



Audio Engineering Society

Convention Paper

Presented at the 126th Convention
2009 May 7–10 Munich, Germany

The papers at this Convention have been selected on the basis of a submitted abstract and extended precis that have been peer reviewed by at least two qualified anonymous reviewers. This convention paper has been reproduced from the author's advance manuscript, without editing, corrections, or consideration by the Review Board. The AES takes no responsibility for the contents. Additional papers may be obtained by sending request and remittance to Audio Engineering Society, 60 East 42nd Street, New York, New York 10165-2520, USA; also see www.aes.org. All rights reserved. Reproduction of this paper, or any portion thereof, is not permitted without direct permission from the Journal of the Audio Engineering Society.

Perceptually-motivated audio morphing: softness

Duncan Williams, Tim Brookes¹

¹ Institute of Sound Recording, University of Surrey, Guildford, Surrey, GU2 7XH, England

ABSTRACT

A system for morphing the softness and brightness of two sounds independently from their other perceptual or acoustic attributes was coded. The system is an extension of a previous one that morphed brightness only, that was based on the Spectral Modelling Synthesis additive/residual model. A Multidimensional Scaling analysis, of listener responses to paired comparisons of stimuli generated by the morpher, showed movement in three perceptually-orthogonal directions. These directions were labelled in a subsequent verbal elicitation experiment which found that the effects of the brightness and softness controls were perceived as intended. A Timbre Morpher, adjusting additional timbral attributes with perceptually-meaningful controls, can now be considered for further work.

1. INTRODUCTION

Audio morphing or sound hybridisation is a technique for creating a new sound with characteristics derived from a pair of existing sounds. The technique is commonly used for innovative sound design and digital effects processing [1]. Current morphing systems interpolate between the spectral, or spectral and temporal, features of each source sound to create a new feature set which can then steer additive resynthesis [2]. Our long-term goal is to develop a more intuitive system allowing perceptually-meaningful control over the feature set being morphed [3].

A pilot study established that a morph along a single perceptual dimension with spectral correlates was possible [4]. The current investigation aims to establish the viability of steering a morph across additional perceptual dimensions, such that the intended timbral attributes change but other timbral characteristics are preserved. A complete timbre morpher should be able to control attributes with spectral, temporal, or spectro-temporal acoustic correlates.

1.1. Selecting a second timbral attribute

Softness was chosen as the next attribute to be implemented in the morpher, following a review of published research into timbral attributes which

suggested that softness was correlated to both inharmonicity and attack time [5] [6] [7], [8], [9].

The literature review did not reveal the exact contribution of these correlates to perceived softness and so three separate listening tests were necessary. Firstly, a preliminary experiment was required in order to determine the contribution of each acoustic correlate to softness. Secondly, the main experiment was designed to establish whether or not the manipulation of these correlates would be perceived as orthogonal to other timbral attributes. Thirdly, a subsequent verbal elicitation experiment was specified to provide appropriate labels for the timbral movements perceived.

2. PERCEPTUAL TESTING

2.1. Preliminary experiment

A preliminary experiment was designed to reveal the contribution of each proposed acoustic correlate to perceived softness, and, if necessary, to derive a softness matrix modelling the relative contribution of those correlates. A stimulus set varying in both inharmonicity and attack time (the proposed correlates) was synthesised using Matlab mathematical modelling software. Stimuli were generated with a fixed fundamental of 200 Hz, normalized peak amplitude, and 4 additional partials, to ensure that upper harmonics would not clash within a critical band, which could otherwise have caused variations in the perceived roughness of the stimuli [16].

Inharmonicity values were measured using the inharmonicity coefficient specified by Rossing [10], shown in equation 1.

$$h[n] = \frac{2(f[n] - n * f[1])}{(\{n^2 - 1\} * n * f[1])}$$

(Equation 1. Inharmonicity coefficient
 h = inharmonicity coefficient
 n = partial number
 $f[n]$ = partial frequency)

To establish a suitable range of inharmonicity values for the main experiment, 10 sounds were generated, varying from totally harmonic to totally inharmonic in 10% increments. This stimulus set was listened to informally and, to the experimenter's ears, a difference of at least 20% was required in order to readily differentiate between the sounds. Therefore, when generating the stimuli for the main experiment, five values of inharmonicity were used: from 0% to 100% in 25% increments.

