

The closer the better: The role of formant positions in timbre similarity perception and timbre classification

Christoph Reuter^{1†}, Isabella Czedik-Eysenberg¹, Saleh Siddiq¹ and Michael Oehler²

¹ Musicological Department, University of Vienna, Vienna, Austria

² Musicological Department, University of Osnabrück, Osnabrück, Germany

†Corresponding author email: christoph.reuter@univie.ac.at



Background

It is a well-known fact in timbre research that the timbral character of a musical instrument as a whole cannot be sufficiently represented by the sound attributes of a single tone (Stumpf, 1926, p. 393; Siedenburg, Jones-Mollerup, & McAdams, 2016, p. 15). Only few concepts for a wider timbral description of an instrument through its entire dynamics and pitch range do exist. Beside Mel Frequency Cepstral Coefficients (MFCCs, e.g. Loughran, Walker, O'Neill, & O'Farrell, 2008) and Modulation Power Spectrum (MPS, e.g. Elliott, Hamilton, & Theunissen, 2013) formant areas were shown to be very useful for a more comprehensive description of musical instruments' timbres (since 1926, e.g. Stumpf, 1926; Schumann, 1929; Mertens, 1975; Meyer, 2009 etc.).

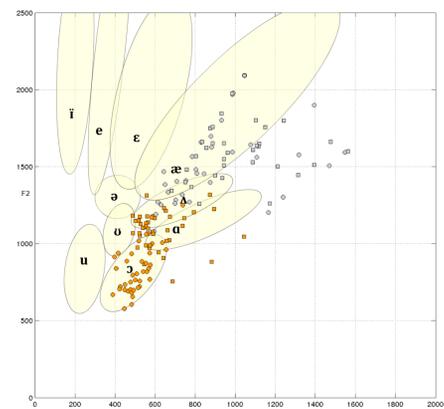


Figure 1. Formant map with the sounds of bassoon (orange) and oboe (grey) in all achievable pitches in *ff* and *pp* (Reuter, Czedik-Eysenberg, Siddiq, & Oehler, 2017).

Therefore, based on the formant areas of 586 wind instrument timbres in all reachable pitches and two dynamic levels we compiled a two-dimensional formant map (X-axis: formant 1, Y-axis: formant 2). The calculations were done using Praat (Boersma, & Weenink, 2013). In this map, the musical instruments (represented as point clouds in different colors) emerged to be clustered together by instrument, dynamics and register (Figure 1). The calculated formant positions correspond to the descriptions found in the literature (e.g. Schumann, 1929; Mertens, 1975; Reuter, 1996; Meyer, 2009 etc.).

Objectives

The aim of this study is to test the possible advantages and suitability of the formant paradigm for musical timbre similarity prediction as well as for timbre classification by investigating the following questions:

- Do shorter distances between the formant positions of two sounds correlate with the impression of greater timbre similarity (and vice versa)?
- Is the numerical vector of formant positions a suitable basis for a precise computational classification of the involved instruments? (Besides, which additional timbre features can improve the classification result?)

Computational description of perceived timbre similarity: Formants and MFCCs in comparison

In a listening test, 22 participants rated the (dis)similarity of 40 loudness-adjusted timbre combinations including wind instruments like flute, oboe, clarinet, bassoon, trumpet, trombone, French horn and tuba on a scale between 1 and 8 (8 = maximum dissimilarity).

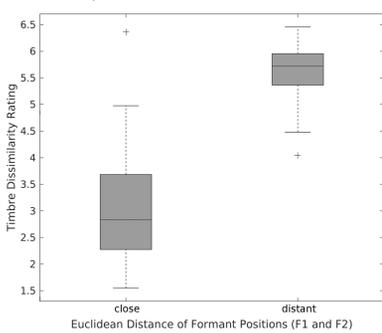


Figure 2. The distance of the formant positions (X-axis: close vs. distant) correlates strongly with ratings of perceived timbre similarity (Y-axis: 1 = very similar; 7 = very dissimilar).

Half of the stimuli contained sounds with extremely close/overlapping formant regions while the other half contained sounds with very distant formant regions. Instrument sounds were taken from the Vienna Symphonic Library (VSL) and adjusted to a matching ANSI-loudness level using the Genesis loudness toolbox in Matlab. The listeners' judgments were tested for correlations with the Euclidean distances between positions in the formant map. Indeed, a very significant correlation between the formant distances and the listeners' similarity judgments was found ($r = 0.759$, t-test with $p < 0.001$, 95% CI [-3.1381, -1.8960]).

