



# Zero-shot Cross-lingual Speech Emotion Recognition: A Study of Loss Functions and Feature Importance

Sneha Das (sned@dtu.dk), Nicole Nadine Lønfeldt<sup>2</sup>, Nicklas Leander Lund<sup>1</sup>,  
Anne Katrine Pagsberg<sup>2,3</sup>, Line H. Clemmensen<sup>1</sup>

<sup>1</sup>Department of Applied Mathematics and Computer Science, Technical University of Denmark

<sup>2</sup>Child and Adolescent Mental Health Center, Copenhagen University Hospital, Capital Region

<sup>3</sup>Faculty of Health, Department of Clinical Medicine, Copenhagen University



## Motivation

- Speech analysis for children with OCD → Conversations in Danish
- Paralinguistic analysis → Speech emotion recognition (SER): inferring emotional state from speech signals.
- Challenges: Very little of data (Danish speech from children)
- Generalization and black-box nature of models → impacts security of systems.

**Zero-shot Cross-lingual SER, with focus on loss-functions.**

## Research questions

- ① Which model is more suitable for learning emotion representations from speech that are transferable over languages, categorical or dimensional model.
- ② Does semi-supervision aid in transferring emotion representations over languages, without labelled data in the target domain?
- ③ How do different learning functions influence the feature attribution scores?

## Methods

- Low complexity architecture, same architecture for all methods.
- Semi-supervision through loss function:
  - ① Cluster-loss → Learning emotion classes
  - ② Continuous metric-loss → Learning dimensional model of emotions → Activation, valence.

## Evaluation

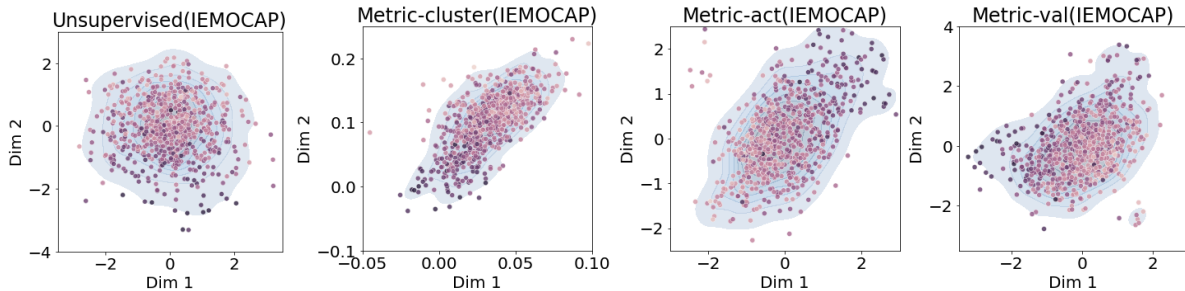
- Datasets: IEMOCAP (Training), SAVEE, Emo-DB, CaFE, URDU, AESD (Transfer)
- Input features: eGeMAPS using OpenSmile
- Preprocessing: remove outliers using z-score normalization ( $-10 > z > 10$ )
- 5-fold cross validation

## Rank correlation analysis

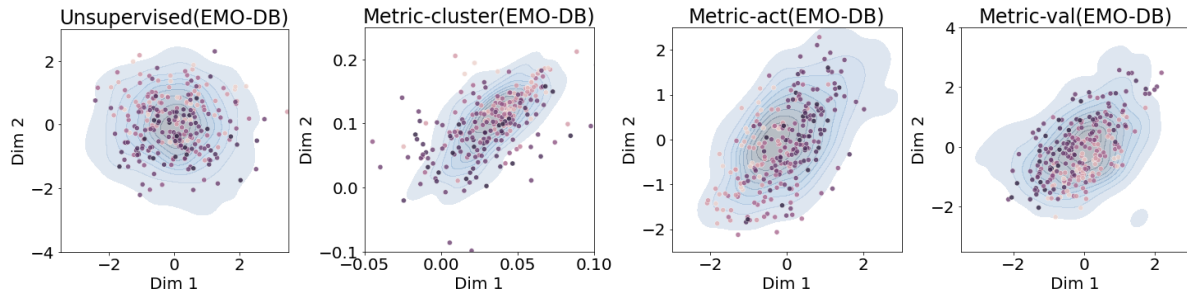
Table: Spearman's rank order correlation for the validation and transfer datasets aggregated over all model runs with different folds and random initial seeds. Higher correlation implies a larger correspondence to the ground truth labels (activation).