With regards to selection of attack time values, in keeping with research on percussion onset times [6, 11], an initial stimulus set comprising 39 sounds varying from instantaneous to 2 second attack times was generated and evaluated informally. To the experimenter's ears, differences of smaller than 25ms were difficult to discriminate between, and sounds with longer attack times than 100ms did not appear to be varying in perceived softness. Furthermore, an instantaneous attack was giving rise to a digital click on some playback systems. Attack time values for this experiment were therefore set in 25ms increments, up to 100ms, though the instantaneous attack value was adjusted to 5ms to avoid the occurrence of a digital click on some playback systems. All attack envelopes were linear.

The full stimulus set is summarised in Table 1.

Stimulus No.	Attack Time	Inharmonicity
1	5ms	100%
2	5ms	75%
3	5ms	50%
4	5ms	25%
5	5ms	0%
6	25ms	100%
7	25ms	75%
8	25ms	50%
9	25ms	25%
10	25ms	0%
11	50ms	100%
12	50ms	75%
13	50ms	50%
14	50ms	25%
15	50ms	0%
16	75ms	100%
17	75ms	75%
18	75ms	50%
19	75ms	25%
20	75ms	0%
21	100ms	100%
22	100ms	75%
23	100ms	50%
24	100ms	25%
25	100ms	0%

Table 1 Stimulus set for preliminary experiment

Moving from stimulus 1 to stimulus 7, to 13, to 19, to 25, gives increased attack time and decreased inharmonicity at each step and so, according to the published research mentioned above, should provide increased softness at each step, regardless of the relative contributions to softness of these two parameters. These five stimuli were therefore chosen as references.

In the listening test, each reference stimulus was presented 5 times: once together with all the 0% harmonic stimuli, once together with all the 25% harmonic stimuli, once with all the 50%, once with all the 75% and once with all the 100%. Thus, for each listener, the test consisted of 25 presentation groups, represented by the 25 rows of buttons on the user interface shown in Figure 1. The left-most button in each row triggered playback of a

reference stimulus; the other 5 buttons triggered playback of each of the comparison stimuli in one presentation group, ordered from left to right in order of increasing attack time. For each presentation group (row) the listener's task was to audition each stimulus in that group and to select the non-reference stimulus whose softness matched most closely that of the reference.

12 listeners, all with some experience of audio engineering, took part in the test. The row order was randomised.

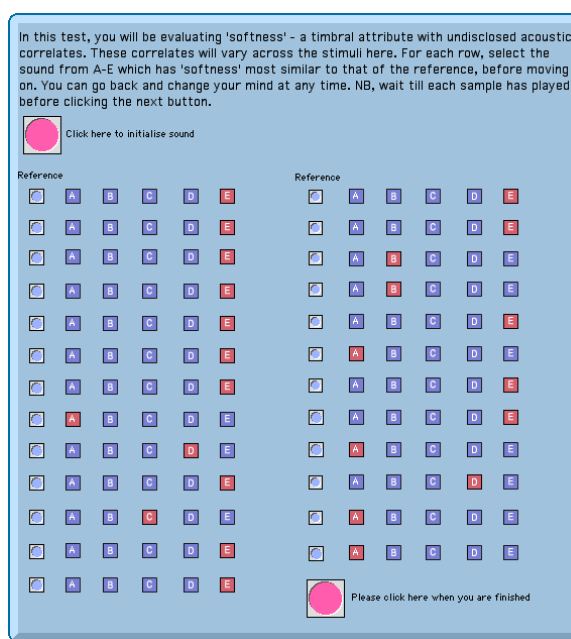


Figure 1 Listener interface of preliminary experiment, showing 25 stimulus presentation groups.

2.1.1. Results from preliminary experiment

Each listener's responses comprised a series of 25 attack time values. For each combination of reference stimulus and comparison inharmonicity (i.e. each presentation group or row), the listener's response revealed the attack time, for a tone of that particular inharmonicity, that they felt gave it a softness most close to that of the reference stimulus. Mean results across all listeners, for each

of the 25 presentation groups, are summarised in Table 2.

