



**HAL**  
open science

# A spectrotemporal modulation application for distinguishing modal and whistled speech

Benjamin O'brien, Anna Marczyk

► **To cite this version:**

Benjamin O'brien, Anna Marczyk A spectrotemporal modulation application for distinguishing modal and whistled speech. *International Journal of Speech Technology*, 2025, <10.1007/s10772-025-10185-1>. <hal-05025403>

**HAL Id: hal-05025403**

**<https://hal.science/hal-05025403v1>**

Submitted on 8 Apr 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

# A spectrotemporal modulation application for distinguishing modal and whistled speech

## Abstract

Machine learning procedures have become efficient at using acoustic information to distinguish modal speech, but their application on other speech modalities remains largely overlooked. Operating under the hypothesis that non-spoken speech classification requires a more robust acoustic analysis in which to train models, this study applies a spectrotemporal modulation (STM) analysis to examine its effectiveness at distinguishing modal and whistled Spanish speech. A linear mixed-effects model revealed that in comparison to modal speech, whistled speech had significantly increased spectrotemporal modulations generally above 1 cycles per octave. Features identified as significantly relevant for distinguishing speech modalities were extracted, and an automatic speech modality classification task was developed. STM- and MFCC-feature models performed similarly, boasting weighted accuracy performances of 92% and 94%, respectively. Although an STM analytical approach had not previously been applied to whistled speech, our results support existing evidence on its acoustic characteristics. Moreover, our findings have key implications for using STM-features for automatic speech modality classification tasks, as it reduced the feature space at little cost to performance.

**Keywords:** automatic speech modality classification, whistled speech, non-speech vocalisations, spectrotemporal modulation

## 1 Introduction

Machine learning procedures have become efficient at using acoustic information to distinguish modal speech, but their application on other speech modalities remains largely overlooked. Fortunately, recent work has applied automatic speaker verification methods to non-typical speech, e.g., pathological speech [1, 2] and non-spoken speech, such as whispered [3] and shouted speech [4]. Nevertheless, despite the productive and perceptual differences between spoken and non-typical speech, often-times the models are trained on traditional acoustic features, i.e., cepstral coefficients, associated with modal speech. It is plausible that models designed to distinguish different speech modalities might enhance performance when trained on acoustic characteristics associated with non-typical speech. However, this depends on

the extent to which these characteristics provide distinct and reliable patterns for distinguishing speech modalities.

One example of a non-typical speech is non-spoken speech, which, in general, can be characterized as a vocalization used to convey meaning through a transformed speech signal. Whistled speech is an example that, unlike its spoken modality, does not require vocal folds. Instead it relies on the flow of compressed air trapped in the anterior cavity by manipulating lips, incisors, or fingers for its production (see [5] for a comprehensive review of whistled languages). Whistled speech is based on a spoken language, however, it expressed through a different modality [6]. Unlike whispering or shouting, the transformation of the speech signal is much more extreme, causing phonetic details to be smoothed or erased entirely [7]. The anterior cavity is the resonator that determines the intensity and high frequencies that are characteristic of whistling. Based on these differences with its spoken modality, whistling speech a good candidate for further study, however, additional work is needed to identify which acoustic domains and dimensions are useful for automatic speech processing.

The purpose of the current study was to identify acoustic features that were capable of distinguishing modal and whistled speech modalities and efficient for modeling. To do so, two experiments were designed. The first experiment relied on statistical analyses to examine a feature space that we hypothesized would best capture the acoustic characteristics of whistled speech. Based on any significant differences between the speech modalities, the second experiment selected the most relevant acoustic features to train machine learning models and tested them on a separate dataset. By reporting our findings, our goal was to add to the discussion surrounding the acoustic characteristics of whistled speech and how derived features played a role on automatic speech modality classification performance.

## 2 Related work

There are approximately 80 world populations that have adapted their local spoken language to a whistled modality [5]. In general, whistled speech can be described as having a strong, clear signal that is used to communicate over long-distances [8]. By not being forced to vibrate, the vocal folds act as a conduit for air flow, which amplifies the speech signal, but also limits the presence of phonemes, thus reducing the frequency space, i.e., higher order formants. More specifically, whistling transposes speech signal frequencies to within 1-4 kHz [5], a range humans are sensitive to and associate with speech [9]. Signal intensity can vary between 75 dB and 120 dB, which, again, can help propagate speech signals over large distances with minimal distortion, e.g., natural obstacles, noise. While there are tonal whistled languages, the current study focused on the non-tonal Silbo Gomero whistled form of Spanish spoken on La Gomera in the Canaries. This decision was based on the volume of available recordings and their use in a previous work involving automatic speech recognition [10].