Here, the first two formants (F1 and F2) have an almost equally strong linear relationship with the similarity scores and also have a very strong intercorrelation ($r = 0.9196$, $p < 0.0001$). In comparison, each of the first three MFCCs correlates weaker with the listeners' similarity scores, while the remaining MFCCs (4-13) show no significant correlations to the listeners' judgements at all. (Formants were calculated in Praat (Boersma & Weenink, 2013) and MFCCs in Matlab with the help of the MIRtoolbox (Lartillot & Toivainen, 2007).)

Timbre feature	r	p	Timbre feature	r	p
F1	0.7514	< 0.0001	F1 and F2	0.7591	< 0.0001
F2	0.7477	< 0.0001	MFCCs 1, 2, and 3	0.6949	< 0.0001
F3	0.4227	< 0.0001	MFCCs 1 and 2	0.6939	< 0.0001
MFCC1	0.6384	< 0.0001	F1, F2, and F3	0.6916	< 0.0001
MFCC2	0.5959	< 0.0001	MFCCs 1-13	0.6812	< 0.0001
MFCC3	0.3513	< 0.05			

Table 1. Correlation of individual formant positions and MFCCs with the perceived timbre similarity.

Table 2. Correlation of the Euclidean distances of formant (F1, F2, and F3) and MFCC combinations with perceived timbre similarity.

A similar result can be seen when considering formant and MFCC combinations. In a two- or three-dimensional formant or MFCC space, the Euclidean distances based on the formants 1 and 2 correlate most strongly with the perceived timbre similarity ($r = 0.759$, $p < 0.0001$, see Fig. 3), while the Euclidean distances based on the MFCCs 1-3 give a slightly weaker but comparable picture ($r = 0.695$, $p < 0.0001$, see Fig. 4).



Figure 3. Scatter plot of the Euclidean distances of the first two formants positions (X-axis) and the perceived timbre similarity (Y-axis) (red: close formant positions, blue: distant formant positions).

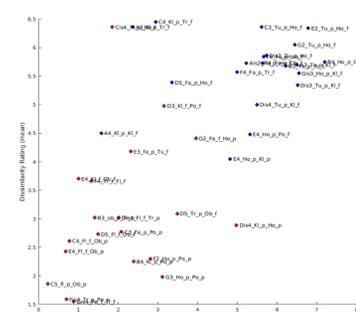


Figure 4. Scatter plot of the Euclidean distances of the first three MFCCs (X-axis) and the perceived timbre similarity (Y-axis) (red: close formant positions, blue: distant formant positions).

In a spider's web representation with the respective axes F1, F2, and F3 or MFCC1, MFCC2, MFCC3 the triangles formed between the axes overlap very closely in case of a perceived timbre similarity (Figure 5) while the distances of the triangle sides increase with increasing timbre similarity (Figure 6).

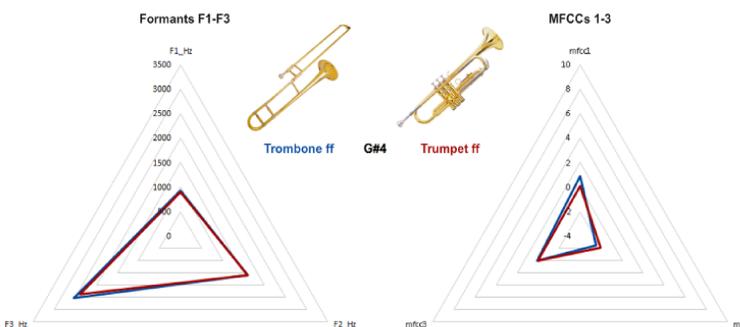


Figure 5. Spider's web representation with the respective axes F1, F2, F3 or MFCC1, MFCC2, MFCC3 of the pair with the most perceived timbre similarity (trombone *ff* and trumpet *ff* on G#4).

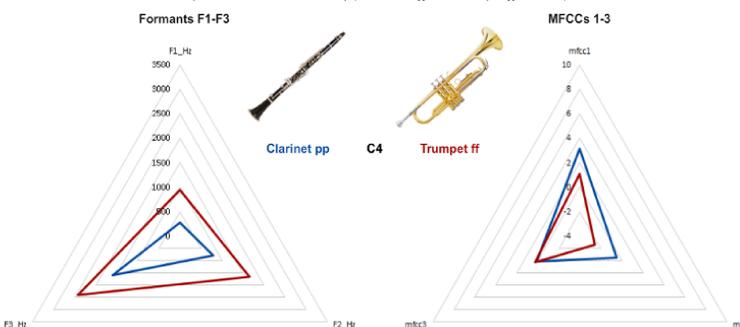


Figure 6. Spider's web representation with the respective axes F1, F2, F3 or MFCC1, MFCC2, MFCC3 of the pair with the least perceived timbre similarity (clarinet *pp* and trumpet *ff* on C4).