Reference Stimulus	Inharmonicity of Comparison Stimuli				
	0%	25%	50%	75%	100%
1. 5ms 100%	17.9	16.3	14.2	16.7	16.1
7. 25ms 75%	40	35.4	29.2	37.5	36.5
13. 50ms 50%	62.5	60.4	50	62.5	54.2
19. 75ms 25%	77.1	72.9	64.6	81.3	77.1
25. 100ms 0%	89.6	85.4	72.9	87.5	79.2

Table 2 Results from preliminary experiment. Numbers in body of table are means of matched stimulus attack times (ms).

These results show little variation in terms of mean matched attack time across the range of inharmonicity values. The full set of result data is plotted as means plus 95% confidence intervals in Figure 2, which can be interpreted as a set of 5 equal-softness contours. Each softness contour shows the combinations of inharmonicity and attack time which will lead to a softness equal to that of the reference stimulus identified in the key.

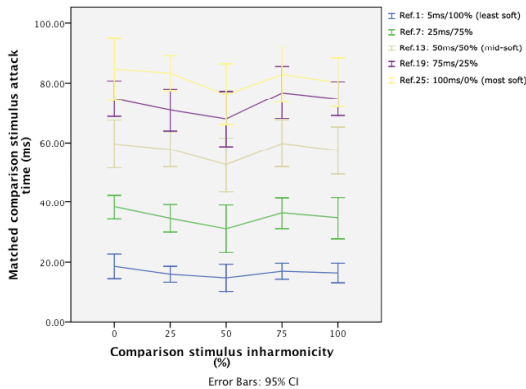


Figure 2 Listener responses plotted as means plus 95% confidence intervals. Data relating to each reference stimulus are connected to show an equal-softness contour.

The data relating to any softness level in Figure 2 can be fitted by a horizontal straight line. This

shows that stimulus inharmonicity makes no significant contribution to the perception of softness.

2.1.2. Specifying the softness morphing routine

In view of the results of the preliminary experiment, the softness morphing algorithm was designed to ignore inharmonicity and to manipulate attack time alone, as shown in Figure 3.

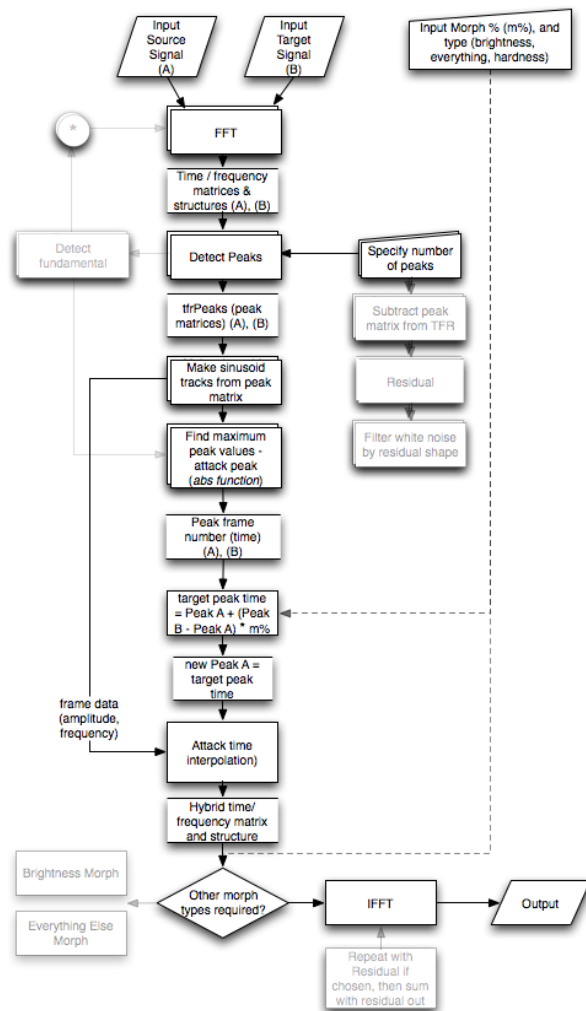


Figure 3 Signal flow of softness morphing algorithm. Lighter shaded blocks are optional if the filtered noise portion of the resynthesis engine is selected [12]

Within this algorithm a simplistic notion of attack time was adopted, whereby a threshold detection procedure finds the maximum peak level across the amplitude values in each frame of the analysis, and subsequently defines the input signal's attack portion as ending at the frame with the maximum amplitude in the signal. More complicated envelopes would require a robust algorithm for detecting the attack portion of a signal but complex attack time measurement is outside the scope of the current project.

