

# ENHANCING SOCIAL INTERACTIONS WITH AI-POWERED SYSTEMS

Nils Achy



## 1. INTRODUCTION

- AI-powered systems like chatbots, robots, and home assistants becoming increasingly popular.
- Conversations with these systems do not provide the same social interaction experience.
- Disparity caused by the absence of voice nuances, strict language use, and limited adaptation.
- Inability to recognize when someone in a group discussion wants to contribute.**
- If robots and conversational agents could detect this behavior, they could create a more engaging environment for discussion. Ex: serve as mediators to ensure that everyone has the opportunity to express their thoughts (information gain).

## 2. RELATED WORK

- Previous work by Litian Li et al. [1] aims to infer instances of speaking intentions by training a model on accelerometer data using fixed window sizes.
- Does not provide a comprehensive understanding of the underlying structure of these intentions.
- When does the intention truly start?

[1] Jing Zhou Litian Li, Jord Molhoek. Inferring intentions to speak using accelerometer data in-the-wild, 2023. Unpublished.

## 3. CONTRIBUTION

- Infer segments (finding the start and end time) of speaking intentions (segmentation) on the same dataset using data captured by a body-worn accelerometer as input data to the model.
- Infer instances of speaking intentions (classification) using varying segment sizes of accelerometer data instead of fixed window sizes and compare with [1].

How can body language, captured by a body-worn accelerometer, be utilized to estimate segments of speaking intentions in time, and does a supervised learning process improve the performance of detecting such cases?

## 4. EXPERIMENT

- REWIND dataset (business networking event) of in-the-wild data of conversational speech for 1.5h in Dutch language (limitation)
- First half of the event participants engage on assigned topics, the second half in free conversation
- 13 people visible on camera during the 10-minute extract (1:00:00 to 1:10:00) wearing an audio recorder and accelerometer device

## 5. METHODOLOGY

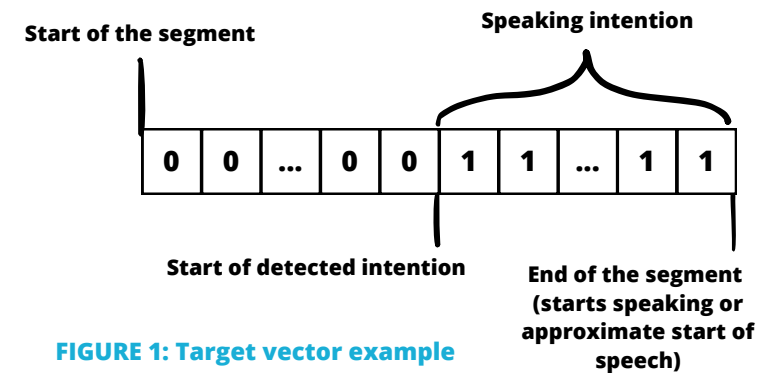
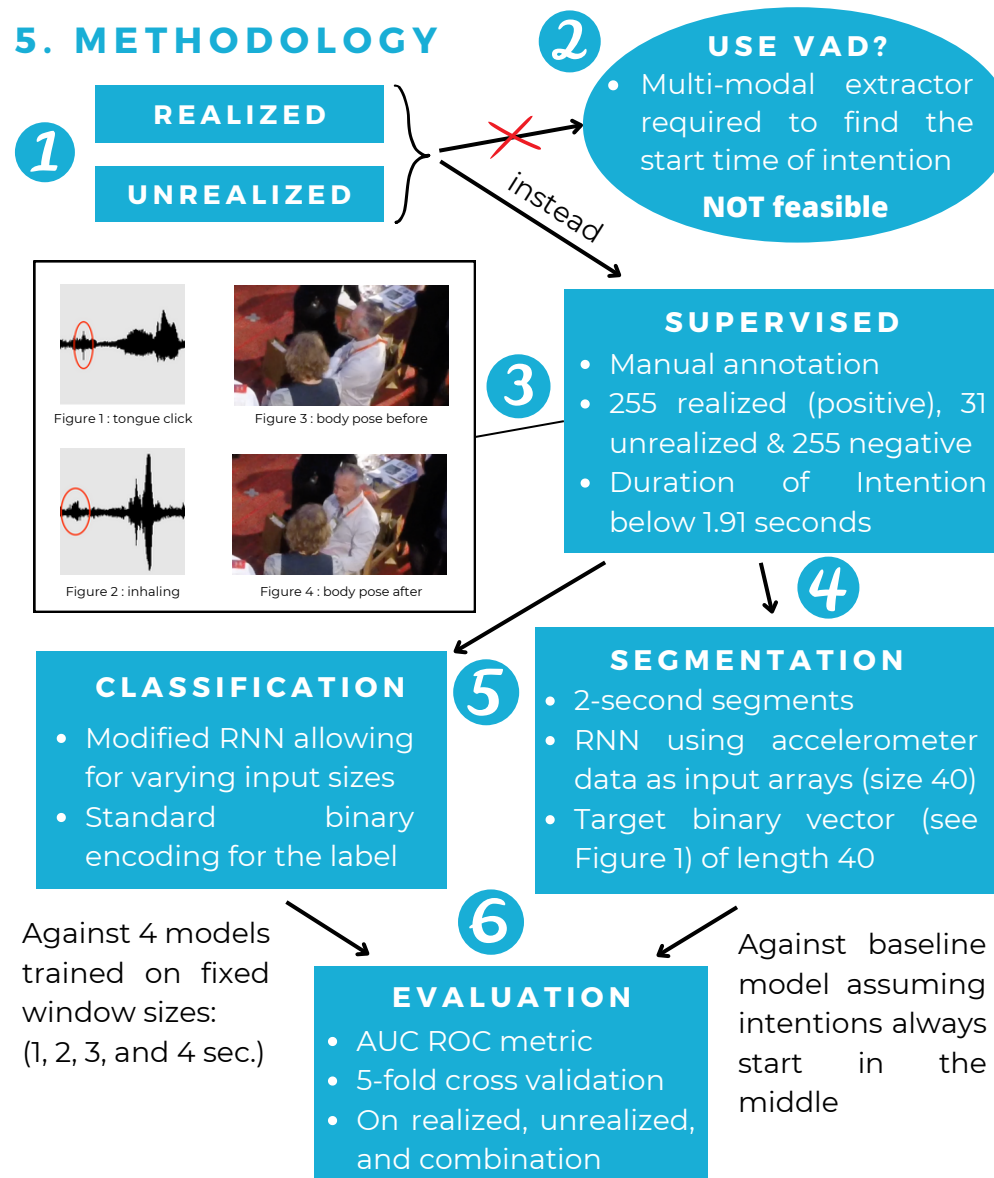
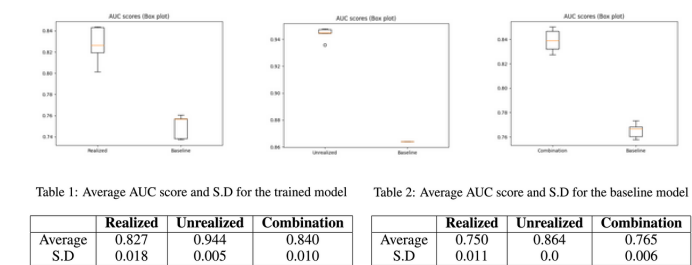