A comparison of these results indicates that formants and MFCCs are roughly similarly suitable for the determination of perceived timbre similarity. This is also suggested by the model of Darch et al. (2005), where formant areas can be derived on the basis of MFCCs.

Computational predictability of perceived timbre similarity: Formants and MFCCs in comparison

Regression models trained via machine learning (using 5-fold cross-validation) allow the prediction of perceived timbre similarities to a certain degree: Based on formants, the best result with a coefficient of determination (R^2) of 0.53 was found when using the absolute differences of F1 and F2 as features (Figure 7). Based on MFCCs a comparable result ($R^2 = 0.56$) can be obtained when using MFCC1 and MFCC2 (Figure 8).

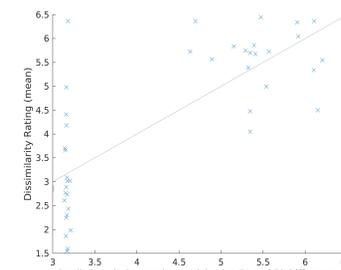


Figure 7. Prediction model based on F1 and F2 ($R^2 = 0.53$, RMSE = 1.08, MSE = 1.16, MAE = 0.86).

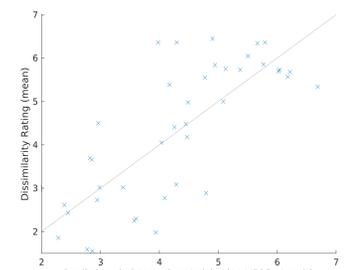


Figure 8. Prediction model based on MFCC1 and MFCC2 ($R^2 = 0.56$, RMSE = 1.05, MSE = 1.10, MAE = 0.81).

Timbre classification based on formants

With the help of machine classification methods, timbre classes can be assigned to musical instruments with a precision of 46.1% based on the first three formants (F1, F2, F3, cubic k-nearest-neighbour classification) (Figure 9).

True Class	Fagott	Flöte	Horn	Klarinette	Oboe	Posaune	Trompete	Tuba	True Positive Rate	False Negative Rate
Fagott	72%	1%	7%	1%	4%	2%		12%	72%	28%
Flöte	12%	42%	3%	14%	13%	4%	9%	3%	42%	58%
Horn	10%	1%	66%	3%	3%	8%	3%	6%	66%	34%
Klarinette	7%	18%	13%	24%	15%	8%	5%	10%	24%	76%
Oboe	3%	24%	1%	10%	34%	7%	19%	1%	34%	66%
Posaune	20%	6%	14%	3%	2%	35%	15%	6%	35%	65%
Trompete	8%	15%	2%	8%	14%	20%	31%	2%	31%	69%
Tuba	24%	3%	4%	1%	1%	5%	1%	59%	59%	41%

Figure 9. Confusion matrix based on F1, F2, and F3 (cubic KNN, classification precision of 46.1%).

When adding a set of complementary timbre features to the formant positions, the precision of the instrument recognition can be increased up to 84.6%. These additional timbre features are attack time, spectral flux, roughness, brightness, spectral entropy, maximum RMS value, spectral centroid and unpleasantness. The last descriptor is implemented based on the results in Reuter, Oehler, & Mühlhans (2013), the rest is calculated with the help of Matlab MIRtoolbox (Lartillot & Toivainen, 2007). The resulting confusion matrix shows plausible parallels to human perception, as common perceived instrument confusions are mirrored strikingly well in the matrix. Examples are the confusions between flute, clarinet and oboe, between trumpet and trombone as well as in between tuba and bassoon timbres (Figure 10).

True Class	Fagott	Flöte	Horn	Klarinette	Oboe	Posaune	Trompete	Tuba	True Positive Rate	False Negative Rate
Fagott	90%		2%	1%	1%	2%		2%	90%	10%
Flöte		79%		3%	9%	3%	7%		79%	21%
Horn	1%		88%	3%	3%	1%	4%		88%	12%
Klarinette		6%		83%	8%	1%			83%	15%
Oboe	3%	6%		1%	79%	4%	6%		79%	21%
Posaune	5%			2%	3%	80%	8%	3%	80%	20%
Trompete			5%		5%	8%	81%		81%	19%
Tuba	1%	1%		1%		4%		92%	92%	8%

Figure 10. Confusion matrix based on F1, F2, and F3 as well as on attack time, spectral flux, roughness, brightness, spectral entropy, maximum RMS value, spectral centroid, and unpleasantness (quadratic SVM, classification precision of 84.6%).