The addition of this algorithm to the brightness morphing system developed previously [4] should allow for six discrete types of audio morph:

- 1) Morph all characteristics
- 2) Morph brightness only (by means of spectral centroid manipulation)
- 3) Morph softness only (by means of attack time manipulation)
- 4) Morph all characteristics other than softness
- 5) Morph all characteristics other than brightness
- 6) Morph all characteristics other than brightness and softness

To morph all characteristics a direct interpolation of the entire acoustic feature set is used. For the brightness morph, a spectral tilt is applied, such that the spectral centroid of the source sound shifts towards that of the target sound. The softness morph timestretches the source sound's attack portion such that its duration moves towards that of the target sound's attack duration. The fourth type of morph is based on a combination of the 'all characteristics' and brightness morph, and performs a full interpolation together with a reverse spectral tilt to move the spectral centroid back to its original value. The other morph types are achieved by carrying out a combination of these three processes.

Note that a seventh type of morph, brightness *and* softness, would also be achievable by daisy-chaining processes 2 and 3. However, in the timbre morpher interface, end users might find one control

that adjusts two attributes counter intuitive and, at this stage, the aim was to establish if softness could be manipulated independently from brightness and the other characteristics. Therefore this seventh morph type is not included at the next experimental stage.

2.2. Main experiment

A complete timbre morpher should be capable of morphing any number of attributes independently of one another. The aim of the main experiment was to determine whether or not the current morphing system, with the addition of the softness algorithm shown in Figure 2, could manipulate two timbral attributes, independently from one another, and also from the other characteristics of the source and target sounds.

A pairwise dissimilarity experiment was designed to reveal the perceptual dimensionality of a stimulus set created using the discrete morph types listed in section 2.1.2. Three outcomes were possible. Firstly, morphs in softness (type 3) and brightness (type 2) reveal no clear pattern, or a pattern which demonstrated that the morphed stimuli were perceived in an incorrect order (for example with the softest sound being perceived as nearer to the hardest sound than to other stimuli), then the code, or the test procedure, would need to be reconsidered. Secondly, it was possible that the stimuli might be perceived to be in broadly the correct order, but that brightness- and softness-morphed stimuli might exhibit a high degree of perceptual correlation, despite having independent acoustic correlates. This would indicate that the chosen acoustic correlates, or the code for their manipulation, required respecification. Thirdly, and ideally, the softness and brightness morphs would produce perceived movement in discrete, orthogonal dimensions against the other characteristics, showing that the code, and its underlying acoustic manipulation, had produced the desired perceptual variation.

Listeners were asked to grade pairs of stimuli using a continuous scale, with end-points labelled ‘most dissimilar’ and ‘the same’, shown in Figure 4.

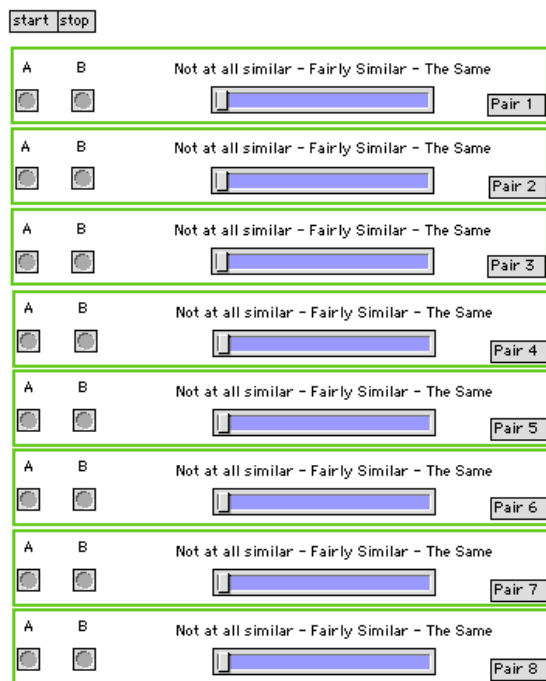


Figure 3 Max/MSP listener interface showing two stimuli per slider. Listeners were asked to rate how similar A sounded to B using the sliders

In order to allow for a suitable range of data for a Multidimensional Scaling (MDS) analysis to potentially reveal movement in up to 4 dimensions, at least sixteen stimuli were required [13]. In order to fulfill this requirement and to have a uniform number of variations in each of the morph types, 18 stimuli were generated, as shown in Table 3.

