

1 Acoustical correlates of perceptual blend in timbre dyads and triads

2 Sven-Amin Lembke

3 Centre for Interdisciplinary Research in Music Media and Technology (CIRMMT),

4 Schulich School of Music, McGill University, Montréal, Québec, Canada

5 Kyra Parker

6 Centre for Interdisciplinary Research in Music Media and Technology (CIRMMT),

7 Schulich School of Music, McGill University, Montréal, Québec, Canada

8 Eugene Narmour

9 Department of Music, University of Pennsylvania, Philadelphia, Pennsylvania, USA

10 Stephen McAdams

11 Centre for Interdisciplinary Research in Music Media and Technology (CIRMMT),

12 Schulich School of Music, McGill University, Montréal, Québec, Canada

Abstract

13

14 Achieving a blended timbre for particular combinations of instruments, pitches, and
15 articulations is a common aim of orchestration. This involves a set of factors that this
16 study jointly assesses by correlating the perceptual degree of blend with the underlying
17 acoustical characteristics. Perceptual blend ratings from two experiments were
18 considered, with the stimuli consisting of: 1) dyads of wind instruments at unison and
19 minor-third intervals and at two pitch levels, and 2) triads of wind and string
20 instruments, including bowed and plucked string excitation. The correlational analysis
21 relied on partial least-squares regression, as this technique is not restricted by the
22 number and collinearity of regressors. The regressors encompassed acoustical
23 descriptors of timbre (spectral, temporal, and spectrotemporal) as well as ones
24 accounting for pitch and articulation. From regressor loadings in principal-components
25 space, the major regressors leading to substantial and orthogonal contributions were
26 identified. The regression models explained around 90% of the variance in the datasets,
27 which was achievable with less than a third of all regressors considered initially. Blend
28 seemed to be influenced by differences across intervals, pitch, and articulation. Unison
29 intervals yielded more blend than did non-unison intervals, and the presence of plucked
30 strings resulted in clearly lower blend ratings than for sustained instrument
31 combinations. Furthermore, prominent spectral features of instrument combinations
32 influenced perceived blend.

33 *Keywords:* timbre, blend, orchestration, instrument dyads/triads, acoustical
34 descriptors, multivariate regression

35 Acoustical correlates of perceptual blend in timbre dyads and triads

36 In orchestration, composers may consider several factors when they intend to
37 achieve a *blended* timbre between two or more instruments playing synchronously.
38 There is the choice of suitable instruments that can yield a blended combination, which
39 depends on the acoustical traits of these instruments. The remaining factors involve
40 more musical considerations: whether instruments will be playing in unison or
41 non-unison, which instrument is assigned to the top voice in non-unison passages, in
42 what registral range the instruments will be playing, and what kind of articulation they
43 will employ (e.g., bowed or plucked string). When it comes to establishing general
44 associations between the perception of timbre blend and its underlying acoustical
45 characteristics, the joint assessment of these factors will assist in predicting the
46 perceived degree of blend for combinations of instruments, pitches, and articulations.

47 Previous research has defined perceived timbre blend as the auditory fusion of
48 concurrent instrumental sounds, where individual sounds become less distinct. The
49 most common method to measure perceived blend employs rating scales (Kendall &
50 Carterette, 1993; Lembke, Levine, & McAdams, in press; Lembke & McAdams, 2015;
51 Sandell, 1995; Tardieu & McAdams, 2012). All studies found that spectral features
52 influence blend, but employed different approaches to spectral description. One
53 approach used the global descriptor *spectral centroid*, i.e., the amplitude-weighted
54 frequency average of a spectrum. The composite (or sum) of the individual sounds'
55 centroids was found to predict blend in unison dyads best (Sandell, 1995; Tardieu &
56 McAdams, 2012), whereas for non-unison dyads, the absolute difference in individual
57 spectral centroids served as the more reliable predictor (Sandell, 1995).

58 Another approach to spectral description has considered the influence of
59 prominent spectral features, such as maxima or *formants*. Similar to the relevance of
60 formants in describing the acoustics of the human voice (Fant, 1960), wind instruments
61 in particular exhibit formant structures that remain largely invariant across pitch
62 (Lembke & McAdams, 2015; D. Luce & Clark, 1967; D. A. Luce, 1975; Schumann,
63 1929). Their identification and description can be achieved through spectral estimations

64 that are aggregated across an instrument's complete pitch range (Lembke & McAdams,
65 2015), and therefore can be considered *pitch-generalized*. Reuter (1996) has argued that
66 similarity between instruments' formant structures can explain blend. Hardly
67 distinguishable instrument pairings can exhibit very similar formant locations (e.g.,
68 horn and bassoon), whereas the strongly pronounced, unique formant structure of the
69 oboe may hinder it from blending with most other instruments.

70 Frequency relationships between the most prominent *main formants* appear to
71 influence blend critically (Lembke & McAdams, 2015). In dyads comprising a recorded
72 wind-instrument sound and a synthesized analogue to that instrument, whose
73 main-formant frequency could be shifted relative to that of the recorded sound, blend
74 decreased drastically as the frequency of the synthesized formant exceeded that of the
75 recorded sound. This relative dependency relates to musical performance, where
76 accompanying musicians adjust their main formants to be lower than when playing as
77 the leading instrument (Lembke et al., in press).

78 Apart from spectral properties, differences between temporal features, such as
79 note attacks or onsets, have been found to explain blend as secondary factors for unison
80 dyads (Sandell, 1995). However, their influence becomes more dominant as attacks turn
81 impulsive: shorter durations and steeper attack slopes lead to reduced blend (Tardieu &
82 McAdams, 2012).

83 With respect to those musical factors unrelated to timbre, blend for unison dyads
84 is perceived as stronger than for non-unison combinations (Kendall & Carterette, 1993;
85 Lembke et al., in press). Furthermore, the assignment of instruments to the upper and
86 lower pitches in non-unison intervals resulted in differences in perceived blend between
87 instrument inversions in one study (Kendall & Carterette, 1993), but lacked a
88 comparable effect in another (Sandell, 1995). All of these studies on blend are limited to
89 dyadic contexts, leaving open how the obtained results and proposed hypotheses would
90 fare in combinations of three or more instruments. Little work has been published on
91 timbre combinations in triadic contexts (Kendall, 2004; Kendall & Vassilakis, 2006,
92 2010), and none of these papers address issues directly related to blend.

93 With the aim of predicting perceived blend between arbitrary instrument
94 combinations, linear correlation or regression can be employed to associate blend
95 measures with single acoustical features (Sandell, 1995; Tardieu & McAdams, 2012),
96 without, however, making it possible to assess how several acoustical descriptors could
97 jointly contribute to the explanation of blend measures. This limitation can be
98 overcome by *multiple linear regression* (MLR). Past attempts have succeeded in
99 explaining up to 63% of the variance in blend ratings for mixed-instrument dyads
100 (Sandell, 1995). Similarly, MLR models also explained up to 87% of the variance in
101 blend ratings across dyads in which the role of local, parametric variations of the
102 main-formant frequency was studied (Lembke & McAdams, 2015).

103 Yet, the MLR approach also has clear limitations. High collinearity among
104 independent variables (regressors) or a low number of cases compared to the number of
105 regressors may both lead to less reliable and less valid results as well as mathematically
106 ill-defined solutions. This becomes problematic given the aim of the current paper,
107 because many spectral descriptors are known to exhibit a high inter-correlation
108 (Peeters, Giordano, Susini, Misdariis, & McAdams, 2011). For conventional MLR, this
109 leaves two options: 1) disregarding the collinearity, at the risk of obtaining less reliable
110 or invalid results or 2) eliminating regressors that are collinear to a reference regressor,
111 i.e., one found to predict blend most strongly in simple linear regression. However, the
112 latter approach risks excluding variables that might perform even better than the
113 selected one once they interact with other regressors.

114 A viable solution to deal with collinearity is to employ a dimension-reduction
115 technique like *principal component analysis* (PCA) that reduces a high quantity of
116 regressors to a small number of substitute or latent variables, i.e., *principal*
117 *components* (PCs), which are orthogonal to one another. These PCs can thereafter
118 serve as regressors that represent the common aspects for groups of collinear descriptors
119 (e.g., Giordano, Rocchesso, & McAdams, 2010).

120 A promising regression method that uses PCA as an integral part is *partial*
121 *least-squares regression* (PLSR), which originates from the discipline of chemometrics,

122 but has more recently been applied within the field of auditory perception (Eerola,
123 Lartillot, & Toivainen, 2009; Kumar, Forster, Bailey, & Griffiths, 2008; Rumsey,
124 Zieliński, Kassier, & Bech, 2005). PLSR allows analysis of complex correlational
125 relationships among perceptual measures and arrays of acoustical or psychoacoustical
126 variables.

127 The current investigation uses PLSR in an attempt to predict blend ratings from
128 perceptual experiments. The perceptual data are collected on a diverse set of variables
129 that affect timbral blend and orchestration, including different instruments, pitches, and
130 unison and non-unison intervals, as well as dyadic and triadic contexts. The set of
131 potential regressors consists of a wide range of acoustical measures that, through several
132 stages of PLSR models, are continually refined to retain only the most relevant
133 regressors and, importantly, ones that are independent of each other.

134 Method

135 Partial least-squares regression (PLSR)

136 Predicting a single measure of blend through a set of regressors relating to
137 acoustical descriptors can be expressed mathematically by associating the column vector
138 of blend ratings y with an $n \times m$ matrix X , which encompasses n cases (e.g., stimulus
139 conditions) across m regressors. Conventional MLR employs the relationship $y = X \cdot b$,
140 with b being a vector of regression coefficients of length m . PLSR represents algorithms
141 that employ an inherent coupling between MLR and PCA (Geladi & Kowalski, 1986),
142 allowing large m relative to n and even collinearity among the m regressors.