One approach to improving our understanding of the differences between modal and whistled speech is to examine the dynamic temporal and spectral modulations underpinning the signals. Numerous studies have focused on unpacking modulation

77 representations of speech signals in the context of assessing intelligibility [11, 12]  
 78 and voice pathology [13, 14], shouted speech [15], and more recently laughter  
 79 [16, 17]. Unlike time-frequency representations, e.g., spectrograms, spectrotempo-  
 80 ral modulation (STM) representations unpack signals in terms of joint spectral and  
 81 temporal modulations, thus providing a combined and more comprehensive charac-  
 82 terization of the sound. Moreover, previous work has applied STM analysis to topics  
 83 such as birdsong [18] and communication in reverberant environments [19], suggest-  
 84 ing its usefulness for further study, given whistled speech is generally used for long  
 85 distance communication in natural settings. To the best of our knowledge, a STM  
 86 analytical method has not yet been used to evaluate the acoustic characteristics of  
 87 whistled speech, although its implementation could offer considerable advantages.  
 88 Given the characteristics of whistled speech as outlined above, our goal was to  
 89 examine whether STM-features were sensitive to differences between these speech  
 90 modalities.

91 Although limited, previous work has relied on STM-features for automatic  
 92 speech processing tasks, including speaker verification in the contexts of reverberant  
 93 environments [19] and voice-spoofing [20]. Typically, models are trained on Mel-  
 94 Frequency Cepstral Coefficients (MFCC), which scale the frequency components in  
 95 audio recordings to the Mel-scale, a domain that is better suited to model the human  
 96 auditory perception. While they are effective at transforming the spectral information  
 97 encoded in an audio recording into compressed, decorrelated vector, it is plausible  
 98 that nuances associated with non-typical speech signals are not represented in the  
 99 vector. As [10] and others suggest, models are then trained on constrained informa-  
 100 tion and, in turn, limiting their effectiveness. As evidence has shown machines and  
 101 humans model speech differently [21], our goal was to examine which STM-features  
 102 distinguished modal and whistled speech and whether this information was effective  
 103 when training automatic speech modality classification models.

## 104 3 Methods

### 105 3.1 Corpus

106 Table 1 describes the speech stimuli used in the experiments. The corpus was divided  
 107 into Train and Test datasets, each containing modal and whistled speech recordings  
 108 with no overlap between them. As described in further detail below, the Train dataset  
 109 was used in Experiments 1 and 2, whereas the Test dataset was used solely to test  
 110 models developed in Experiment 2.

**Table 1:** Description of the corpus

	Train		Test	
	Modal	Whistle	Modal	Whistle
Samples (#)	436	463	180	44
Duration (s)	$3.8 \pm 0.9$	$7.2 \pm 2.9$	$5.36 \pm 1.59$	$1.89 \pm 0.76$
Corpus	Common Voice	[10]	Common Voice	[7, 22, 23]

111 The Train dataset contained 463 whistled speech recordings produced by four  
 112 male whistlers (see [10] for corpus details). The Test dataset contained 44 whis-  
 113 tled speech recordings produced by a different group of 4 male speakers (see the  
 114 Supplementary Materials of [7, 22, 23] for more information). Following similar pro-  
 115 cedures described in [10], the Train and Test modal speech recordings were obtained  
 116 from the Mozilla Common Voice’s Spanish Common Voice Corpus 4<sup>1</sup> data set.  
 117 Although Common Voice provides speaker sex information, speaker identifications  
 118 were anonymised. Modal Spanish speech recordings in the Train dataset (N = 436)  
 119 were the same as those used in [10], while the Test dataset (N = 180) contained  
 120 randomly selected Spanish speech recordings made by male speakers. Although the  
 121 Modal:Whistle speech ratio in the Train dataset was approximately 1:1, the decision  
 122 to have a 4:1 ratio for the Test dataset was based on our interest in developing a  
 123 dataset that reflected the natural unbalance between speech modalities.

124 Prior to data processing, all stimuli were down-sampled to 8 kHz and normal-  
 125 ized such that the maximal amplitude of each recording was adjusted to a target of  
 126 100% of the signal dynamic. Finally, to ensure MPS representations had the same  
 127 dimensions, all recordings were zero-padded (see [24]).