These stimuli were generated using various percentages of morphing between one source and one target sound (synthesised, so that amplitude, fundamental frequency, envelope, and duration could remain constant). The source was a sawtooth with an attack time of 5ms; the target was a triangle wave with a 250ms attack time. These particular waves were chosen to provide clear and simply quantifiable differences in attack time,

spectral centroid and harmonic structure, allowing a wide range of morphed stimuli. Stimuli throughout the experiment were loudness equalised according to the experimenter’s ear.

Label	% trgt softness?	% target brightness?	% target in <i>all</i> characteristics?
S1	0	0	0
S2	50	0	0
S3	100	0	0
B1	0	30	0
B2	0	60	0
B3	0	100	0
eeSB1	0	0	30
eeSB2	0	0	60
eeSB3	0	0	100
E1	30	30	30
E2	50	60	60
E3	100	100	100
eeB3	0	0	30
eeB2	50	0	60
eeB3	100	0	100
eeS1	0	30	30
eeS2	0	60	60
eeS3	0	100	100

Table 3 Stimulus set for main experiment. Labels refer to object space plots used in MDS analysis (see section 2.2.1)

For each stimulus, the percentage of morphing towards the target was chosen to provide (to the experimenter’s ears) perceptually even steps across each intended dimension. The stimuli were stored in PCM Wave format, at 16 Bit, 44.1Khz, in Mono, and played back on Sennheisser HD550 headphones.

A panel of 14 listeners, each with some experience of critical listening to recorded audio, took part in the test. Each listener was presented with a series of paired stimuli and was asked to rate the similarity of the two stimuli in each pair using a continuous scale with end-points labeled ‘most dissimilar’ to ‘the same’. 174 pairs, randomly ordered, were presented for comparison over two tests, and rated by means of the sliders on a 100

point hidden dissimilarity scale. The 174 pairs consisted of 171 pairing each sound against each other, plus 3 repeats (S1-S3, EES1-EES3, and B1-B3) which were included to enable assessment of listener consistency if required.

2.2.1. Results from main experiment

The listener responses were used to create a dissimilarity matrix, and analysed by MDS INDSCAL [15] analysis.

The ‘measures-of-fit’ calculated by the MDS analysis were examined, showing that the ‘s-stress’, decreased from ~ 0.2032 with a 1-D fit, to ~ 0.0227 with a 2-D fit, ~ 0.0034 with a 3-D fit, and ~ 0.0001 with a 4-D solution, as shown in Table 4.

Dimensionality	S-stress
1-D	~ 0.2032
2-D	~ 0.0227
3-D	~ 0.0034
4-D	~ 0.0001

Table 4 ‘s-stress’ decrease against dimensionality in MDS solutions

This indicated a good fit to the 3-D solution, as the unaccounted for variance did not decrease significantly between the 3-D and 4-D solutions, though the ‘scree plot’ shown in Figure 5 does not show a readily interpretable knee at 3-D.

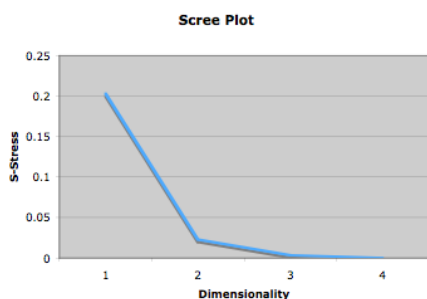


Figure 5 - Scree plot of s-stress values for unaccounted for variance in 1-D to 4-D MDS solutions

It has been previously shown that if RSQ, or squared correlation factor, does not improve by more than ~ 0.05 per dimension, this is a good indicator that the optimal fit for the data has been found by the solution [14].

The RSQ found by the MDS analysis here increased from ~ 0.8838 with a 1-D fit, to ~ 0.9481 with a 2-D fit, ~ 0.9999 with a 3-D fit. Combined with the low stress values, this RSQ improvement at 3-D also indicated that the 3-D solution was the most appropriate. RSQ values for each of the dimensional models are shown in Table 5, indicating that the 3-D plot is the most appropriate for the data.