143 PLSR decomposes X into k principal components (PCs), yielding the relationship
144 $X = T \cdot P'$, with T representing an $n \times k$ matrix of *scores* and P an $m \times k$ matrix of
145 *loadings*. Unlike computing a PCA on X independently and inputting the obtained
146 scores T into MLR, PLSR achieves the component decomposition by maximization of
147 the inherent covariance between y and X , leading to a better predictive relationship.
148 The loadings P can be understood as vectors for the m regressors in k -dimensional
149 space, describing the degree to which regressors contribute to individual PCs and also

150 showing the collinearity or independence among regressors. The PLSR technique used
151 here is SIMPLS (de Jong, 1993), in its implementation for MATLAB.¹

152 **Performance, predictive power, and reliability.** Regression performance
153 evaluates the variation in y that is explained by the model. The measure R^2 describes
154 both the global and component-wise performance, with the latter quantifying the
155 relative contribution of PCs. With increasing k , however, models are prone to
156 overfitting the data, at the cost of predictive power when applied to other data sets. In
157 order to assess the predictive power of models, sixfold cross validation (CV) is
158 employed, partitioning the n cases into six subsets of similar size, building models based
159 on five subsets, assessing the error in predicting the remaining, excluded subset, and
160 repeating the last two steps for all permutations of subsets. CV allows the computation
161 of an alternative measure of explained variance, Q^2 (Wold, Sjöström, & Eriksson, 2001).
162 Just as R^2 evaluates the sum of squared deviations between the *fitted* and actual y , Q^2
163 evaluates the sum of squared deviations between the *predicted* and actual y , with these
164 predictions made for the excluded subsets across all CV permutations. Together, R^2
165 and Q^2 can be taken as the upper and lower benchmarks of the model, respectively, in
166 terms of explaining the data and assessing the degree of predictive power. The selection
167 of the optimal number of components k considers two independent criteria: 1) the
168 component-wise gain in R^2 , and 2) the component-wise decrease in CV prediction error,
169 with k being chosen when both measures cease to exhibit substantial improvements for
170 additional PCs.

171 **Identifying relevant and independent regressors.** The current PLSR
172 analysis aims to reduce the number of investigated regressors in X to those of greatest
173 utility in explaining y as well as further reducing it to a selection of regressors that are
174 relatively independent of each other. The chosen approach consists of three stages of
175 sequentially evaluating and refining PLSR models: 1) An initial model is obtained for
176 the original matrix X_{orig} of all regressors considered (see Table 3). 2) Based on the
177 loadings P_{orig} from the first PLSR model, one half of the variables are identified that
178 act as the strongest predictors. More specifically, only those regressors are retained for

179 which the Euclidean distances across k dimensions exceed the distribution median
180 (Q_{50}), leading to the computation of another model based on the reduced matrix $X_{Q_{50}}$.
181 3) The following stage distills the regressors down to those that explain y through
182 ideally independent contributions. Such independence or orthogonality is achieved for
183 loadings P that point in perpendicular directions in k -dimensional space. To this aim,
184 imagine the resulting loadings $P_{Q_{50}}$ rotated so as to align the most dominant variable
185 loading along one axis of a Cartesian coordinate system. Relative to this variable,
186 loadings aligned along the remaining $k - 1$ axes exhibit maximal independence. To
187 obtain an ideally independent set of variables, the selection is constrained to variables
188 such that the angles ϕ_i between variable loadings and the i th axis are less than 22.5° .
189 This constraint yields an approximately orthogonal set of regressors X_{ortho} , on which
190 the final PLSR model is computed.

191 **Perceptual data sets**

192 The regression analysis considers two data sets that originate from listening
193 experiments in which participants provided blend ratings for dyads or triads. The two
194 experiments were unrelated with respect to their original motivation and experimental
195 design, yet they employed similar blend ratings, with the medians across participants
196 taken as the dependent variable y to be modeled through PLSR.

197 The stimuli were presented over a standard two-channel stereophonic loudspeaker
198 setup inside an Industrial Acoustics Company double-walled sound booth. They were
199 drawn from recorded instrument samples from the Vienna Symphonic Library² (VSL),
200 supplied as stereo WAV files (44.1 kHz sampling rate, 16-bit amplitude resolution). In
201 separate pilot experiments, all stimuli had been both adjusted for perceptual synchrony
202 between sounds constituting the dyads and triads and equalized for loudness within the
203 dyad and triad sets independently. Adjustments for synchrony were based on consensus
204 by three people for dyads and two for triads. The loudness equalization was conducted
205 subjectively, anchored to a global reference for all dyad or triad conditions. The
206 equalization was conducted by five people for dyads and six for triads. Gain levels were

207 determined that equalized stimulus loudness to the global reference. These gain levels
208 were based on median values across participants; all corresponding interquartile ranges
209 were less than 4 dB.

210 For the main experiments, participants with varying degrees of musical experience
211 were recruited from the McGill University community. All participants passed a
212 standardized pure-tone audiogram (ISO 389–8, 2004; Martin & Champlin, 2000)
213 ensuring that thresholds at all audiometric frequencies were less than or equal to
214 20 dB HL. Informed consent was obtained, and both studies were certified for ethical
215 compliance by the McGill University Research Ethics Board II.

216 **Dyads.**

217 ***Participants.*** Nineteen people took part in the experiment (12 female and
218 seven male) with a median age of 21 years (range: 18–46). Among the participants, nine
219 considered themselves amateur musicians, two as professional musicians, and eight as
220 non-musicians. All were compensated financially for their participation in the hour-long
221 experiment.

222 ***Stimuli.*** The stimulus set comprised a total of 180 dyads that resulted from the
223 combination of several factors. Six wind instruments, namely, (French) horn, bassoon,
224 oboe, C trumpet, B \flat clarinet, and flute, formed the 15 possible non-identical-instrument
225 pairs listed in Table 1. These instrument pairs occurred at two pitch levels: C4
226 ($f_0 = 261.6$ Hz) and G4. Furthermore, dyads comprised both unison and minor-third
227 intervals, including the inverse voicing of instruments for the latter, resulting in a total
228 of three interval conditions. Based on the two pitch levels, minor thirds occurred at the
229 pitches C4-E \flat 4 and G4-B \flat 4.

230 All VSL samples were sustained, non-vibrato recordings, performed at *mezzoforte*
231 dynamics, and were limited to the signal in the left channel. Both instruments were
232 simulated as being captured by a stereo main microphone at spatially distinct locations
233 inside a mid-sized, moderately reverberant room. Encompassing a volume of 600 m³,
234 the relatively absorbent room yielded a reverberation time $T_{30} = 0.4$ s; due to the
235 configuration inside the room being fully symmetric, identical frequency responses

236 applied to both instruments.³ The spatial locations of instruments included both
237 possible orientations (e.g., horn left of bassoon and vice versa). Overall, this resulted in
238 the full-factorial combination of 15 pairs \times 2 pitches \times 3 intervals \times 2 orientations =
239 180 dyads. All stimuli had a duration of 1200 ms, with artificial offsets imposed by a
240 100-ms linear amplitude ramp. A set of 12 representative dyad stimuli can be found in
241 the Supplemental Material Online section.

242 [— Insert Table 1 about here. —]

243 **Procedure.** Participants heard individual dyads in randomized order and were
244 asked to rate their degree of blend, employing a continuous slider scale with the verbal
245 anchors *most blended* and *least blended* visualized on a computer screen. Ahead of the
246 main experiment, participants had been familiarized with the degree of possible
247 variation in blend among all dyads and had completed 15 practice trials on a separate
248 but comparable stimulus set.

249 **Triads.**

250 **Participants.** Twenty people (15 female and five male) with a median age of
251 21 years (range: 19–64) completed the experiment. Thirteen participants classified
252 themselves as amateur musicians, with the remaining seven being non-musicians. All
253 were remunerated for the hour-long experiment.

254 **Stimuli.** The stimuli comprised 20 triads, representing only a selection of the
255 vast multiplicity of possible instrument and pitch combinations. In order to focus on
256 timbral characteristics, all triads formed the same chord with pitches C4, F4, and B \flat 4,
257 thus controlling for contextual effects with pitch register, chroma, and height. This
258 chord choice of stacked perfect-fourth intervals avoided standard major and minor
259 triads and generated the same consonances (perfect fourths) and a single dissonance
260 (minor seventh). Such quartel chords were neutral enough not to draw attention to any
261 one melodic voice, while allowing ‘inside’, middle voices to be easily heard.

262 In terms of instrumentation, the triads were composed of flute, oboe, B \flat clarinet,
263 tenor trombone and cello sounds, corresponding to the instrument families woodwinds,
264 brass, and strings. The instrument selection for triads (see Table 2) comprised mixtures

265 between two or three instrument families. Furthermore, the selection included all
266 woodwind reed types (air jet, single and double reed) and two different excitation types
267 for strings (bowed and plucked excitation; *arco* and *pizzicato*, respectively), with each
268 distinction represented by a single instrument (e.g., oboe for double reed, cello for string
269 instrument). Instruments would only take on pitches based on conventional voice
270 assignments given a particular mixture. For instance, the cello only occurred at the two
271 lower pitches, whereas the flute was always highest in pitch relative to other woodwinds.
272 Each instrument appeared in from six to 10 triads (counting different excitation types
273 as separate instances).

274 All samples were taken as stereo files from VSL, with woodwind samples
275 comprising sustained sounds at *mezzoforte* dynamics and without vibrato. The
276 trombone samples were similar, but at *mezzopiano* dynamics. The *arco* cello samples
277 were recorded at *mezzoforte* dynamics. Unlike the wind instruments, they decayed after
278 just a brief bow stroke, in order to be more similar to the *pizzicato* versions, which
279 occurred at *forte* to allow for a longer sound decay. All cello sounds contained vibrato.
280 The total duration for all triads was limited to 850 ms by applying an artificial 100-ms
281 linear amplitude-decay ramp. A set of 10 representative triad stimuli can be found in
282 the Supplemental Material Online section.