### 128 **3.2 Spectrotemporal modulation features**

129 The spectrotemporal modulation domain can be characterized in terms of the Mod-  
 130 ulation Power Spectrum (MPS). Previous work [11, 25, 26] has defined MPS as a  
 131 two-dimensional Fourier transform of the time-frequency representation of an audio  
 132 signal. Equation 1 provides the formal definition of the MPS, where  $s$  and  $r$  are spec-  
 133 tral and temporal modulations, respectively, and  $Y(t, f)$  is the amplitude extracted  
 134 from the Fourier transform:

$$\text{MPS}(s, r) = \int \int Y(t, f) e^{-2\pi i s f} e^{-2\pi i r t} df dt \quad (1)$$

135 Similar procedures described in [14] and [16] were used to obtain MPS repre-  
 136 sentations of the speech recordings. All processing was done in MATLAB 2021a  
 137 (MathWorks Inc, USA) and based on adaptations to scripts described in [12]. Time-  
 138 frequency representations were obtained using a gammatone filter bank summation  
 139 method (128 full-width half-maximum Gaussians with center frequencies logarith-  
 140 mically spanning the frequency domain). Hilbert transforms were then used to extract  
 141 the analytical amplitudes from these filter outputs. The *fft2* MATLAB function  
 142 transformed the time-frequency representations into the modulation domain, where  
 143 temporal and spectral modulations were measured in increments of 0.05 Hz and 0.22  
 144 cycles per octave ( $c/o$ ), respectively. The amplitudes of positive and negative tempo-  
 145 ral and spectral modulations were averaged ( $A$ ) and transformed into decibels, i.e.,  
 146  $20 * \log_{10}(A/Ref)$  where  $Ref$  was set to the peak value of 16-bit audio (32,678).

147 Temporal and spectral modulation boundaries ranged from 0 to 20 Hz and 0 to 3  
 148 ( $c/o$ ), respectively. The former was selected due to findings reported in [15], which  
 149 associated  $< 20$  Hz temporal modulations with communicative meaning, while the  
 150 latter was identified as speech perception threshold in [11, 12]. Thus, for each speech

---

<sup>1</sup><https://commonvoice.mozilla.org/>

151 recording we obtained a 280-feature MPS representation, i.e., a 2-D matrix of ampli-  
 152 tudes associated with temporal modulations ( $N = 20$ ) by spectral modulation ( $N =$   
 153 14).

## 154 4 Experiments

### 155 4.1 Experiment 1

#### 156 4.1.1 Experimental setup

157 The goal of Experiment 1 was exploratory, aiming to identify and validate whether  
 158 spectrotemporal modulation (STM) features were useful when distinguishing modal  
 159 and whistled speech modalities. Developed as a two-step experiment, the Train  
 160 dataset was split into two parts, where 80% ( $N = 718$ ) was used to extract STM-  
 161 features and model speech modalities, henceforth the *exploratory* dataset, and the  
 162 remaining 20% ( $N = 181$ ) was used for validation, henceforth *validation* dataset.

163 First, a linear mixed-effects model (LMEM) was applied to the exploratory  
 164 dataset (*lm* from the *lme4* R-package). Model 2 set the spectrotemporal modulation  
 165 amplitude (Amp) as the dependent variable and Modality (modal, whistled) as a fixed  
 166 effect. To examine any differences between speech modalities and STM-features, our  
 167 model included the interaction between Modality and Rate, which expressed a coordi-  
 168 nate on the 2-D spectrotemporal modulation feature space, e.g., coordinate (2, 2)  
 169 represented 1 Hz temporal modulations and 0.44 c/o spectral modulations. As our  
 170 model only included Rate in interaction with Modality, it was treated as a nested  
 171 factor. Following [15], Bonferroni correction was applied to correct for multiple  
 172 comparisons ( $N = 50,679$ ), where  $\alpha = 0.001$ . Estimated marginal means (*emmeans*  
 173 R-package) were used to carry out pairwise comparisons.

$$174 \quad \text{lm}(\text{Amp} \sim \text{Type} + \text{Type}:\text{Rate}, \text{data} = \text{data}) \quad (2)$$

174

175

176 Based on these results (see Section 4.1.2), adjacent STM-features that were iden-  
 177 tified as significantly relevant for distinguishing speech modalities were bounded into  
 178 separate regions and averaged. Thus, for each recording in the validation dataset,  
 179 a mean amplitude was calculated for each region. A linear mixed-effects model  
 180 (Model 3) was designed to evaluate any differences between regions across speech  
 181 modalities, where the Mean amplitude of the region was set as a dependent vari-  
 182 able and Modality (modal, whistled) was set as a fixed factor. To examine any  
 183 differences between speech modalities across regions, our model included the inter-  
 184 action between Modality and Region, which was treated as a nested factor. Similarly,  
 185 post-hoc Bonferroni-adjusted t-tests were carried out using the *emmeans* R-package.