Dimensionality	RSQ	Improvement
1-D	0.88381	
2-D	0.94813	0.06432
3-D	0.99994	0.05181
4-D	1.00000	0.00006

Table 5 Improvement in RSQ values against dimensionality in MDS solutions to listener response data

The 3-D object spaces (rotated to an appropriate angle) were then examined to reveal the positions of stimuli and directions of morph-produced change. When all of the stimuli were plotted in the object space, these plots were extremely cluttered, therefore for the purpose of clarity the Figures presented here include only the key stimuli. Figure 5 shows the object space at a default rotation, whilst Figure 6 shows the key stimuli in an object space rotated to 85 degrees.

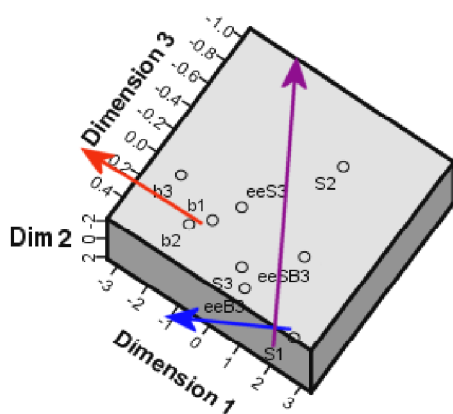


Figure 6 - Key stimuli shown at a default rotation. Movement in the ‘everything else’ dimension is shown in purple, the ‘brightness’ dimension in red, and ‘softness’ in blue (see Table 3 for key)

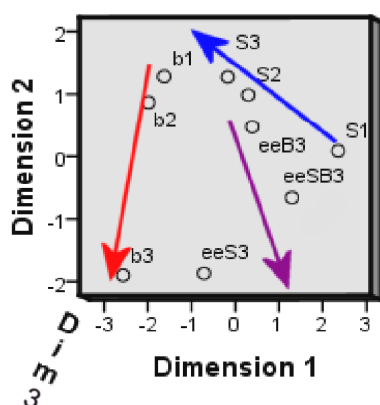


Figure 7 - key stimuli shown at a rotation of 85 degrees – as before, movement in the ‘everything else’ dimension is shown in purple, the ‘brightness’ dimension in red, and ‘softness’ in blue (see Table 3 for key)

Both figures show softness morphed stimuli (S1-3), brightness morphed stimuli (B1-3), and ‘everything apart from softness or brightness’ morphed stimuli (eeSB1-3). In the default rotation, shown in Figure 6, it is difficult to visualise the movement, from one stimulus to the next, in each direction – the brightness morphed stimuli are close together, and the full softness morph (S3) seems to be closer to S1 than S2 appears. The

spacing of the stimuli, and the direction of movement, were generally easier to see when the object space is rotated to 85 degrees, as in Figure 7. Combining both rotations was useful to help visualize the orthogonal movement across the three dimensions in the dataset, which was not clear in a single object space due to the limitations of a 2D plot.

At the 85 degree rotation, the brightness, softness, and ‘everything else’ stimuli appeared both more evenly spaced, and in broadly orthogonal directions, indicating that listeners perceived the morphs in these directions to be independent.

Together, these results demonstrate the unidimensionality of the brightness and softness morphing routines and their independence from other timbral attributes, collectively labeled ‘everything else’ here.

2.3. Subsequent verbal experiment

MDS analysis is unable to provide names for the dimensions revealed, and a subsidiary verbal elicitation experiment was therefore necessary to establish whether listeners perceived the morpher to be manipulating the intended timbral attributes.

Label	% softness?	% brightness?	% everything else?
S1 (A1)	0	0	0
S2 (A2)	50	0	0
S3 (A3)	100	0	0
S1 (B1)	0	0	0
B2 (B2)	0	30	0
B3 (B3)	0	100	0
S1 (C1)	0	0	0
eeSB2 (C2)	0	0	30
eeSB3 (C3)	0	0	100

Table 6 Stimulus set for subsequent verbal experiment. Labels in brackets were used in the test interface (see Figure 8)

Selected stimuli from the main experiment, shown

in Table 6, were presented to new listeners. These particular stimuli were chosen to demonstrate, as clearly as possible, timbral changes in each of the three dimensions discussed above. Stimuli A1-A3 differed only in terms of the intended softness dimension, stimuli B1-B3 only in terms of intended brightness, and stimuli C1-C3 only in terms of intended ‘everything else’. These three triplets of stimuli were presented to listeners using the interface shown in Figure 8.