283 [— Insert Table 2 about here. —]

284 **Procedure.** Participants were asked to sort all triads based on their relative
285 degree of blend along a scale continuum with the verbal anchors *most blended* and *least*
286 *blended*. At the beginning, visual icons for all triads were randomly arranged on a
287 computer screen and could be dragged around or clicked on to trigger sound playback.
288 Participants were first asked to identify two triads perceived as exhibiting the highest or
289 lowest blend, to assign them to the extremes of the visualized continuum and then to
290 position all remaining triads along the continuum. The sorting was conducted twice,
291 the first counting as a practice round meant to familiarize participants with both the
292 experimental task and the triads, the second serving as the main experiment.

293 **Acoustical descriptors**

294 For each data set, a collection of acoustical measures constitute the regressors in
 295 matrix X . The measures include spectral, temporal, and spectrotemporal acoustical
 296 descriptors of timbre, as well as other potentially relevant features such as differences in
 297 fundamental frequency. Table 3 lists all the investigated descriptors, specifying how
 298 individual descriptor values were associated with dyads and triads.

299 **Descriptor relationships within dyadic and triadic contexts.** As dyads
 300 and triads consist of several constituent sounds, their individual descriptor values need
 301 to be summarized to a single regressor value per stimulus by an association of some
 302 kind. For dyads with the constituent sounds a and b and the acoustical descriptor x ,
 303 the association considers the *difference* measure $\Delta x = |x_a - x_b|$ and the *composite*
 304 measure $\Sigma x = x_a + x_b$. Three associations are computed for triads with sounds a , b , and
 305 c , whose relationship along descriptor x is $x_a \leq x_b \leq x_c$. The triad *difference* considers
 306 the range between maximum and minimum, i.e., $\Delta x = x_c - x_a$. The *composite* sums all
 307 three values, i.e., $\Sigma x = x_a + x_b + x_c$. In addition, a third measure relates the
 308 *distribution* of the intermediate value x_b relative to the extremes, i.e.,
 309 $\Xi x = 2 \cdot (x_b - x_a) / \Delta x - 1$. Ξx varies from -1 to 1 . It is -1 when $x_a = x_b$, 0 when x_b is
 310 halfway between x_a and x_c , and 1 when $x_b = x_c$. These three regressor types apply to
 311 most of the investigated acoustical descriptors but not all, based on whether the
 312 association is appropriate or not, as indicated in Table 3.

313 [— Insert Table 3 about here. —]

314 **Timbre descriptors.**

315 **Spectral descriptors.** These descriptors assess properties associated with a
 316 time-averaged spectral representation. The investigated descriptors are computed on
 317 the output of one of two spectral-analysis methods: 1) analyses of the audio signals (\sim)
 318 for individual instrument samples from VSL (e.g., oboe at G4) by use of the *Timbre*
 319 *Toolbox* (Peeters et al., 2011) employing harmonic analysis, and 2) pitch-generalized ($^\circ$)
 320 spectral envelopes (Lembke & McAdams, 2015), which are estimated by fitting a curve

321 to partial tones aggregated across all available pitches from VSL (e.g., oboe from Bb3 to
 322 G6) and therefore allow for the characterization of an instrument’s pitch-generalized
 323 formant structure. Furthermore, the spectral descriptors can be distinguished as
 324 quantifying *global* and *local* spectral properties, as listed in Table 3. The global
 325 descriptors (S) include measures of spectral centroid (amplitude-weighted mean
 326 frequency), spectral slope (linear regression on the spectral envelope), spectral skewness
 327 (asymmetry in the spectral distribution), spectral kurtosis (peakier or flatter deviation
 328 from a Gaussian distribution), spectral spread (standard deviation of the spectral
 329 distribution), spectral roll-off (95th percentile of spectral-energy distribution), spectral
 330 decrease (spectral slope with low-frequency emphasis) and noisiness of the signal. They
 331 are described in detail in Peeters et al. (2011). The local, formant-related descriptors
 332 (F) require some elaboration.

333 [— Insert Figure 1 about here. —]

334 The formant structure derived from the pitch-generalized spectral envelope for a
 335 horn is shown in Figure 1. A set of frequencies indicate formant maxima (solid red
 336 lines; two formants are identified for this horn) and delineate their extent through lower
 337 and upper bounds (dotted lines) at which the magnitude has decreased by 3 dB. Note
 338 in the case of the horn in Figure 1 that there is no lower bound for the second formant,
 339 because the minimum point in the envelope between the two formants is not at least
 340 3 dB below the maximum of the second formant. In this study, the focus lies on the
 341 main formant F_1 , with F_{max} and F_{3dB} characterizing the frequency at its maximum
 342 magnitude and at the 3 dB upper bound, respectively (the latter appears to be more
 343 perceptually relevant; Lembke & McAdams, 2015). Two related measures, F_{slope} and
 344 $F_{\Delta mag}$, assess the relative importance of the main formant compared to the
 345 spectral-envelope regions lying above it. F_{slope} evaluates the (linear) spectral slope (grey
 346 diagonal) above the main formant in Figure 1, whereas $F_{\Delta mag}$ quantifies the level
 347 difference between the main-formant peak and the averaged magnitude of the spectral
 348 envelope above it (black arrow).

349 [— Insert Figure 2 about here. —]

350 Furthermore, the degree to which wind instruments are characterized by formant
351 structure varies, being strongest for oboe but much weaker for clarinet and flute. The
352 measure F_{prom} quantifies the prominence of up to two formants based on a cumulative
353 score that increases with the number and width of formant features. As illustrated in
354 Figure 2, the larger the total area covered by existing formant bounds (shaded
355 rectangles), the higher the prominence of an instrument’s formant structure, reflected in
356 F_{prom} being considerably higher for oboe (blue) than for flute (red). Two additional
357 difference measures, F_{freq} and F_{mag} , relate magnitude and frequency differences
358 between the formants of constituent instruments. More specifically, F_{mag} quantifies the
359 cumulative magnitude deviation between the constituent instruments’ spectral
360 envelopes at all formant frequencies they exhibit (vertical lines projected on the far
361 right). F_{freq} evaluates the cumulative frequency difference (horizontal line projected at
362 the top) between formants of the same order (e.g., main formant with main formant), if
363 they exist for both sounds.

364 **Temporal descriptors.** Three descriptors characterize the time course of the
365 amplitude envelope with respect to the attack (A) or onset portions of sounds,
366 considering attack time and attack slope descriptors (Peeters et al., 2011).

367 **Spectrotemporal descriptors.** A pair of descriptors account for spectral
368 variation across time, which the (static) spectral descriptors leave unaddressed.
369 Previous research has not reported specific spectrotemporal (ST) descriptors as being
370 relevant to blend, although temporal modulation of spectral components has been
371 discussed in the context of blend (Reuter, 2009). In the interest of using a
372 comprehensive set of timbre-related descriptors, two descriptors are included that
373 involve the commonly reported *spectral flux* (ST_{flux} , Peeters et al., 2011) and the
374 alternative measure *spectral incoherence* ($ST_{incoher}$), which quantifies the aggregate
375 deviations of spectral magnitude between successive time frames (Horner, Beauchamp,
376 & So, 2009).

377 **Other descriptors and variables.** The experimental designs involved factors
378 that were likely to explain variance in median blend ratings but were not related to or

379 not reliably measured through timbre features. Their relevance as potential regressors is
380 assessed by several categorical variables (C), in addition to acoustical descriptors that
381 could serve as equivalent predictors in application scenarios lacking a priori knowledge
382 of categorical distinctions, e.g., by quantifying pitch relationships or the loudness
383 balance between combined sounds. The categorical variables make binary or ternary
384 distinctions and for the use with PLSR are expressed as *dummy variables* (Martens,
385 Høy, Westad, Folkenberg, & Martens, 2001). A categorical variable is represented by as
386 many dummy variables as there are categories, with each category’s dummy variable set
387 to 1 for cases matching the category and 0 if not. As a result, these regressors yield
388 multiple loadings. For example, a binary categorical variable yields two loadings in
389 opposing orientations, with “-D1” or “-D2” being appended to the variable name to
390 symbolize the first and second categories of the variable, respectively (or -D0, -D1, and
391 -D2 for a ternary variable).

392 For triads, a strong distinction was expected beforehand for the presence (D1)
393 versus absence (D2) of *pizzicato* string sounds (C_{pizz}), as they are highly impulsive.
394 Similarly, the distinction between unison (D1) and non-unison (D2) dyads was also
395 expected to yield higher ratings for the former (C_{unison} and Δf_0). Additional regressors
396 account for the lower (D1) or higher (D2) pitch level (C_{pitch}) and difference between
397 pitches expressed in ERB units ($f_{0|ERB}$; Moore & Glasberg, 1983), interval type
398 ($C_{interval}$; D0: unison, D1: instrument A low, instrument B high, D2: instrument B low,
399 instrument A high), and instrument position at left (D1) or right (D2) in stereo space
400 ($C_{position}$). In addition, the production of dyads and triads also involved determining
401 relative mix or scaling ratios between the amplitudes of the constituent sounds forming
402 dyads or triads, which are also quantified to assess their possible influence on the blend
403 ratings (x_{mix}). For dyads, x_{mix} concerned a fractional gain value between 0 and 1 that
404 applied to one instrument, while $x_{mix} - 1$ scaled the other instrument. For triads, x_{mix}
405 concerned the sound-level balance between the constituent sounds (e.g., negative slope
406 for monotonously decreasing sound levels with ascending pitch).