$$186 \quad \text{lm}(\text{Mean} \sim \text{Modality} + \text{Modality}:\text{Region}, \text{data} = \text{data}) \quad (3)$$

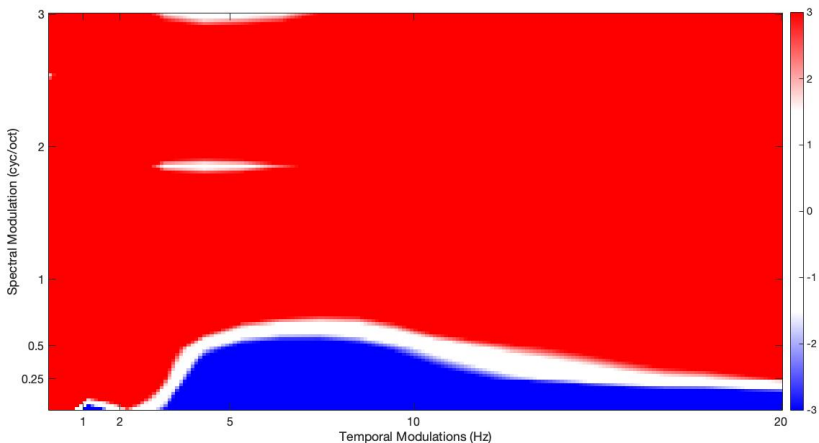
186

187

### 4.1.2 Results

Model 2 was statistically significant  $F_{559,200480} = 3995$ ,  $SE = 2.8$ ,  $p < 0.001$  (Intercept:  $\beta = 28.3$ ,  $SE = 0.15$ ,  $t = 191.96$ ,  $p < 0.001$ ) with a strong fit to the data (adjusted  $R^2 = 0.92$ ). We observed a significant fixed effect on Modality  $\beta = 3.4$ ,  $SE = 0.21$ ,  $t = 16.56$ ,  $p < 0.001$ . Posthoc pairwise tests showed Whistled speech ( $-16.74 \pm 0.01$  dB) had a significantly increased spectrotemporal modulation energy in comparison to Modal speech ( $-18.91 \pm 0.01$  dB),  $p < 0.001$ .

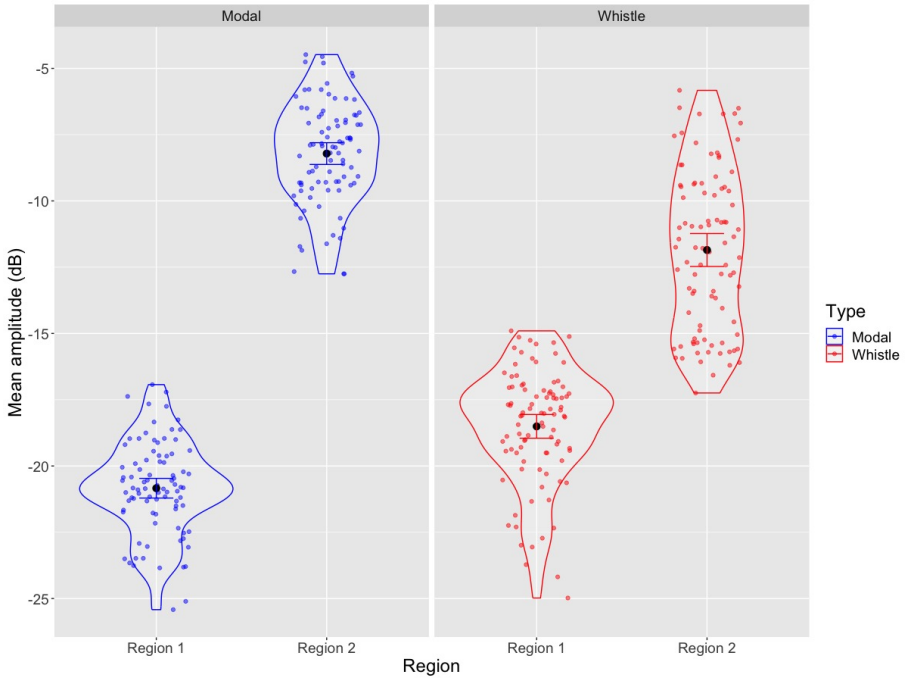
Figure 1 illustrates the  $t$ -values based on pairwise comparisons between Modal and Whistled speech across spectrotemporal modulations. In comparison to Modal speech, Whistled speech had a much more expansive, but nuanced, spectrotemporal modulation profile, largely operating at spectral modulations above 1 cyc/oct and across the temporal modulation domain (0-20 Hz, red in Fig. 1), henceforth *Region 1*. Alternatively, Modal speech exhibited a relatively uniform increase in temporal modulations ranging from 4-20 Hz and spectral modulations ranging from 0 to 0.22-0.44 cyc/oct (blue in Fig. 1), henceforth *Region 2*. Region boundaries were set at  $t > 3$  (Region 1) and  $t < -3$ , respectively. Overall, the 280 STM-feature space was reduced to 248 features when combining Region 1 ( $N = 217$ ) and Region 2 ( $N = 31$ ).