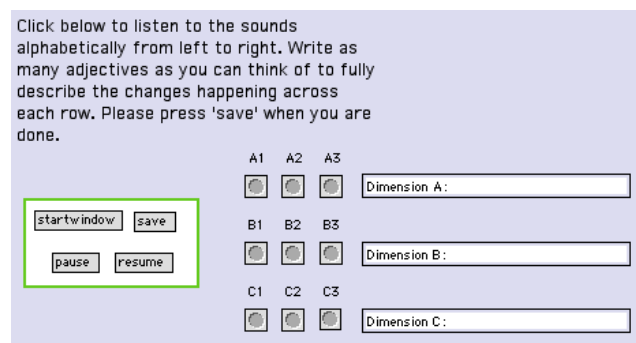


Figure 8 - Listener interface for verbal experiment.

Taking each triplet in turn, the listeners were asked to listen to its three component stimuli and to describe fully the inter-stimulus differences they perceived, using as many adjectives as they felt necessary. Twenty listeners undertook the test.

2.3.1. Results from subsequent verbal experiment

Eighty listener responses were elicited across the three types of morph. Taking the responses to each morph type in turn, an independent academic grouped together similar responses based on, for example, synonyms, antonyms and commonalities. The number of listener responses in each group was summed to determine a weighting factor, indicating the perceptual importance of each group, by dividing the number in each group by the total number of responses, as shown in Figures 9-11.

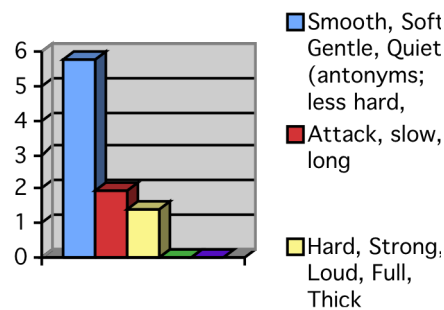


Figure 9 - Verbal groupings and weighting for the terms used to describe changes between softness-morphed stimuli

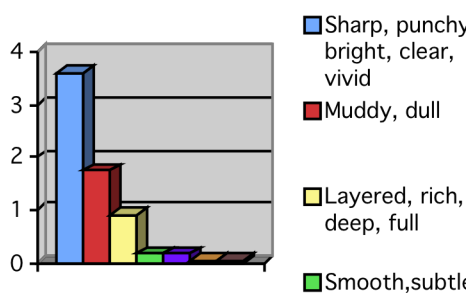


Figure 10 - Verbal groupings and weighting for the terms used to describe changes between brightness-morphed stimuli

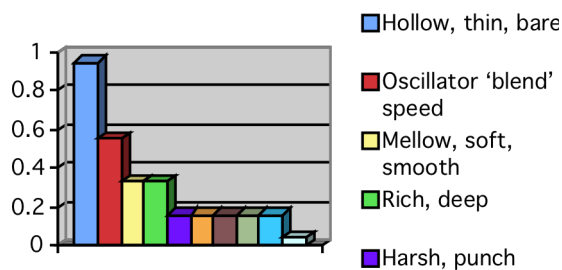


Figure 11 - Verbal groupings and weighting for the terms used to describe changes between everything-else-morphed stimuli

The first morph type (A1-3) was labeled as soft, smooth, gentle, or quiet, by a large margin. Quiet was an unexpected adjective to find in the most prominent group for this type of morph, particularly as the stimuli were loudness equalised.

It could be that that this listener was associating a softer attack with a musical note being performed quietly (a bowed string for example). Nevertheless, listener perception of this dimension was broadly as intended.

The second morph type (B1-3) was labeled as bright, sharp, punchy, clear, or vivid, by a good margin. Antonyms of bright and clear, such as dull, and muddy, were also popular choices. It is surprising that some listeners perceived these sounds in the opposite way to the majority of listeners, but, this could be due to misinterpretation of the instructions, whereby these listeners simply labeled the wrong end of the scale. Again the listener perception of the stimuli in this dimension seems to be broadly correct.