Conclusion

The timbre similarity of wind instruments corresponds well to the proximity of the first two formants. A similarity prediction based on formant positions yielded results comparable to using MFCCs. Constructing classification models, we showed that a formant-driven instrument identification is possible, but does not lead to particularly high recognition rates. However, adding a small number of additional timbre features yields an effective classification model that even corresponds to human perception. While formants, as physical descriptors of perceived timbre similarity, strongly correlate with each other, individual MFCCs are more linearly independent of each other. However, the advantages of formants over MFCCs are: (1) formants need only two values to compactly describe a distinctive audible and therefore intuitively accessible spectral content and (2) they provide a solid foundation through more than 100 years of research history. For these reasons, formants can be considered as useful alternatives to MFCCs when analysing (wind) instrument timbres in music information retrieval.

Acknowledgements. This project is supported by funds of the Österreichische Nationalbank (OeNB, Anniversary Fund, project number: 16473). Special thanks to the Vienna Symphonic Library for supporting our work with two free "Vienna Super Package" licenses.

Reference

Boersma, P., & Weenink, D. (2013). Praat: Doing Phonetics by Computer [Computer program]. Version 5.3.51, retrieved 2 June 2013 URL: <http://www.praat.org/>.
 Darch, J., Miller, B., Shao, X., Vasighi, S., & Tan, Q. (2005). Predicting Formant Frequencies from MFCC Vectors. In Proceedings of the 2005 IEEE International Conference on Acoustics, Speech, and Signal Processing, Volume 1 of V. (pp. 1941-1944). Pennsylvania, USA: Pennsylvania Convention Center/Marriott Hotel.
 Elliott, T., Hamilton, L., & Theunissen, F. (2013). Acoustic Structure of the Five Perceptual Dimensions of Timbre in Orchestral Instrument Tones. Journal of the Acoustical Society of America, 133(1), 389-404.
 Lartillot, O., & Toivainen, P. (2007). A Matlab Toolbox for Musical Feature Extraction from Audio. In Proceedings of the 10th International Conference on Digital Audio Effects (DAFx) (pp. 237-244). France, Bordeaux.
 Loughran, R., Walker, L., O'Neill, M., & O'Farrell, M. (2008). Musical Instrument Identification using Principal Component Analysis and Multi-layered Perceptrons. In Proceedings of the International Conference on Audio, Language and Image Processing (CALP 2008) (pp. 643-648). China: Shanghai University.
 Mertens, P.-H. (1975). Die Schumannschen Klangtemperamenten und ihre Bedeutung für die Übertragung von Sprache und Musik (The Schumann Laws of Timbre and their Meaning for the Transmission of Voice and Music). Frankfurt: Borchinsky.
 Meyer, J. (2009). Acoustics and the Performance of Music. Manual for Acousticians, Audio Engineers, Musicians, Architects and Musical Instrument Makers. New York: Springer.
 Reuter, C. (1996). Die auditive Diskrimination von Orchesterinstrumenten (Auditive Separability of Orchestral Instruments). Frankfurt: Lang.
 Reuter, C., Czedik-Eysenberg, I., Siddiq, S., & Oehler, M. (2017). Formanten als hilfreiche Timbre-Deskriptoren für die Darstellung von Blasinstrumententönen (Formants as helpful Timbre Descriptors for the Depiction of Wind Instruments Timbres). In Proceedings of the 43rd German Annual Conference on Acoustics 'Fortschritte der Akustik'. DAGA 2017 (pp. 190-193). Germany: University of Kiel.
 Reuter, C., Oehler, M., & Mühlhans, J. (2014). Physiological and acoustical correlates of unpleasant sounds. In Proceedings of the Joint Conference ICMI3/APCOMS (p. 9). South Korea: Yonsei University, Seoul.
 Schumann, K.E. (1929). Physik der Klangfarben (Physics of Timbres). (Habilitation Thesis). University of Berlin, Berlin, Germany.
 Siedenburg, K., Jones-Mollerup, K., & McAdams, S. (2016). Acoustic and Categorical Dissimilarity of Musical Timbre: Evidence from Asymmetries between Acoustic and Categorical Sounds. Frontiers in Psychology, 6:1977. doi:10.3389/fpsyg.2015.01161.
 Stumpf, C. (1926). Die Sprachlaute (The Sounds of Voice). Berlin: Springer.