**Fig. 1:** Color map representing the posthoc pairwise comparisons ( $t$ -value) between modal and whistled speech across the spectrotemporal modulation domain (temporal modulations: 0-20 Hz; spectral modulations: 0-3 cyc/oct). Region 1 (red) and Region 2 (blue) represent all spectrotemporal modulations where  $t$ -values were  $t > 3$  and  $t < -3$ , respectively.

Model 3 was statistically significant  $F_{3,358} = 583$ ,  $SE = 2.28$ ,  $p < 0.001$  (Intercept:  $\beta = -20.84$ ,  $SE = 0.24$ ,  $t = -85.66$ ,  $p < 0.001$ ) with a strong fit to the data (adjusted  $R^2 = 0.83$ ). Our statistical analysis revealed a significant fixed effect on Modality  $\beta = 2.33$ ,  $SE = 0.34$ ,  $t = 6.87$ ,  $p < 0.001$ . Post-hoc pairwise comparisons showed

209 differences between Modal and Whistled speech for Region 1 (Modal:  $-20.84 \pm 0.24$   
 210 dB; Whistled:  $-18.51 \pm 0.24$  dB), and Region 2 (Modal:  $-8.21 \pm 0.24$  dB; Whistled:  
 211  $-11.85 \pm 0.24$  dB),  $p < 0.001$ . Figure 2 illustrates the distribution and means across  
 212 Modality and Region factors.



**Fig. 2:** Mean amplitude (dB) for regions across speech modalities. Significant differences were observed between speech modalities for Region 1 and Region 2,  $p < 0.001$ .

## 213 4.2 Experiment 2

### 214 4.2.1 Experimental procedure

215 The goal of Experiment 2 was to examine whether STM-features were efficient  
 216 for modeling modal and whistled speech. To do so, an automatic speech modality  
 217 classification task was developed and performed by two separate models. The STM  
 218 model used features identified in Experiment 1. To compare its performance, a second  
 219 model was developed that was trained on traditional Mel-Frequency Cepstral coefficients  
 220 (MFCC). To obtain the latter, we extracted 13 MFCCs and the first and second  
 221 derivatives resulting in 39-feature vector (window: 25 ms; step-size: 10 s) using the  
 222 *librosa* Python package. Each model was fitted using the Python scikit-learn library.  
 223 A 4-fold cross-validation design was performed on the Train dataset, whereupon the  
 224 models were tested on the Test dataset. A Support Vector Machines (SVM) model

225 was selected, as our goal was to identify the best margin between modal and whis-  
 226 tled speech modalities. Using hyperparameterization, both models were set with 128  
 227 coefficients and linear kernels.

## 228 4.2.2 Results

229 Table 2 illustrates the results of our STM and MFCC-feature models. Both mod-  
 230 els achieved high performances when distinguishing speech modalities, however,  
 231 the MFCC-feature model enhanced performance slightly (weighted average: 94%) in  
 232 comparison to the STM-feature model (92%).

**Table 2:** Results of Spectrotemporal Modulation (STM) and Mel-Frequency Cepstral Coefficients (MFCC) feature models. Precision, Recall, and F1-scores are reported across speech modalities and weighted averages (WA).