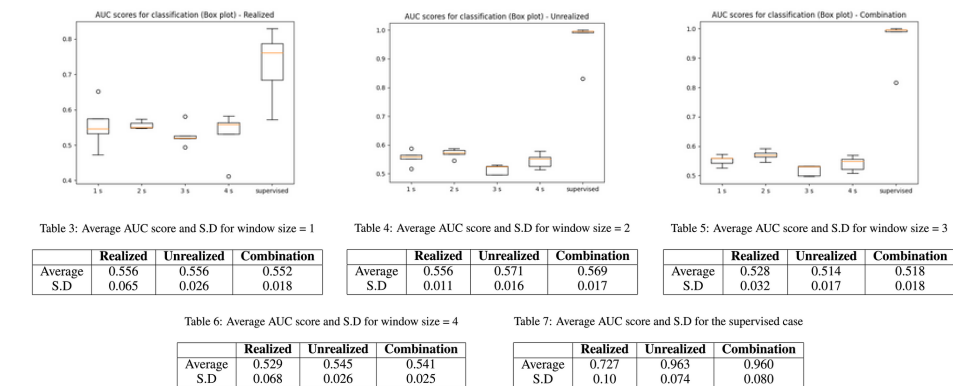
FIGURE 1: Target vector example

## 6. SEGMENTATION RESULTS



- Trained model outperforms the baseline model (label structure explains high scores)
- Results show a small standard deviation across the 5 folds.
- Lower performance was observed in realized intentions compared to unrealized intentions.
- Possible reason for lower performance is the inclusion of negative samples in the test set of realized intentions.
- Model may have reduced accuracy when dealing with negative samples, leading to a higher rate of false positives.
- Model tends to overestimate the occurrence of speaking intentions.
- Despite overestimation, the model still performs better than the baseline.

## 7. CLASSIFICATION RESULTS



- Supervised model outperforms models with fixed window sizes in all three evaluation criteria and in all 4 window sizes.
- Results show consistency and low standard deviation.
- Supervised learning provides benefits in intention inference.
- Alignment with findings from [1] for unrealized intentions (2-second window size yields the most promising outcomes).
- Misalignment with findings from [1] for realized intentions (2-second window size in the present research, 1-second for [1]).

## 8. CONCLUSION & LIMITATIONS

- Model trained on accelerometer data demonstrates effective segmentation capability within the 2-second segment. Performs better on positive instances but achieves higher AUC scores on all 3 criteria compared to the baseline model with low standard deviation.
- Supervised learning brings significant improvement to the classification task; suggesting that qualitative data can play a role in building more realistic estimators for inferring speaking intentions.
- Potential future work includes the use of a larger dataset, incorporating data from more languages, and employing a rule-based approach to build a multi-modal extractor (to use a VAD).
- Manual annotations of speaking intentions are subjective, recommending multiple individuals for accuracy and reliability
- Fine-tuning the model and assessing performance using different metrics are important considerations.
- Exploring these ideas can lead to advancements in our ability to infer speaking intentions and improve Human-Computer interactions.