Results

407

408 As mentioned under Method, PLSR analysis of a particular data set involves three
409 stages, beginning with the original set of regressors X_{orig} , then restricting the selection
410 to X_{Q50} , and finally attaining an approximately orthogonal selection of regressors
411 X_{ortho} . Although statistics for all three stages are reported in Tables 4 and 5, only the
412 results for the final stage X_{ortho} are presented in detail. In some of the following
413 visualizations, data points representing dyads or triads use the labels in Tables I and II,
414 respectively. A further distinction between dyads or triads containing instruments that
415 blend more strongly and those that blend weakly is made with color to help assess how
416 a given acoustical descriptor separates these instruments. For dyads, the instruments
417 clarinet, bassoon, and horn lead to the highest blend ratings of comparable magnitude,
418 whereas the trombone leads to the highest blend ratings for triads. Therefore, the horn
419 and trombone were chosen to represent instruments that blend well with others (colored
420 green) in the dyad and triad sets, respectively, as both brass instruments' spectral
421 descriptions also resemble each other. Furthermore, the oboe was chosen as the
422 exemplary instrument leading to poor blend (colored grey) in both sets.

423 Dyads

424 PLSR models predicting median blend ratings for dyads initially involved 46
425 regressors (X_{orig}). Elimination of loadings in P_{orig} that fall below the median threshold
426 yielded 23 regressors in X_{Q50} . As listed in Table 4, a three-PC model explains 93% of
427 the variance for X_{Q50} . Refining the regressors to an approximately orthogonal set, the
428 resulting X_{ortho} consists of 14 regressors, again, leading to a three-PC model explaining
429 93% of the variance. The model fit in y for X_{ortho} , displayed in Figure 3, shows the
430 variation in median blend ratings to be represented well. However, the blend ratings
431 (x-axis) exhibit two distinct groups of data points, corresponding to unison dyads
432 (circles) leading to substantially greater blend than non-unison dyads (diamonds).
433 Furthermore, non-unison dyads involving horn (green) yielded slightly greater blend
434 overall than those with oboe (grey), whereas no such distinction is observable for unison

435 dyads.

436 [— Insert Table 4 about here. —]

437 [— Insert Figure 3 about here. —]

438 Figure 4 visualizes the loadings P_{ortho} (vectors) and the scores T_{ortho} (symbols)
 439 across the first two PCs. Larger symbols for scores correspond to higher blend ratings.
 440 Likewise, longer vectors represent loadings that contribute more strongly, while the
 441 vector orientation illustrates which PCs the contribution primarily affects.

442 [— Insert Figure 4 about here. —]

443 Reflecting the main distinctions in median blend ratings, the scores T_{ortho} also
 444 form two distinct groups for unison and non-unison dyads, with the corresponding
 445 categorical variable C_{unison} describing this distinction most accurately along PC 1. The
 446 acoustical descriptor Δf_0 predicts the same distinction comparably well. PC 2 appears
 447 to be influenced by two factors: 1) an additional grouping of dyads based on low and
 448 high pitch levels, described by the categorical variable C_{pitch} and the acoustical
 449 descriptor $f_{0|ERB}$, and 2) a collinear set of spectral descriptors, falling slightly oblique to
 450 the PC axis. The distinction across interval types (horizontal) and pitch levels (vertical)
 451 yields four subgroups. Along each of these obliquely aligned groups the influence of
 452 spectral features appears to lead to similar dyad constellations.

453 Figure 5 suggests that the spectral and pitch influence is independent (orthogonal)
 454 on the plane spanning PCs 2 and 3. The spectral regressors involve several composite
 455 (Σ) as well as difference (Δ) measures for S_{centr}° and formant-related descriptors. With
 456 regard to the resulting scores, T_{ortho} yields a grouping of dyads into those containing
 457 either horn or oboe (green/low-left vs. grey/top-right), for both unison and non-unison
 458 dyads.

459 [— Insert Figure 5 about here. —]

460 Overall, the dyad data exhibit a complex structure of underlying factors, involving
 461 interval type, pitch level, and spectral features. Across all investigated models, their
 462 performance (R^2) is remarkably well matched by their predictive power (Q^2). Given the
 463 relatively large number of cases, $N = 180$, further PLSR analyses on subsets separated

464 by interval type are conducted, yielding $N = 60$ for unison and $N = 120$ for non-unison
 465 dyads. Separate analyses allow an assessment of whether certain spectral and pitch
 466 trends are specific to only one of the interval types.

467 **Unison.** A three-PC model on X_{Q50} involving 22 regressors leads to 46%
 468 explained variance in median blend ratings for unison dyads, exhibiting a substantially
 469 lower predictive power of only 17% explained variance. Due to a fairly wide variation in
 470 P_{Q50} orientations, the angular threshold ϕ_i determining X_{ortho} had to be increased to
 471 $|\phi_i| < 30^\circ$ to ensure that the reduction to an approximately orthogonal set would lead
 472 to a meaningful number of contributing regressors. The resulting model with nine
 473 regressors yields a two-PC model explaining 27% of the variance, which appears a more
 474 realistic estimate of the true predictive relationship between median blend ratings and
 475 X_{ortho} , as the discrepancy between model performance and the predictive power is
 476 substantially reduced.

477 As shown in Figure 6, the y_{unison} fit appears a closer fit to the diagonal than for
 478 the complete dyad data (Figure 3), but the blend ratings only span a relatively narrow
 479 scale range. This may result from a reduction in the perceptual resolution among the
 480 unison dyads due to the dominant distinction between unison and non-unison dyads.
 481 The reduced resolution also makes it more likely for the variation in median blend
 482 ratings to contain increased noise levels, supported by the large discrepancy between R^2
 483 and Q^2 in the initial models.

484 [— Insert Figure 6 about here. —]

485 PC 1 explains 22% of the variance and, as shown in Figure 7, appears to be linked
 486 to spectral composite (Σ) descriptors for main formant location (e.g., F_{max} , F_{3dB}) as
 487 well as centroid (e.g., S_{centr}°), which also distinguishes low register and high register
 488 instrument dyads (e.g., HB vs. OF). PC 2 accounts for another 5% of the variance,
 489 involving a distinction between instrument dyads with similar formant structure and
 490 those with divergent structures (e.g., HB vs. BF and HF), explained by the
 491 formant-related descriptors ΔF_{slope} and ΔF_{freq} .

492 [— Insert Figure 7 about here. —]

493 **Non-unison.** Twenty-three regressors in X_{Q50} yield a three-PC model
 494 explaining 55% of the variance in median blend ratings for non-unison dyads, with the
 495 predictive power corresponding to 47% of the variance explained. The reduction to
 496 X_{ortho} yields 11 regressors and a three-PC model explaining 48% of the variance, with a
 497 lower predictive power accounting for 35% of the variance. The model fit in $y_{non-unison}$
 498 for X_{ortho} , shown in Figure 8, improved compared to the one for the complete dyad set
 499 (Figure 3), showing a better approximation to the ideal fit (dashed line).

500 [— Insert Figure 8 about here. —]

501 As shown in Figure 9, PC 1 clearly reflects a grouping of dyads based on pitch
 502 level (C_{pitch} and $f_{0|ERB}$), accounting for 33% of the explained variance. At the same
 503 time, the composite of the spectral slope S_{slope}^{\sim} appears to covary with pitch change. All
 504 remaining spectral regressors appear relatively independent (orthogonal) to the pitch
 505 influence. Figure 10 illustrates that across the plane spanning PCs 2 and 3, two
 506 seemingly independent contributions of spectral regressors occur: 1) an implied triangle
 507 between the composite (Σ) regressors F_{3dB} , S_{centr}° , and the difference (Δ) descriptor
 508 F_{3dB} distinguishes dyads into those containing horn (bottom-left) and those involving
 509 oboe (top-right); 2) perpendicular to this orientation, difference in spectral slope S_{slope}^{\sim}
 510 and composite in noisiness S_{noise} contribute somewhat more weakly. Together, PCs 2
 511 and 3 account for 8% and 7% of the variance, respectively.

512 [— Insert Figure 9 about here. —]

513 [— Insert Figure 10 about here. —]

514 **Triads**

515 The PLSR analysis of triads first involved 61 regressors, which reduced to
 516 30 regressors in X_{Q50} , leading to a two-PC model explaining 89% of the variance in
 517 median blend ratings and with a predictive power explaining 73% of the variance. The
 518 subsequent reduction to X_{ortho} yields another two-PC model with 15 regressors that
 519 again explains 89% of the variance, notably, gaining in predictive power compared to
 520 the previous models. As shown in Figure 11, the model fit for y appears satisfactory,

521 given the smaller number of cases for triads ($n=20$). Still, a compact cluster involving
 522 *pizzicato* cello (squares, bottom-left) stands in contrast to more spread out ratings for
 523 sounds lacking them (circles, right half). A trend for triads involving trombone (green)
 524 to be the most blended is apparent in each subgroup.

525 [— Insert Table 5 about here. —]

526 [— Insert Figure 11 about here. —]

527 The main distinction found in Figure 12 along PC 1, which accounts for 85% of
 528 the variance, concerns the presence or absence of *pizzicato* cello sounds (the categorical
 529 variable C_{pizz}), with the acoustical difference in attack slopes A_{slope} performing similarly
 530 well. Apart from A_{slope} , the composite and difference descriptors for spectrotemporal
 531 incoherence $ST_{incoher}$ and noisiness S_{noise} are somewhat correlated with C_{pizz} . This
 532 could result from both the transient attack and rapid decay of the temporal envelope of
 533 *pizzicato* sounds contributing to more noise and more spectral change over time,
 534 respectively. In addition, the inclusion of two other spectral descriptors, difference in
 535 $F_{\Delta mag}$ and distribution in S_{skew} , could be explained by differences in spectral-envelope
 536 shape for the two articulations of the cello.