Model	Modality	Precision	Recall	F1
STM	Modal	0.98	0.92	0.95
	Whistle	0.74	0.91	0.82
	WA	0.93	0.92	0.92
MFCC	Modal	0.98	0.94	0.96
	Whistle	0.80	0.91	0.85
	WA	0.94	0.94	0.94

## 233 5 Discussion

234 The overall goal of these works was to examine whether spectrotemporal modula-  
 235 tion (STM) features were capable and efficient at distinguishing modal and whistled  
 236 speech. Our Experiment 1 findings revealed two regions of the STM-feature space as  
 237 characterizing differences between speech modalities. Based on these observations,  
 238 Experiment 2 was designed to compare automatic speech modality classification  
 239 performance by models trained on these specific STM-features and on traditional  
 240 cepstral coefficients. In general our findings suggest STM-features are effective at  
 241 distinguishing speech modalities. Moreover, the results of Experiment 2 suggest that  
 242 models trained on STM-features perform similarly to those trained on traditional  
 243 acoustic metrics, despite the limited number of features required.

244 Our Experiment 1 results revealed two regions of spectrotemporal modulation  
 245 (STM) features where modal and whistled speech differed. On one hand, whistled  
 246 speech exhibited an increase in modulations in Region 1 in comparison to modal  
 247 speech. In general, this region can be characterized as spectrotemporal modulations  
 248 above 1 cyc/oct. As described in [11], higher spectral modulations are associated  
 249 with pitch, whereas lower spectral modulations represent formants and higher order  
 250 harmonics. This observation is consistent with previous acoustic characterizations  
 251 of whistling, where the absence of vibrating vocal folds reduces the production of

252 phonemes and, in turn, the frequency space [5]. It is important to note that this dif-  
253 ference between speech modalities occurs mainly in the spectral modulation range  
254 associated with formants, where modal speech may be richer in comparison to whis-  
255 tled speech. On the other hand, in Region 2 (temporal modulations: ranging from 4  
256 to 20 Hz; spectral modulations: ranging from 0 to approximately 0.5 cyc/oct), we  
257 observed that modal speech exhibited an increase in energy in comparison to whis-  
258 tled speech. It is in this region that phonemes exist, as slow spectral modulations  
259 with temporal modulations in the 4-8 Hz and 8-16 Hz ranges have been associated  
260 with syllabic rhythm and articulation, respectively (see [27]). Numerous studies have  
261 attributed this region as critical for speech comprehensibility (see [11, 12], among  
262 others). These findings support previous work suggesting certain phonetic details  
263 may be lost during whistled speech production [7].

264 The results of Experiment 2 suggest both STM- and MFCC-feature models  
265 performed similarly. Although the MFCC-feature model slightly outperformed the  
266 STM-feature model, two major takeaways were observed. First, the size of the fea-  
267 ture spaces were not the same. While the dimensionality of the STM-feature vector  
268 was 248, the average number of MFCC-features was 21801 (vector dimensions: 39  
269 by  $559.1 \pm 228.0$ ). Thus, despite being nearly 0.01 the scale of the MFCC-feature  
270 vectors, the STM-feature vectors were quite sensitive to these differences in speech  
271 modalities. Although STM-feature vector were more compact in size, it is difficult  
272 to attribute the efficiency of their use, given the pre-processing required in compar-  
273 ison to the ease in which MFCCs can be extracted from audio recordings. Second,  
274 although the spectrotemporal modulation feature space is generally complex and dif-  
275 ficult to unpack, our findings are consistent with previous acoustic characterizations  
276 of whistled speech [5, 7], making it easier to interpret in the context of automatic  
277 speech processing. Alternatively, MFCCs provide a broadband description of the  
278 spectral information encoded in audio recordings, making them useful for automatic  
279 speech recognition tasks, i.e., speaker or speech emotion recognition, but not nec-  
280 essarily when connecting model efficiency with the acoustic information used for  
281 training.

282 Although our findings suggest spectrotemporal modulation metrics are sensi-  
283 tive to differences between modal and whistled speech, the evaluation of automatic  
284 speech modality classification performance could be improved in a number of ways.  
285 First we relied on Support Vector Machines, when other models, e.g., Decision Trees,  
286 Random Forests, may yield different results. Second, given how uncommon whistled  
287 speech is and volume of available recordings, we were unable to validate our models  
288 on a large dataset. A dataset of diverse speakers performing both modal and whis-  
289 tled speech is required. Of particular interest would be to obtain a dataset where each  
290 speaker performs the same work, phrase, or sentence in both modal and whistled  
291 speech. Finally it is difficult to evaluate the robustness of our method on another or  
292 mixed-language dataset, however, we hypothesize that similar STM-feature regions  
293 might emerge, given their acoustic characteristics match those previously reported in  
294 whistled speech [5]. Our future work aims to create such a dataset in order to per-  
295 form not only phonetic analysis, but also examine performance of different automatic  
296 speech processes, e.g., speaker and language recognition.

