Recognizing Speech with no Microphone: using Hidden Markov Models on Intraoral Sensor Data

Nathan Klumpenaar (M.J.N.Klumpenaar@student.tudelft.nl) Supervisors: Przemysław Pawełczak, Vivian Dsouza Part of the 2025 CSE3000 Research Project

Keyword Recognition

and train set.

accuracv.

selected features.

on the accuracy.

Component count: speech

IMU-X Zero crossing count

IMU-Z Standard deviation ALS First diff. mean

ALS Zero crossing count MU-Y Standard deviation

IMLI-Y First diff mean

rometer First diff. mean IMU-X First diff. mean

IMU-Z Zero crossing count

IMU-X Standard deviation

IMU-Y Zero crossing count

IMU-Z First diff, mean 💋

IMU-Y Mean

Prediction HMM Snee

Labeled coughing

2.5 5.0

Prediction HMM No-speed

rometer Zero crossing count

IMU-X Mean ALS Standard deviation

IMU-Z Mean

ALS Mean Barometer Mean

Component count: no speech

3

UDelft

Introduction

The Densor [1] is an intraoral sensor platform designed to capture unique data inside the human mouth. While primarily focused on sleep analytics, this thesis explores a different use case: its use for voice activity detection (VAD) and keyword recognition.

An initial literature survey and data exploration showed the success and ease of use of Hidden Markov Models (HMMs), leading to the research question:

Is it possible to detect and recognize speech with Hidden Markov Models using data gathered by the intraoral Densor sensing platform?

Feature Extraction Pipeline



Five sensor channels

- Barometer
- Ambient Light Sensor (ALS)
- Inertial Measurement Unit (IMU)
 - Reports X, Y, Z acceleration Fig. 2: Feature extraction

Limitations & Future Work

Limited variety in dataset: Prevented most session/user independent testing. • leading to limited reliability of results.

.

Mean value

Standard deviation

· Amount of zero crossings

First differential mean value

- Controlled position & environment: All recordings were done with a stationary user. Results would presumably be worse in a dynamic environment.
- Future work: A more (varied) dataset and exploring more advanced models.

Conclusions

- An F1-Score of 0.73 and a keyword recognition accuracy of up to 72% show that Densor-recorded data contains enough speech information for basic speech detection and recognition.
- These results are limited by the small and unvaried dataset. It is not expected that these models perform well on real world data.
- A bigger and more varied dataset can improve the models and possibly support more advanced models as well.

[1] Dsouza, V., Pronk, J., Peppelman, C., Madariaga, V. I., Pereira-Cenci, T., Loomans, B., & Pawełczak, P. (2024) Densor: An Intraoral Battery-Free Sensing Platform. Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous_Technologies, 8(4), 1-30. https://doi.org/10.1145/3699746

[2] Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). Optuna: A Next-generation Hyperparameter Optimization Framework. Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2623-2631. https://doi.org/10.1145/3292500.3330701

Methodology

Voice Activity Detection



Partition the training data into 1second segments labelled either speech or no-speech.

Automatically select the highestperforming féatures, based on *F1-Score.* 3

Train the final HMMs: one speech, one no-speech, using only the selected features.

Evaluate the performance based 5 on recall, precision and F1-Score.

Results

2

Voice Activity Detection

The best performing result, on a sessiondependent dataset, using all sensors:

- Precision: 0.74
- 0.73 Recall: F1-Score: 0.73













Non-boolean feature

Feature enabled

Feature disabled

Relative Importance

Fig 5. VAD automated feature selection results

10.0 12.5 15.0 17.5 20.0

Time (s)

7.5

	food	3	1	0	0	0	0			
	goodbye	0	4	0	0	0	0			
True	hello	0	0	7	1	0	0			
	start	0	0	0	2	5	1			
	stop	0	0	0	0	8	0			
	water	0	1	0	0	1	2			
70	food goodbye hello start stop wate Predicted									
	12% accuracy on a session-dependent									

uii a sessiuii-uei dataset using all sensors.

True	food		0	0	0	0	1
	goodbye	0	1	0	3	0	0
	hello	0	1	2	1	0	0
	start	0	0	0	4	0	0
	stop	0	0	0	0	3	1
	water	0	0	0	1	0	3
		food	goodbye	hello Pred	start icted	stop	water

67% accuracy on a session-independent. user-dependent dataset using all sensors. Fig. 6: Subset of keyword recognition results



Fig. 1: Photograph of a Densor, by Dsouza et al., 2024 [1].

feature

values