537 [— Insert Figure 12 about here. —]

538 PC 2 explains the remaining 3% of the variance, appearing to relate to the
 539 distribution (Ξ) of a number of formant descriptors. These comprise frequency
 540 measures of the main formant, F_{3dB} and F_{max} , measures of balance between main
 541 formants and the remaining spectral envelope, F_{slope} and $F_{\Delta mag}$, and the overall
 542 prominence of formant structure, F_{prom} . Furthermore, the relevance of two difference
 543 measures related to formants, F_{mag} and F_{prom} , suggests that the most pronounced
 544 differences among three descriptor values could also be of importance.

545 Discussion

546 Previous research has associated blend with acoustical measures describing
 547 spectral features, as well as temporal features like the attacks or onsets of sounds under
 548 certain circumstances. The current investigation pursued a correlational analysis using

549 PLSR, modeling two perceptual data sets involving dyads and triads. PLSR loadings
550 allowed us to evaluate the extent to which regressors were collinear or independent of
551 each other. This approach helped select the most effective regressors. Applied to the
552 complete data sets for both dyads and triads, the final models based on optimized
553 regressor sets explain around 90% of the variance in median blend ratings. Notably,
554 these levels of explained variance were still achieved after the elimination of
555 non-essential regressors, i.e., more than two thirds from the original set. The variation
556 in both data sets is best explained by a dominant factor that is unrelated to spectral
557 features.

558 For dyads, the distinction between unison and non-unison intervals explains 91%
559 of the variance, with the fundamental-frequency difference Δf_0 representing a reliable
560 acoustical predictor. That unison dyads would lead to higher blend than for non-unison
561 had been anticipated, given that similar effects have been found in other studies
562 (Kendall & Carterette, 1993; Lembke et al., in press). The pronounced difference
563 obtained in the current results, however, seems to exceed those previously reported,
564 which could be related to the current study being the only one in which unison and
565 non-unison were presented in a common stimulus set, whereas in other studies both
566 interval types had been grouped into separate experimental blocks (Kendall &
567 Carterette, 1993; Lembke et al., in press) or had even been tested in separate
568 experiments (Sandell, 1995).

569 In addition, even the second-most important factor in explaining the variation
570 among dyads, $f_{0|ERB}$, is unrelated to spectral features, as it reflects differences in pitch
571 height, accounting for 2% in all dyads and 33% when considering only the non-unison
572 dyads. The fact that the contribution of the pitch level is limited to the non-unison
573 dyads implies that it may not affect blend of unison dyads. For non-unison dyads, it is
574 also worth noting that inverting the assignment of instruments to the two pitches had
575 no effect on blend ratings. This negative finding goes counter to many claims in
576 orchestration treatises that the order of pitch assignment affects blend. It thus supports
577 the conclusion by Sandell (1995) that timbral inversion does not appear to influence

578 blend; only a single finding argues in its favor (Kendall & Carterette, 1993).

579 With regard to triads, the presence of a *pizzicato* cello evoked a strong decrease in
580 blend ratings, whereas even triads including cello sounds excited by a single, brisk bow
581 stroke led to substantially more blend. Again, this distinction had been anticipated,
582 given that increasingly impulsive sounds have been associated with comparable
583 decreases in blend (Tardieu & McAdams, 2012). Regarding the description of onset
584 articulations, the difference in attack slopes A_{slope} is strongly collinear with the
585 categorical distinction C_{pizz} , explaining about 85% of the variance; additional
586 collinearity with spectrotemporal or noise features can be assumed to co-occur as
587 byproducts of the abrupt changes in temporal envelopes.

588 With both data sets being dominantly influenced by pitch or temporal features
589 (e.g., attack), spectral descriptors only occur as secondary or even tertiary sources of
590 variation in the modeled blend ratings. In perceptual tasks comparable to those
591 employed in these experiments, participants may focus their attention on the dominant
592 distinctions across stimuli at the cost of perceptual resolution for the less pronounced
593 differences.

594 As the spectral factors likely only affected blend ratings in these regions of
595 reduced perceptual resolution, the possible role of behavioral noise needs to be
596 considered. Indeed, clear discrepancies between model performance R^2 and predictive
597 power Q^2 indicate that the initial PLSR models could have been overfitting to noise
598 artifacts instead of systematic factors of variation. For example, stripped of the
599 dominant factor, the unison and non-unison subsets of data account for no more than
600 50% of the variance (R^2). The unison-dyad data suggest that the true performance is
601 substantially lower as the predictive power is generally quite low and likely reflects
602 random variation or factors not captured by the tested regressors. In summary, the
603 identified tendencies for spectral regressors can be assumed to be valid for the obtained
604 proportions of explained variance, but they should be considered preliminary until
605 confirmed in additional datasets yielding greater resolution in the perceptual ratings.

606 Three spectral descriptors stand out in explaining the PLSR models for both data

607 sets, namely, the centroid of the pitch-generalized spectral envelope S_{centr}° and the two
608 main-formant descriptors F_{max} and F_{3dB} , notably representing spectral features that
609 have previously been found to be relevant (Lembke et al., in press; Lembke &
610 McAdams, 2015; Reuter, 1996; Sandell, 1995; Tardieu & McAdams, 2012). Differences
611 exist concerning the types of association between descriptor values of the instruments
612 constituting dyads or triads. For unison dyads, the composite (Σ) measures for all three
613 descriptors became relevant in explaining 22% of the variance, which is in agreement
614 with the same association explaining other perceptual results for unison dyads (Sandell,
615 1995; Tardieu & McAdams, 2012).

616 Non-unison dyads yield a more complex relationship and involved the composite
617 for S_{centr}° and F_{3dB} complemented by the difference in F_{3dB} , overall contributing 15% of
618 the variance. The relevance of the difference measure (Δ) is in agreement with the
619 absolute spectral-centroid difference having previously been reported as the strongest
620 predictor for non-unison dyads (Sandell, 1995). The particular combination of
621 composite and difference measures suggests that as S_{centr}° and F_{3dB} increased, so did the
622 divergence of F_{3dB} between the individual instruments, with both possibly contributing
623 to decreased blend. For instance, oboe paired with horn yields a higher composite
624 centroid due to the oboe's higher main formant, which at the same time increases the
625 frequency distance to the horn's low main formant, whereas for horn and bassoon, both
626 main formants are relatively low and, moreover, practically coincide in frequency.

627 The results for triads expand previous knowledge beyond dyadic contexts. Even if
628 spectral features only account for 3% of the variance, some new insight is gained from
629 the *distribution* (Ξ) of three descriptor values for several formant measures serving as
630 the strongest predictor, suggesting that relative position of the sound having an
631 intermediate descriptor value among all three sounds may indeed be useful in describing
632 instrument combinations with more than two instruments.

633 Overall, the global descriptor S_{centr}° and the main-formant location F_{3dB} indicate
634 that prominent spectral-envelope properties represent reliable correlates to blend across
635 various instruments, pitches, and polyphonic combinations. Being the first investigation

636 to test for the relevance of global and local spectral descriptors jointly, both domains
 637 seem helpful as regressors in a predictive application. Across all datasets, the descriptor
 638 loadings P confirmed that most spectral descriptors were partially correlated, at the
 639 same time, allowing the identification of descriptors that appeared independent of S_{centr}°
 640 and F_{3dB} , namely, the spectral slope, S_{slope}^{\sim} , and noisiness, S_{noise} (Figure 10), as well as
 641 the formant-based spectral slope, F_{slope} , and formant frequency difference between
 642 constituent sounds, F_{freq} (Figure 7). These additional descriptors could be of special
 643 interest in achieving more complete prediction models, although their relevance seems
 644 to depend on the stimulus context. A similar analysis approach on a wider data set is
 645 needed to confirm these trends, and possibly even to give further insight into the role of
 646 associations (Σ , Δ , Ξ) relevant for different musical scenarios. Furthermore, the
 647 apparent utility of pitch-generalized descriptors, i.e., all F descriptors and S_{centr}° as
 648 opposed to S_{centr}^{\sim} , implies that a case-by-case signal analysis on individual pitches may
 649 not be necessary, but instead, a prediction application could rely on a comprehensive,
 650 offline database of pitch-generalized instrument descriptions alone, which would
 651 significantly facilitate computation.

652 When considering the relative locations of instrument combinations along the PCs
 653 that correlate with spectral features, a recurring pattern of dyads or triads including
 654 oboe (grey), on the one side, opposed to combinations involving horn or
 655 trombone (green), on the other, becomes apparent. Dyads or triads containing oboe are
 656 often less blended, whereas combinations with horn or trombone (e.g., bassoon and
 657 horn, clarinet and horn, trombone and trombone) are among the most blended ones. If
 658 we consider the notion of *blendability* of a particular instrument, the oboe should be
 659 considered a poor ‘blender’, which can be explained spectrally by its prominent and
 660 unique formant structure. Similar observations linking oboe to poor blend have been
 661 made in previous perceptual investigations (Kendall & Carterette, 1993; Reuter, 1996;
 662 Sandell, 1995; Tardieu & McAdams, 2012) as well as ‘prescriptions’ found in
 663 orchestration treatises (Koechlin, 1954; Reuter, 2002). On the other hand, the horn is
 664 generally considered an easily blendable instrument, again reflected in perceptual results

665 (Reuter, 1996; Sandell, 1995). The relatively ‘dark’ timbre of the horn could support a
666 general hypothesis of lower centroids leading to more blend (Sandell, 1995), at the same
667 time supporting the argument that similar main-formant locations explain the good
668 blend obtained between horn and bassoon (Lembke et al., in press; Reuter, 1996).