## 6 Conclusion

This study investigated whether spectrotemporal modulation metrics were sensitive to distinguishing modal and whistled speech. The following lists take-away messages from the study: (1) statistical analysis was used to identify spectrotemporal modulation features that distinguished modal and whistled speech; (2) these features formed two separate regions and significant differences between speech modalities were observed; and (3) STM-feature and MFCC-feature models both exhibited high automatic speech modality classification performance. The findings support previously reported acoustic characterizations of whistled speech and suggest spectrotemporal modulation features are useful for distinguishing whistled and modal speech in automatic speech processing applications.

**Funding.** The authors did not receive support from any organization for the submitted work.

## References

- [1] Mohammed, M., Abdulkareem, K., Mostafa, S., Maashi, M., Zapirain, B., Oleagordia, I., Al-Dhief, F., Alhakami, H.: Voice pathology detection and classification using convolutional neural network model. *Applied Sciences* **10**, 3723 (2020). <https://doi.org/10.3390/app10113723>
- [2] Alhussein, M., Muhammad, G.: Voice pathology detection using deep learning on mobile healthcare framework. *IEEE Access* **6**, 41034–41041 (2018). <https://doi.org/10.1109/ACCESS.2018.2856238>
- [3] Vestman, V., Gowda, D., Sahidullah, M., Alku, P., Kinnunen, T.: Speaker recognition from whispered speech: A tutorial survey and an application of time-varying linear prediction. *Speech Communication* **99**, 62–79 (2018). <https://doi.org/10.1016/j.specom.2018.02.009>
- [4] Pohjalainen, J., Raitio, T., Yrttiaho, S., Alku, P.: Detection of shouted speech in noise: Human and machine. *The Journal of the Acoustical Society of America* **133**(4), 2377–2389 (2013). <https://doi.org/10.1121/1.4794394>
- [5] Meyer, J.: *Whistled Languages: A Worldwide Inquiry on Human Whistled Speech*. Springer (2015). <https://doi.org/10.1007/978-3-662-45837-2>
- [6] Busnel, R.-G.: Recherches experimentales sur la langue sifflee de kuskoy (experimental research on the whistling language of kuskoy). *Revue de phonétique appliquée* (1970)
- [7] Meyer, J., Magnasco, M.O., Reiss, D.: The relevance of human whistled languages for the analysis and decoding of dolphin communication. *Frontiers in Psychology* **12** (2021). <https://doi.org/10.3389/fpsyg.2021.689501>

- 333 [8] Classe, A.: The whistled language of la gomera. *Scientific American* **196**(4),  
334 111–120 (1957). <https://doi.org/10.1038/scientificamerican0457-111>
- 335 [9] Moore, B., Glasberg, B.: A revision of zwicker’s loudness model. *Acta Acustica*  
336 *united with Acustica* **82**, 335–345 (1996)
- 337 [10] Jakubiak, A.: Whistle-to-text: Automatic recognition of the Silbo Gomero whis-  
338 tled language. In: *Proc. INTERSPEECH 2023*, pp. 3402–3406 (2023). <https://doi.org/10.21437/Interspeech.2023-989>
- 339
- 340 [11] Elliott, T.M., Theunissen, F.E.: The modulation transfer function for speech  
341 intelligibility. *PLOS Computational Biology* **5**(3), 1–14 (2009). <https://doi.org/10.1371/journal.pcbi.1000302>
- 342
- 343 [12] Flinker, A., Doyle, W., Mehta, A., Devinsky, O., Poeppel, D.: Spec-  
344 trotemporal modulation provides a unifying framework for auditory cortical  
345 asymmetries. *Nature Human Behaviour* **3** (2019). <https://doi.org/10.1038/s41562-019-0548-z>
- 346
- 347 [13] Moro-Velázquez, L., Gómez-García, J., Godino Llorente, J., Andrade-Miranda,  
348 G.: Modulation spectra morphological parameters: A new method to assess  
349 voice pathologies according to the grbas scale. *BioMed Research International*  
350 **2015** (2015). <https://doi.org/10.1155/2015/259239>
- 351 [14] Marczyk, A., O’Brien, B., Tremblay, P., Woisard, V., Ghio, A.: Correlates of  
352 vowel clarity in the spectrotemporal modulation domain: Application to speech  
353 impairment evaluation. *The Journal of the Acoustical Society of America* **152**,  
354 2675–2691 (2022). <https://doi.org/10.1121/10.0015024>
- 355 [15] Arnal, L., Flinker, A., Kleinschmidt, A., Giraud, A.-L., Poeppel, D.: Human  
356 screams occupy a privileged niche in the communication soundscape. *Current*  
357 *biology : CB* **25** (2015). <https://doi.org/10.1016/j.cub.2015.06.043>
- 358 [16] Mazzocconi, C., O’Brien, B., Bodur, K., Fourtassi, A.: Do children laugh  
359 like their parents? conversational laughter mimicry occurrence and acoustic  
360 alignment in middle-childhood. *Journal of Nonverbal Behavior* (2025). <https://doi.org/10.1007/s10919-025-00478-z>
- 361
- 362 [17] Ludusan, B., Wagner, P.: An Evaluation of Manual and Semi-Automatic Laugh-  
363 ter Annotation. In: *Proc. Interspeech 2020*, pp. 621–625 (2020). <https://doi.org/10.21437/Interspeech.2020-2521>
- 364
- 365 [18] Woolley, S.M.N., Fremouw, T.E., Hsu, A., Theunissen, F.E.: Tuning for spectro-  
366 temporal modulations as a mechanism for auditory discrimination of natural  
367 sounds. *Nature Neuroscience* **8**(10), 1371–1379 (2005). <https://doi.org/10.1038/nn1536>
- 368