Method	IEMOCAP ( $\mu \pm \sigma$ )		EMO-DB ( $\mu \pm \sigma$ )		CAFE ( $\mu \pm \sigma$ )		URDU ( $\mu \pm \sigma$ )		AESD ( $\mu \pm \sigma$ )	
	Transfer	Supervised	Transfer	Supervised	Transfer	Supervised	Transfer	Supervised	Transfer	Supervised
Unsupervised	0.26 $\pm$ 0.17	0.26 $\pm$ 0.17	0.31 $\pm$ 0.22	0.31 $\pm$ 0.22	0.24 $\pm$ 0.14	0.24 $\pm$ 0.14	0.12 $\pm$ 0.1	0.1 $\pm$ 0.07	0.18 $\pm$ 0.11	0.16 $\pm$ 0.09
Metric-cluster	0.19 $\pm$ 0.14	0.19 $\pm$ 0.14	0.23 $\pm$ 0.16	0.28 $\pm$ 0.19	0.12 $\pm$ 0.08	0.07 $\pm$ 0.04	0.07 $\pm$ 0.06	0.09 $\pm$ 0.07	0.12 $\pm$ 0.06	0.11 $\pm$ 0.05
Metric-act	<b>0.76 <math>\pm</math> 0.05</b>	<b>0.76 <math>\pm</math> 0.05</b>	<b>0.53 <math>\pm</math> 0.08</b>	<b>0.61 <math>\pm</math> 0.04</b>	<b>0.35 <math>\pm</math> 0.04</b>	<b>0.39 <math>\pm</math> 0.03</b>	<b>0.38 <math>\pm</math> 0.05</b>	<b>0.39 <math>\pm</math> 0.05</b>	<b>0.31 <math>\pm</math> 0.01</b>	<b>0.31 <math>\pm</math> 0.01</b>
Metric-val	0.29 $\pm$ 0.11	0.29 $\pm$ 0.11	-0.05 $\pm$ 0.03	0.27 $\pm$ 0.24	0.31 $\pm$ 0.09	0.32 $\pm$ 0.1	0.03 $\pm$ 0.08	0.07 $\pm$ 0.1	0.01 $\pm$ 0.05	0.14 $\pm$ 0.1