669 In addition, Figure 12 illustrates that the distribution (Ξ) along formant
670 descriptors (F) distinguishes triads with two identical instruments (e.g., two trombones
671 plus clarinet, two *pizzicati* or two *arco* celli plus clarinet) from more diverse
672 combinations, without, however, directly reflecting how these combinations vary in
673 blend (visualized size of the symbols for scores predicted by the models). Nevertheless,
674 it does imply that timbral similarity, if not identity, aids blending. In summary, once
675 factors related to pitch intervals or onset articulation are taken into account, spectral
676 features do seem to represent the main underlying factor governing whether instrument
677 combinations blend or not, with pitch-generalized spectral descriptions possibly
678 conveying the timbral signature traits of instruments.

679 Conclusion

680 The present investigation shows that the perception of blended timbres in dyadic
681 and triadic contexts correlates with a number of acoustical factors. Analyses using
682 PLSR converged on an apparently reliable selection of independent predictors. The
683 importance of factors such as pitch interval type, pitch, and articulation (e.g., impulsive
684 vs. gradual note attack) became apparent. In addition, a group of spectral descriptors
685 that exhibit the strongest predictive abilities could be identified from a wide range of
686 descriptors, namely, the global spectral centroid and the upper frequency bound of main
687 formants, which may represent the relevant features informing instrumentation choices.
688 This wide variety of predictors suggests that in blend-prediction applications aimed at
689 realistic musical scenarios, all factors should be taken into account. Given an
690 appropriate acoustical characterization of instruments and details of how they are
691 combined and employed musically (e.g., in unison or non-unison, the articulation and
692 dynamic markings), these properties could suffice to predict the associated degree of

693 blend.

694 One main challenge for future research is determining the effective weighting
695 between these different factors of influence. Whether the clear dominance of interval
696 type or impulsiveness of attacks over spectral features, which became apparent in the
697 current investigation, would extend to more complex musical contexts remains to be
698 explored. It can be assumed that the growing complexity that a listening scenario
699 involving musical contexts would present, given the simultaneous presence of other
700 musical parameters, could significantly alter the relative importance of factors found in
701 listening experiments employing isolated dyadic or triadic stimuli.

702 For instance, a composer may assign a unison blend between two instruments to a
703 melodic voice while juxtaposing this against a chordal, non-unison accompaniment layer
704 whose instruments are chosen to blend amongst themselves into a homogeneous timbre.
705 On another level, the melody may become more distinct from the accompaniment due
706 to the distinction between unison and non-unison, which may also be desired. This case
707 scenario illustrates that blend-related factors need not stand in competition with each
708 other like they do in the investigated perceptual data, but instead could operate on
709 independent levels, fulfilling separate functions within the musical context.

710 For the composer, working with blend is not a matter of favoring unison intervals
711 over non-unison intervals, but being able to employ it at individual levels of the musical
712 scene (e.g., melody, accompaniment, or contrasting the two). Within each level, blend is
713 achieved by relying on the same principles, i.e., similarity in spectral description as well
714 as articulatory features (e.g., note attacks). This hypothetical scenario encourages
715 future work on blend-prediction models to rely on perceptual data obtained from
716 stimuli involving musical contexts (Kendall & Carterette, 1993; Lembke et al., in press;
717 Reuter, 1996), as it provides a more realistic setting from which weights between
718 blend-related factors could be estimated. We thus propose the need for a
719 meta-analytical investigation into a diverse range of perceptual blend data, in an
720 attempt to move toward generally applicable blend-prediction techniques.

721

Footnotes

722

723

724

725

¹MATLAB, Mathworks. The *plsregress* function from the Statistics and Machine Learning Toolbox was used. URL: <https://uk.mathworks.com/products/matlab.html>. Last accessed: August 21, 2017.

726

²URL: <http://vsl.co.at/>. Last accessed: August 21, 2017.

727

³See Appendix C in Lembke (2015) for details.

728

Supplemental Material

729

730

731

732

Representative examples for the dyad and triad stimuli are available online as supplemental material, which can be found as part of the online version of this article at <http://msx.sagepub.com>. Access the sound files through the hyperlink “Supplemental material”. Their filenames follow the naming convention found in Tables 1 and 2.

733

Author Note

734

735

736

737

738

739

The authors would like to thank Bennett K. Smith for programming the control software for both experiments and Emma Kast for the recruitment and running of participants for the experiment involving triads. We also thank two anonymous reviewers and the editor for their valuable feedback on an earlier version of this article. The findings reported in this article were conducted as part of the first author’s doctoral research and are also featured in his thesis (Lembke, 2015).

740

741

742

Correspondence concerning this article should be directed to Sven-Amin Lembke, currently at De Montfort University, Clephan Building, Room 00.07b, Leicester LE1 9BH, United Kingdom; email: sven-amin.lembke@dmu.ac.uk.

743

Funding

744

745

This research was partly funded by an ACN CREATE undergraduate research award to Kyra Parker as well as a Canadian National Sciences and Engineering

⁷⁴⁶ Research Council grant (RGPIN 312774-2010) and a Canada Research Chair to Stephen
⁷⁴⁷ McAdams.

References

748

- 749 de Jong, S. (1993). SIMPLS: An alternative approach to partial least squares
750 regression. *Chemometrics and Intelligent Laboratory Systems*, 18(3), 251–263.
- 751 Eerola, T., Lartillot, O., & Toiviainen, P. (2009, October). Prediction of
752 multidimensional emotional ratings in music from audio using multivariate
753 regression models. In *Proc. 10th International Society of Music Information*
754 *Retrieval (ISMIR) Conference* (pp. 621–626). Kobe, Japan.
- 755 Fant, G. (1960). *Acoustic theory of speech production*. The Hague: Mouton.
- 756 Geladi, P., & Kowalski, B. R. (1986). Partial least-squares regression: A tutorial.
757 *Analytica Chimica Acta*, 185, 1–17.
- 758 Giordano, B. L., Rocchesso, D., & McAdams, S. (2010). Integration of acoustical
759 information in the perception of impacted sound sources: The role of information
760 accuracy and exploitability. *Journal of Experimental Psychology: Human*
761 *Perception and Performance*, 36(2), 462–476.
- 762 Horner, A. B., Beauchamp, J. W., & So, R. H. Y. (2009). Detection of time-varying
763 harmonic amplitude alterations due to spectral interpolations between musical
764 instrument tones. *Journal of the Acoustical Society of America*, 125(1), 492–502.
- 765 ISO 389–8. (2004). *Acoustics: Reference zero for the calibration of audiometric*
766 *equipment—Part 8: Reference equivalent threshold sound pressure levels for pure*
767 *tones and circumaural earphones* (Tech. Rep.). Geneva, Switzerland:
768 International Organization for Standardization.
- 769 Kendall, R. A. (2004, August). Musical timbre in triadic contexts. In *Proc. 8th*
770 *International Conference on Music Perception and Cognition (ICMPC)* (pp.
771 600–602). Evanston, Illinois.
- 772 Kendall, R. A., & Carterette, E. C. (1993). Identification and blend of timbres as a
773 basis for orchestration. *Contemporary Music Review*, 9(1), 51–67.
- 774 Kendall, R. A., & Vassilakis, P. (2006). Perceptual acoustics of consonance and
775 dissonance in multitimbral triads. *Journal of the Acoustical Society of America*,
776 120(5), 3276–3276.

- 777 Kendall, R. A., & Vassilakis, P. N. (2010). Perception and acoustical analyzes of
778 traditionally orchestrated musical structures versus non-traditional counterparts.
779 *Journal of the Acoustical Society of America*, 128(4), 2344–2344.
- 780 Koechlin, C. (1954). *Traité de l'orchestration: En quatre volumes (Treatise of*
781 *orchestration: In four volumes)*. Paris: M. Eschig.
- 782 Kumar, S., Forster, H. M., Bailey, P., & Griffiths, T. D. (2008). Mapping
783 unpleasantness of sounds to their auditory representation. *Journal of the*
784 *Acoustical Society of America*, 124(6), 3810–3817.
- 785 Lembke, S.-A. (2015). *When timbre blends musically: perception and acoustics*
786 *underlying orchestration and performance* (PhD thesis). McGill University.
- 787 Lembke, S.-A., Levine, S., & McAdams, S. (in press). Blending between bassoon and
788 horn players: An analysis of timbral adjustments during musical performance.
789 *Music Perception*, 35(2).
- 790 Lembke, S.-A., & McAdams, S. (2015). The role of spectral-envelope characteristics in
791 perceptual blending of wind-instrument sounds. *Acta Acustica united with*
792 *Acustica*, 101(5), 1039–1051.
- 793 Luce, D., & Clark, J. (1967). Physical correlates of brass-instrument tones. *Journal of*
794 *the Acoustical Society of America*, 42(6), 1232–1243.
- 795 Luce, D. A. (1975). Dynamic spectrum changes of orchestral instruments. *Journal of*
796 *the Audio Engineering Society*, 23(7), 565–568.
- 797 Martens, H., Høy, M., Westad, F., Folkenberg, D., & Martens, M. (2001). Analysis of
798 designed experiments by stabilised PLS Regression and jack-knifing.
799 *Chemometrics and Intelligent Laboratory Systems*, 58(2), 151–170.
- 800 Martin, F., & Champlin, C. (2000). Reconsidering the limits of normal hearing.
801 *Journal of the American Academy of Audiology*, 11(2), 64–66.
- 802 Moore, B. C., & Glasberg, B. R. (1983). Suggested formulae for calculating
803 auditory-filter bandwidths and excitation patterns. *Journal of the Acoustical*
804 *Society of America*, 74(3), 750–753.
- 805 Peeters, G., Giordano, B. L., Susini, P., Misdariis, N., & McAdams, S. (2011). The