- 369 [19] Avila, A., Fraga, F., Sarria-Paja, M., Falk, T.: Investigating the use of mod-  
370 ulation spectral features within an i-vector framework for far-field auto-  
371 matic speaker verification, pp. 1–5 (2014). [https://doi.org/10.1109/ITS.2014.](https://doi.org/10.1109/ITS.2014.6948012)  
372 [6948012](https://doi.org/10.1109/ITS.2014.6948012)
- 373 [20] Suthokumar, G., Sethu, V., Wijenayake, C., Ambikairajah, E.: Modulation  
374 dynamic features for the detection of replay attacks, pp. 691–695 (2018).  
375 <https://doi.org/10.21437/Interspeech.2018-1846>
- 376 [21] Park, S.J., Yeung, G., Vesselinova, N., Kreiman, J., Keating, P., Alwan,  
377 A.: Towards understanding speaker discrimination abilities in humans and  
378 machines for text-independent short utterances of different speech styles. *The*  
379 *Journal of the Acoustical Society of America* **144**, 375–386 (2018). <https://doi.org/10.1121/1.5045323>  
380
- 381 [22] Meyer, J.: Whistled Speech as a tool for shepherds (demonstration in Whis-  
382 tled Spanish at the ‘House of shepherds’, Ecrins National Park, France). *The*  
383 *World Whistles/Maison Du Berger/ Adaparle Project (UGA)*. [https://vimeo.](https://vimeo.com/327303614)  
384 [com/327303614](https://vimeo.com/327303614)
- 385 [23] Meyer, J., Ridouane, R.: Parler en sifflant dans les montagnes de l’Atlas : dia-  
386 logue sifflé à longue distance en langue berbère tachelhit. (Speaking while  
387 whistling in the Atlas : whistled dialogs in Berber Tashlhiyt). A co-production  
388 CNRS - GIPSA-Lab/CNRS – LPP. <https://vimeo.com/361557210>
- 389 [24] Venezia, J.H., Hickok, G., Richards, V.M.: Auditory “bubbles”: Efficient classi-  
390 fication of the spectrotemporal modulations essential for speech intelligibility.  
391 *The Journal of the Acoustical Society of America* **140**(2), 1072–1088 (2016).  
392 <https://doi.org/10.1121/1.4960544>
- 393 [25] Singh, N.C., Theunissen, F.E.: Modulation spectra of natural sounds and etho-  
394 logical theories of auditory processing. *The Journal of the Acoustical Society*  
395 *of America* **114**(6), 3394–3411 (2003). <https://doi.org/10.1121/1.1624067>
- 396 [26] Thoret, E., Caramiaux, B., Depalle, P., Mcadams, S.: Learning metrics  
397 on spectrotemporal modulations reveals the perception of musical instru-  
398 ment timbre. *Nature Human Behaviour* **5** (2021). [https://doi.org/10.1038/](https://doi.org/10.1038/s41562-020-00987-5)  
399 [s41562-020-00987-5](https://doi.org/10.1038/s41562-020-00987-5)
- 400 [27] Giraud, A.-L., Poeppel, D.: Cortical oscillations and speech processing: Emerg-  
401 ing computational principles and operations. *Nature neuroscience* **15**, 511–7  
402 (2012). <https://doi.org/10.1038/nn.3063>