## Rank correlation analysis

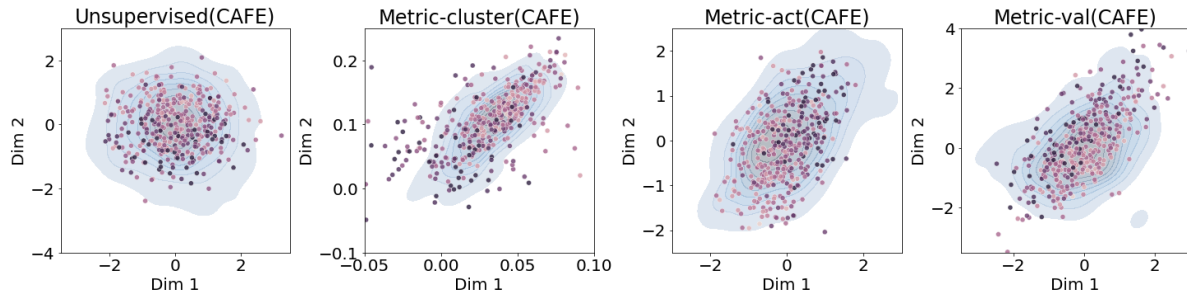


# Rank correlation analysis

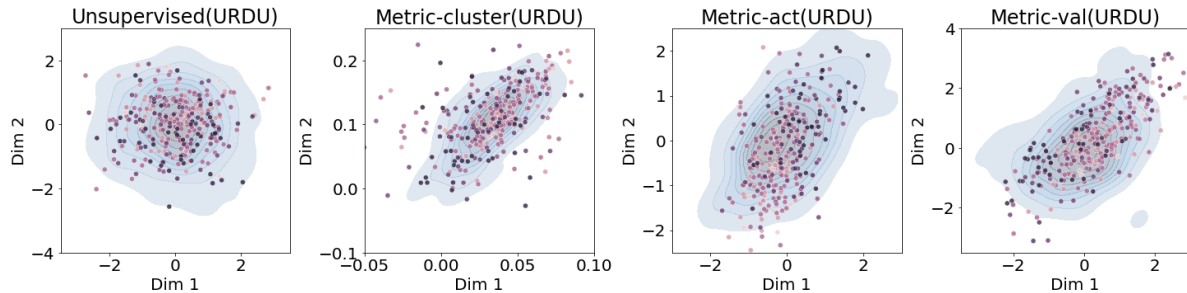




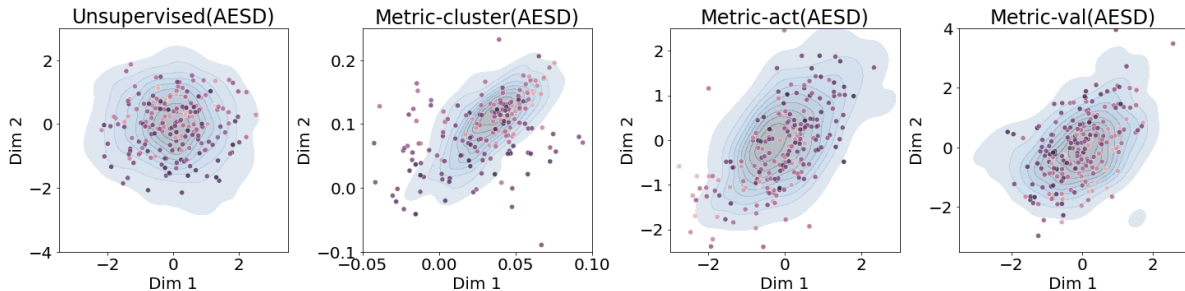
# Rank correlation analysis



## Rank correlation analysis



# Rank correlation analysis



## Rank correlation analysis

- ① Metric-act consistently shows the highest Spearman's correlation coefficient; the Unsupervised model seems to perform better relative to Metric-cluster and Metric-val.
- ② Distribution of the embedding in the latent space is most consistent for metric-act and closely followed by the Unsupervised model over the transfer dataset.
- ③ The correlation coefficient is relatively lower for CAFE, URDU and AESD.

## Classification Accuracy

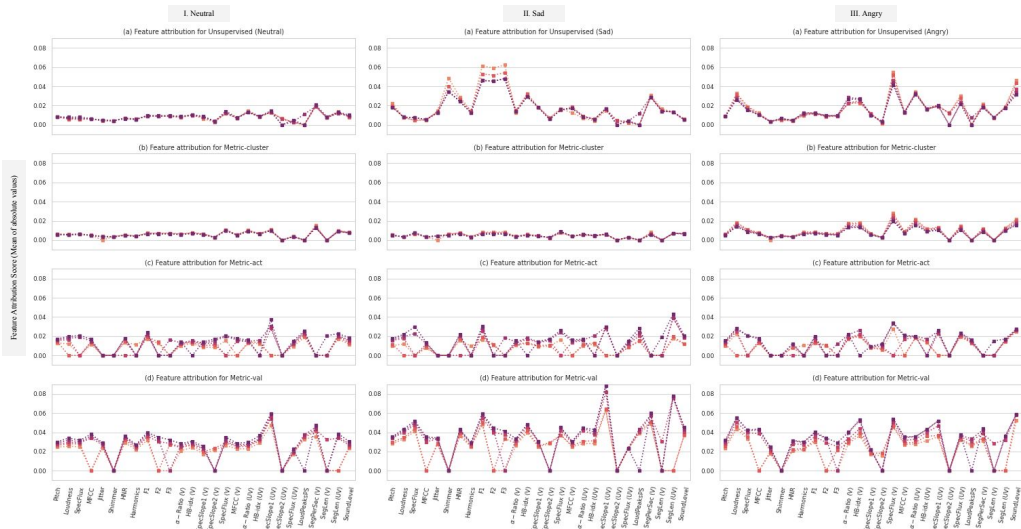
Table: Balanced classification accuracy over 4-class emotion classification. The model is trained on the IEMOCAP dataset only.