- 806 Timbre Toolbox: Extracting audio descriptors from musical signals. *Journal of*
807 *the Acoustical Society of America*, 130(5), 2902–2916.
- 808 Reuter, C. (1996). *Die auditive Diskrimination von Orchesterinstrumenten -*
809 *Verschmelzung und Heraushörbarkeit von Instrumentalklangfarben im*
810 *Ensemblespiel (The auditory discrimination of orchestral instruments: Fusion and*
811 *distinguishability of instrumental timbres in ensemble playing)*. Frankfurt am
812 Main, Germany: P. Lang.
- 813 Reuter, C. (2002). *Klangfarbe und Instrumentation: Geschichte–Ursachen–Wirkung*
814 *(Timbre and instrumentation: History–causes–effect)*. Frankfurt am Main,
815 Germany: P. Lang.
- 816 Reuter, C. (2009). The role of formant positions and micro-modulations in blending
817 and partial masking of musical instruments. *Journal of the Acoustical Society of*
818 *America*, 126(4), 2237–2237.
- 819 Rumsey, F., Zieliński, S., Kassier, R., & Bech, S. (2005). On the relative importance of
820 spatial and timbral fidelities in judgments of degraded multichannel audio quality.
821 *Journal of the Acoustical Society of America*, 118(2), 968–976.
- 822 Sandell, G. J. (1995). Roles for spectral centroid and other factors in determining
823 “blended” instrument pairings in orchestration. *Music Perception*, 13(2), 209–246.
- 824 Schumann, K. E. (1929). *Physik der Klangfarben (Physics of timbres)* (professorial
825 dissertation). Universität Berlin.
- 826 Tardieu, D., & McAdams, S. (2012). Perception of dyads of impulsive and sustained
827 instrument sounds. *Music Perception*, 30(2), 117–128.
- 828 Wold, S., Sjöström, M., & Eriksson, L. (2001). PLS-regression: A basic tool of
829 chemometrics. *Chemometrics and Intelligent Laboratory Systems*, 58(2), 109–130.

Table 1

Fifteen dyads across pairs of the six investigated instruments.

Dyad	Instrument pair	
HB	horn	bassoon
HO	horn	oboe
HT	horn	trumpet
HC	horn	clarinet
HF	horn	flute
BO	bassoon	oboe
BT	bassoon	trumpet
BC	bassoon	clarinet
BF	bassoon	flute
OT	oboe	trumpet
OC	oboe	clarinet
OF	oboe	flute
TC	trumpet	clarinet
TF	trumpet	flute
CF	clarinet	flute

Table 2

Twenty triads and their constituent instruments and assigned pitches.

Triad	Instruments & pitches		
	C4	F4	Bb4
AAF	cello (<i>arco</i>)	cello (<i>arco</i>)	flute
AAC	cello (<i>arco</i>)	cello (<i>arco</i>)	clarinet
PPC	cello (<i>pizz.</i>)	cello (<i>pizz.</i>)	clarinet
PPO	cello (<i>pizz.</i>)	cello (<i>pizz.</i>)	oboe
PAF	cello (<i>pizz.</i>)	cello (<i>arco</i>)	flute
PAO	cello (<i>pizz.</i>)	cello (<i>arco</i>)	oboe
ACF	cello (<i>arco</i>)	clarinet	flute
AOF	cello (<i>arco</i>)	oboe	flute
ACO	cello (<i>arco</i>)	clarinet	oboe
PCO	cello (<i>pizz.</i>)	clarinet	oboe
TTF	trombone	trombone	flute
TTC	trombone	trombone	clarinet
TTO	trombone	trombone	oboe
TCO	trombone	clarinet	oboe
PTT	cello (<i>pizz.</i>)	trombone	trombone
PAT	cello (<i>pizz.</i>)	cello (<i>arco</i>)	trombone
ATF	cello (<i>arco</i>)	trombone	flute
ATC	cello (<i>arco</i>)	trombone	clarinet
PTC	cello (<i>pizz.</i>)	trombone	clarinet
PTO	cello (<i>pizz.</i>)	trombone	oboe

Table 3

Acoustical descriptors investigated for dyads and/or triads (marked by ‘x’ in the rightmost columns), related to the global spectrum (S), formants (F), the attack portion of the temporal envelope (A), spectrotemporal variation (ST), as well as categorical variables (C). Descriptor values for individual sounds forming dyads or triads were associated with a single regressor value by difference Δ , composite Σ , distribution Ξ (triads only) or as specified otherwise.

Abbrev.	Description	Unit	Association	Dyad	Triad
S_{centr}^{\sim}	spectral centroid ^a	Hz	Δ, Σ, Ξ	x	x
S_{centr}°	spectral centroid ^b	Hz	Δ, Σ, Ξ	x	x
S_{slope}^{\sim}	spectral slope ^a	Hz ⁻¹	Δ, Σ, Ξ	x	x
S_{slope}°	spectral slope ^b	Hz ⁻¹	Δ, Σ, Ξ	x	x
S_{skew}	spectral skew ^a	-	Δ, Σ, Ξ	x	x
S_{kurtos}	spectral kurtosis ^a	-	Δ, Σ, Ξ	x	x
S_{spread}	spectral spread ^a	Hz	Δ, Σ, Ξ	x	x
S_{roll}	spectral roll-off ^a	Hz	Δ, Σ, Ξ	x	x
$S_{decrease}$	spectral decrease ^a	-	Δ, Σ, Ξ	x	x
S_{noise}	noisiness ^a	-	Δ, Σ, Ξ	x	x
F_{max}	main-formant maximum ^b	Hz	Δ, Σ, Ξ	x	x
F_{3dB}	main-formant upper bound ^b	Hz	Δ, Σ, Ξ	x	x
F_{slope}	spectral slope above main formant ^b	Hz ⁻¹	Δ, Σ, Ξ	x	x
$F_{\Delta mag}$	level difference F_1 vs. above ^b	dB	Δ, Σ, Ξ	x	x
F_{prom}	formant prominence ^b	-	Δ, Σ, Ξ	x	x
F_{freq}	formant frequency deviations ^b	Hz	Δ	x	x
F_{mag}	formant magnitude deviations ^b	dB	Δ	x	x
A_{time}	attack time	s	Δ, Ξ	x	x
$A_{log(time)}$	log. attack time	s	Δ, Ξ	x	x
A_{slope}	attack slope	s ⁻¹	Δ, Ξ	x	x
ST_{flux}	spectral flux ^a	-	Δ, Σ, Ξ	x	x
$ST_{incoher}$	spectral incoherence ^a	-	Δ, Σ, Ξ	x	x
C_{unison}	unison or non-unison	-	binary	x	
Δf_0	f_0 difference	Hz	Δ	x	
C_{pitch}	pitch level	-	binary	x	
$f_0 _{ERB}$	f_0 , auditory scaling	ERB ^c rate	C4 or G4	x	
$C_{interval}$	interval type	-	ternary	x	
$C_{position}$	instrument positions	-	binary	x	
C_{pizz}	including <i>pizzicato</i> or not	-	binary		x
x_{mix}	amplitude mix or balance	-	scaled value	x	x

Note.

^a S^{\sim} based on signal analysis for individual pitches.

^b S° based on pitch-generalized spectral-envelope estimate.

^cERB: equivalent rectangular bandwidth (Moore & Glasberg, 1983).

Table 4

Performance (R^2) and predictive power (Q^2) of the PLSR model for dyads as well as component-wise contribution along the first three PCs for the three stages X_{orig} , X_{Q50} , X_{ortho} involving a sequential reduction of the number of regressors m .

y dyads	X regressors	m	R^2	Q^2	PC 1	PC 2	PC 3
all	X_{orig}	46	.94	.91	.88	.04	.01
	X_{Q50}	23	.93	.92	.90	.03	<.01
	X_{ortho}	14	.93	.93	.91	.02	<.01
unison	X_{orig}	44	.56	.18	.33	.14	.10
	X_{Q50}	22	.46	.17	.26	.12	.09
	X_{ortho}	9	.27	.16	.22	.05	-
non-unison	X_{orig}	45	.60	.40	.42	.10	.08
	X_{Q50}	23	.55	.47	.39	.14	.03
	X_{ortho}	11	.48	.35	.33	.08	.07

Table 5

Performance (R^2) and predictive power (Q^2) of the PLSR model for triads as well as component-wise contribution along up to three PCs. Three stages X_{orig} , X_{Q50} , X_{ortho} involve a sequential reduction of the number of regressors m .

y triads	X regressors	m	R^2	Q^2	PC 1	PC 2	PC 3
	X_{orig}	61	.90	.64	.86	.04	-
all	X_{Q50}	30	.89	.73	.84	.05	-
	X_{ortho}	15	.89	.76	.85	.03	-

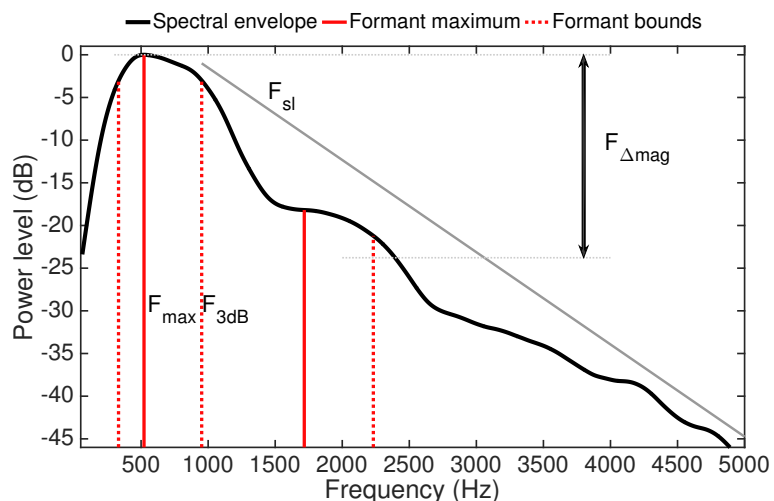


Figure 1. Pitch-generalized spectral envelope of a horn with identified frequencies for formant maxima and 3 dB bounds (see Lembke & McAdams, 2015). F_{max} and F_{3dB} characterize the maximum and upper bound, respectively, for the dominant, lower *main* formant. F_{slope} represents the spectral slope above the main formant. $F_{\Delta mag}$ quantifies the magnitude difference between F_{max} and the average magnitude above F_{3dB} .