The third morph type (C1-3) was labeled as hollow, thin, or bare, but this dimension featured more terms and a smaller margin between the most prominent group and the other groupings. This dimension did not have an intended label but, for the particular source and target sounds used in this experiment, a change in hollowness might be expected from the 'everything else' morph since the acoustic change involved was to the relative intensities of odd vs even harmonics. Therefore this dimension would seem to have been perceived as intended, but by a smaller majority than dimensions one or two.

3. CONCLUSIONS

The preliminary experiment established that the contribution to perceived softness of inharmonicity is negligible, and that attack time is the main acoustic correlate of softness, according to the listening panel. A softness matrix showing the relative contribution of inharmonicity and attack time to softness was therefore not required in the softness morphing stage of the timbre morpher.

The main experiment showed, by means of MDS analysis of listener responses to a pairwise dissimilarity experiment, that, within the stimulus

set comprised of several types of morphed output, there were three broadly orthogonal dimensions.

The subsequent verbal elicitation experiment revealed that listeners perceived movement along these orthogonal dimensions to be broadly as intended – changes to brightness, softness, and everything else. This experiment also demonstrated that the SMS-based morphing system is adaptable to a timbral attribute with a temporal acoustic correlate (softness was manipulated using changes to attack time).

These results indicate that a complete “Timbre Morpher”, based on manipulating specific underlying acoustic correlates of chosen attributes, is viable

Further work should deal initially with attributes with timbral and acoustical overlap, in order that methods for dealing with such overlap can be devised, implemented within the morpher, and tested.

4. REFERENCES

- [1] M Cowell M Slaney, B Lassiter, 'Automatic Audio Morphing' *In Proc. ICASSP, Atlanta, Georgia* pp1001-4 (1996).
- [2] D Furlong C Hope, 'Time-Frequency Distributions for Timbre Morphing: The Wigner Distribution Versus the STFT' *Proc. SBCMIV, Brasilia, Brasil* pp99-110 (1997).
- [3] E Tellman B Holloway, L Haken, 'Timbre Morphing of Sounds with Unequal Numbers of Features' *Journal of the Audio Engineering Society* **43(9)**, pp678-689 (1995).
- [4] Williams D Brookes T, 'Perceptually-motivated audio morphing: Brightness' *122nd Audio Engineering Society Convention*, 7035 (2007).
- [5] Howard DM Disley AC, 'Spectral correlation of timbral descriptors of the pipe organ' *Proceedings of the Baltic-Nordic Acoustics Meeting BNAM-04*, (2004).

- [6] DJ Freed, 'Auditory correlates of perceived mallet hardness for a set of recorded percussive sound events' *J. Acoustic. Soc. Am.* **87(1)**, p311 (1989).
- [7] S Lakatos, 'A common perceptual space for harmonic and percussive timbres' *Perceptual Psychophysics* **26**, p1426 (2000).
- [8] NH Fletcher, 'The Nonlinear Physics of Musical Instruments' *Reports on Progress in Physics* **62**, p723 (1999).
- [9] J Stepanek, 'Musical Sound Timbre: Verbal Description and Dimensions' *Proc. of the 9th Int. Conference on Digital Audio Effects, Montreal, Canada* (2006).
- [10] TR Rossing, 'The science of sound' p290. Addison-Wesley. (1990).
- [11] J Escribe X Rodet, S Durigon, 'Improving score to audio alignment: Percussion alignment and precise onset estimation' *Ircam Centre Pompidou, Analyse Synthese CNRS UMR 9912*, (2004).
- [12] X Serra, 'Musical sound modelling with sinusoids plus noise' *Musical Signal Processing*. p92. Swets &Zeitlinger. (1997).
- [13] R Anderson J Hair, R Tatham, W Black, 'Multivariate Data Analysis'. Prentice Hall, Upper Saddle River, New Jersey (1998).
- [14] BR Astill, 'Humanised Statistics: An investigation of social values in a senior secondary school milieu'. University of Adelaide. (1994).
- [15] JO Ramsay, 'Maximum likelihood estimation in multidimensional scaling' *Psychometrika* **42**, p241 (1977).
- [16] H Fastl, E Zwicker, 'Psychoacoustics Facts and Models'. Springer, pp173,240-242, (2006).