Method	IEMOCAP ( $\mu \pm \sigma$ )	EMO-DB ( $\mu \pm \sigma$ )	SAVEE ( $\mu \pm \sigma$ )	CAFE ( $\mu \pm \sigma$ )	URDU ( $\mu \pm \sigma$ )	AESD ( $\mu \pm \sigma$ )
	N-S-H-A	N-S-H-A	N-S-H-A	N-S-H-A	N-S-H-A	S-H-A
Unsupervised	$0.45 \pm 0.01$	$0.39 \pm 0.03$	$0.46 \pm 0.02$	$0.25 \pm 0.02$	$0.31 \pm 0.02$	$0.17 \pm 0.02$
Metric-cluster	<b><math>0.57 \pm 0.03</math></b>	<b><math>0.47 \pm 0.05</math></b>	$0.45 \pm 0.05$	<b><math>0.33 \pm 0.02</math></b>	$0.32 \pm 0.03$	$0.31 \pm 0.05$
Metric-act	$0.49 \pm 0.01$	$0.44 \pm 0.04$	$0.43 \pm 0.03$	$0.26 \pm 0.03$	<b><math>0.33 \pm 0.03</math></b>	$0.24 \pm 0.03$
Metric-val	$0.52 \pm 0.02$	$0.45 \pm 0.04$	<b><math>0.47 \pm 0.03</math></b>	$0.31 \pm 0.03$	$0.31 \pm 0.03$	$0.25 \pm 0.04$

## Classification accuracy

- ① Metric-cluster shows highest classification performance among all the methods.
- ② Cluster loss function was successful in discriminating between the classes on the training dataset. However, from the scatter points we observe that the remaining datasets have relatively larger covariances that overlap between emotions. This could indicate lower transferability of emotion representations with the metric-cluster loss function.
- ③ Latent samples from the Metric-act method are better contained within the kernel density plot of the training dataset. However, discrimination between the class labels are much lower relative to Metric-cluster, as can be seen from the higher overlap between datasets over the emotion classes.
- ④ *Anger* seems to be the most discernible from other classes, for most methods and datasets.

# Feature attribution analysis



## Feature attribution analysis

- Attribution scores are relatively lower for Unsupervised and Metric-cluster, in contrast to the Metric-act and Metric-val methods. Generally, the amplitude order of the mean absolute scores are  $\text{Metric-val} > \text{Metric-act} > \text{Unsupervised} > \text{Metric-cluster}$ .
- More number of feature groups for Metric-act seem to have no significant mean difference between emotion classes.
- Neutral (vs) Angry: Loudness, spectral flux and sound level are the feature groups that seem to have higher mean attribution scores for the emotion *Angry* over most of the methods considered.
- Neutral (vs) Sad: Formants and segment length show higher attribution scores.
- Metric cluster shows least changes in the attribution scores over different emotion classes.



## Conclusions

- Activation is more transferable to unseen cross-lingual datasets.
- Semi-supervision helps!
- Identified feature groups with influence.

### Limitations:

- Methods employed for transferability inspection increasing difficult for high dimension representations.
- Feature attribution methods here → systems with feature extractors.

## References



# Thankyou!

Email: [sned@dtu.dk](mailto:sned@dtu.dk); Twitter: [@dassneh](https://twitter.com/dassneh)