaki. M, Speech and natural language processing for clinical in-the-wild data 2023.

# Affective-states from speech[1,2,5,6] Applications:

- Entrainment

- Vocalization

- Behavioral coding

---

[1] S. Das et al, Towards Transferable Speech Emotion Representation: On loss functions for cross-lingual latent representations, ICASSP 2022.

[2] S. Das et al, Continuous Metric Learning For Transferable Speech Emotion Recognition and Embedding Across Low-resource Languages, NLDL 2022.
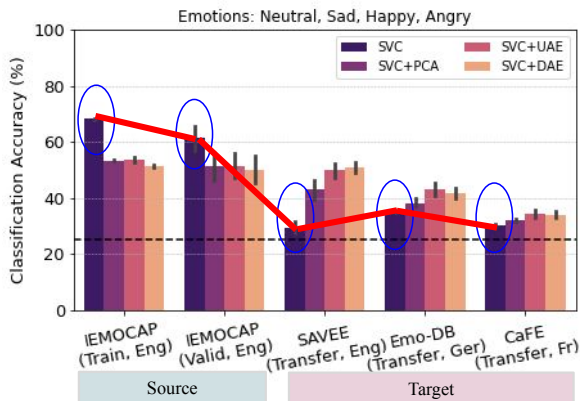
[5] S. Das et al, Zero-shot Cross-lingual Speech Emotion Recognition: A Study of Loss Functions and Feature Importance, ISCA SPSC Symposium 2022.

[6] Clemmensen et al, Associations between OCD severity and vocal features in children and adolescents: A statistical and machine learning analysis plan, JMIR Protocols.
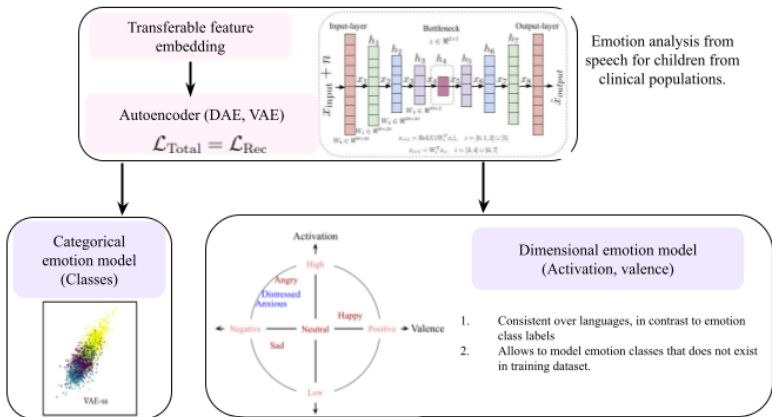
## Audio-features → (Simple!) Emotion-recognition

- Input-features: descriptive features of speech features ($f_0$, tonality, intonation, etc) -features $R^{88 \times 1}$ → Support vector machine (SVM) [Das, S, et al. 2021]


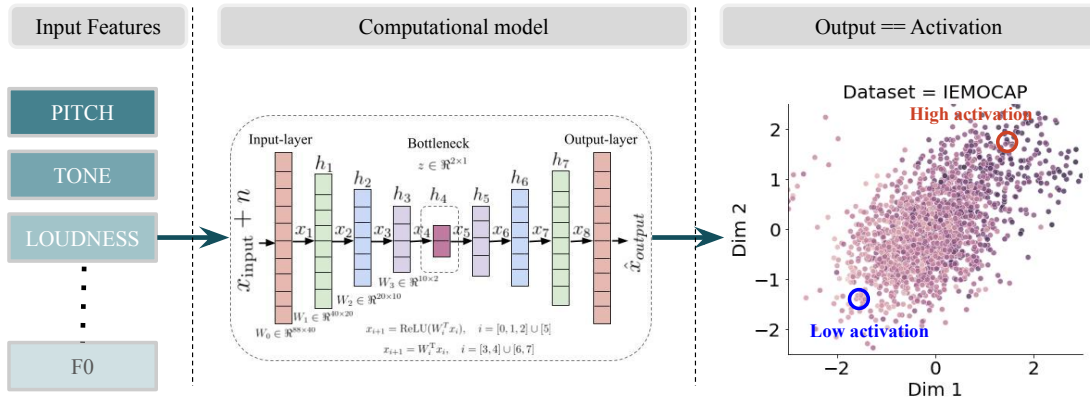
Emotions: Neutral, Sad, Happy, Angry

## Transferability: What variable to condition on?

Emotion class (discrete) or dimensional model (continuous)?

## Universal emotion representation

## Summary

❶ Speech-processing in psychiatry and psychology $\rightarrow$ accelerate and aid

❷ Challenges:

- Models are sensitive to language, age...
- Lack of resources (data, labels)

❸ How to adapt ASR modelled on adults to children with above challenges?

- Augmentation
- Continual learning $\rightarrow$ Elastic weight consolidation.

❹ Performance on adults maintained

❺ Performance on children improved by $\rightarrow$ 80%, 5%

❻ (Large!) room to improve.

# Thankyou!
Email: sned@dtu.dk; Twitter: @dassneh