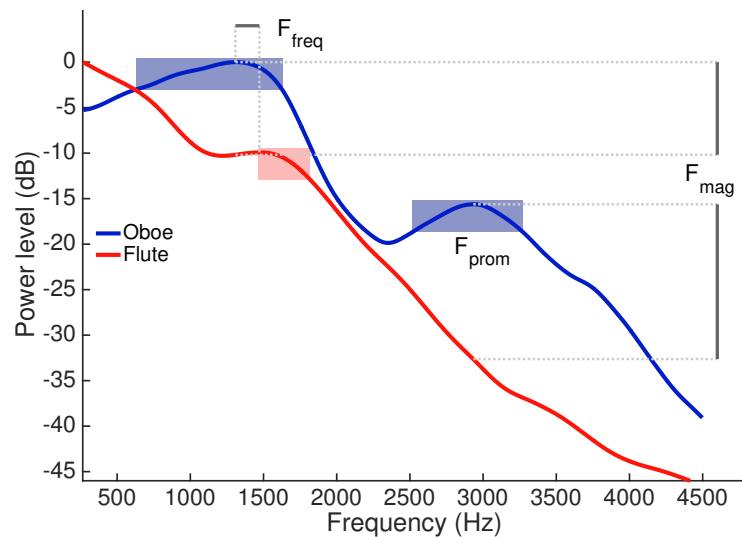


Figure 2. Pitch-generalized spectral envelopes of oboe (blue) and flute (red). F_{prom} characterizes the existence and clarity of formant features (e.g., 3 dB formant bounds); the larger the total shaded area, the more prominent the instrument's formant structure. F_{mag} evaluates the total magnitude difference between spectral envelopes at formant frequencies. F_{freq} quantifies the deviation between formant frequencies.

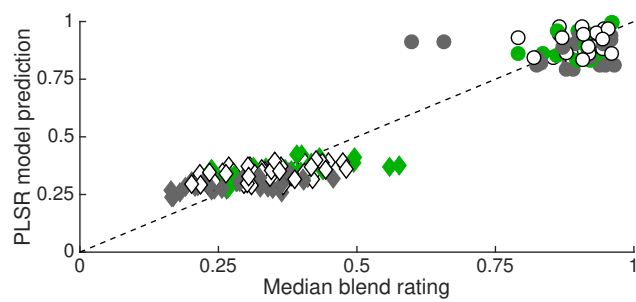


Figure 3. Dyad model fit of y variables for X_{ortho} . Legend: circles, unison; diamonds, non-unison; grey involves oboe; green involves horn (excl. HO).

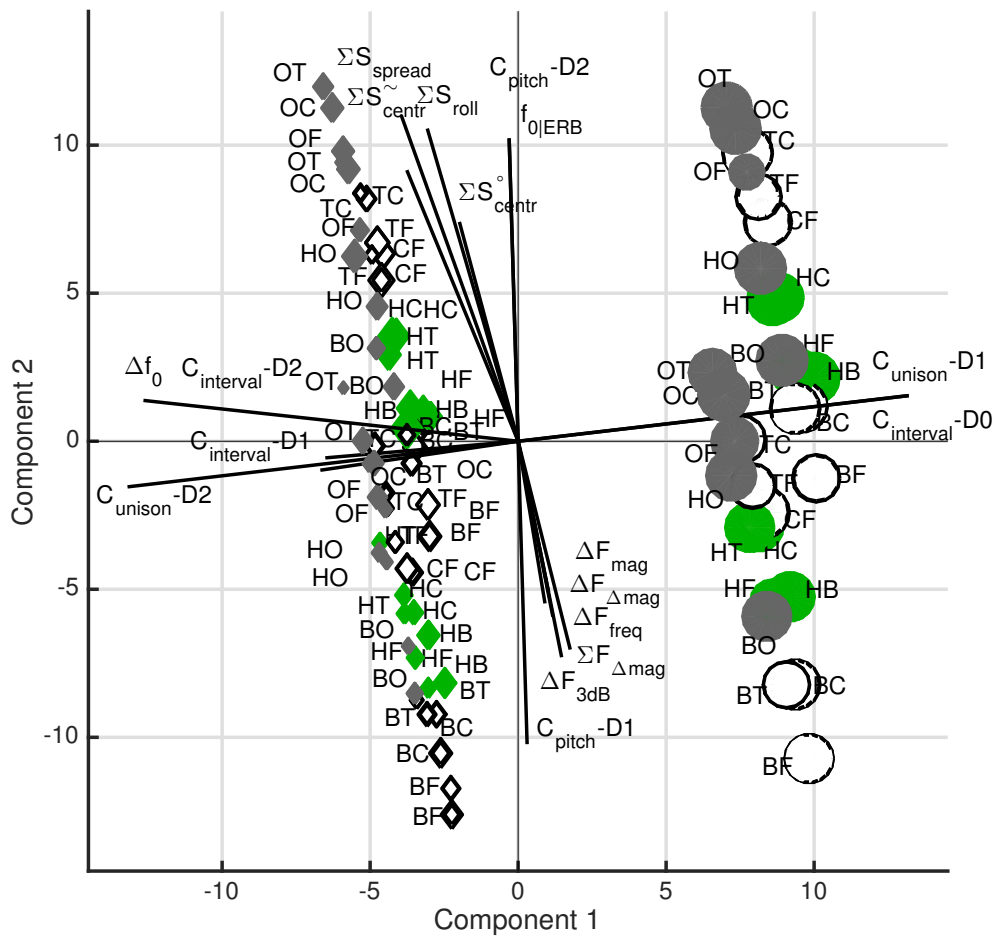


Figure 4. PLSR loadings P_{ortho} (vectors) and scores T_{ortho} (symbols) for PCs 1 and 2 with dyads. Legend: circles, unison; diamonds, non-unison; their size represents the relative degree of blend; dark grey involves oboe; green involves horn (excl. HO).

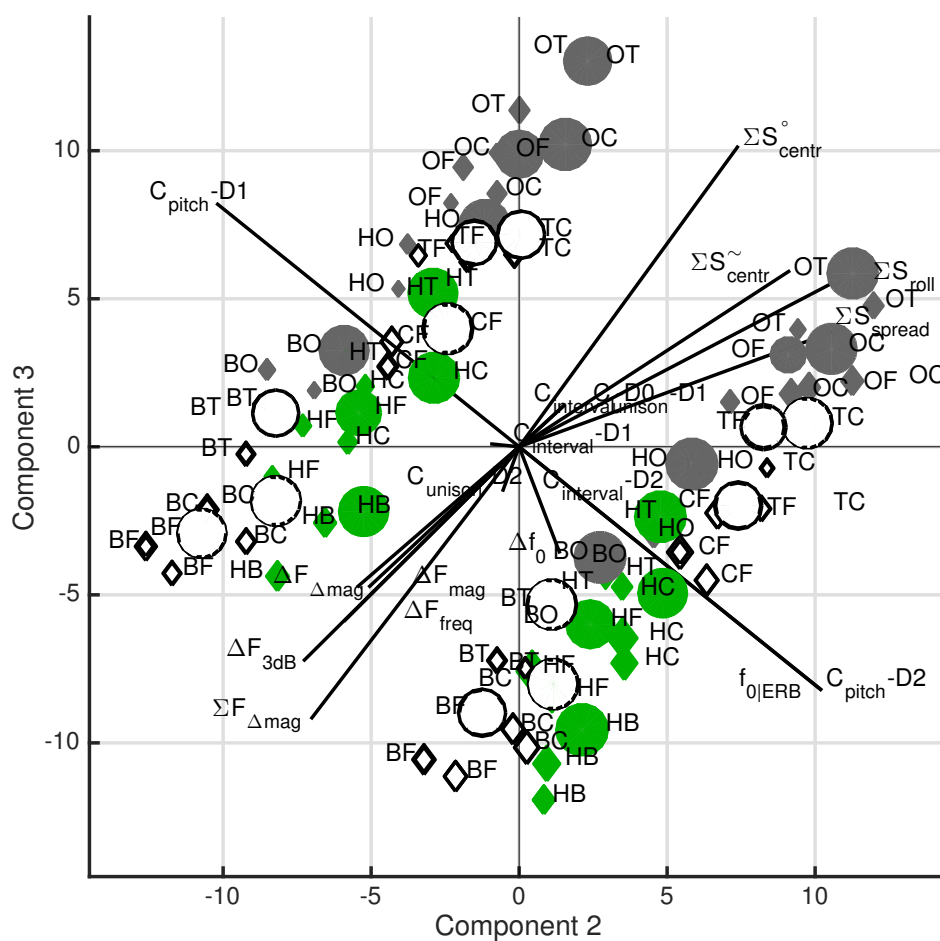


Figure 5. Dyad P_{ortho} and T_{ortho} for PCs 2 and 3. See Figure 4 for legend.

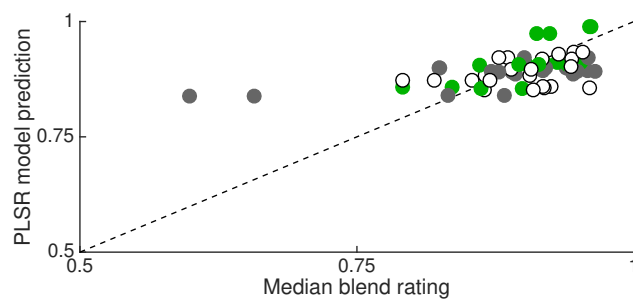


Figure 6. Unison-dyad model fit of y variables for X_{ortho} . See Figure 3 for legend.

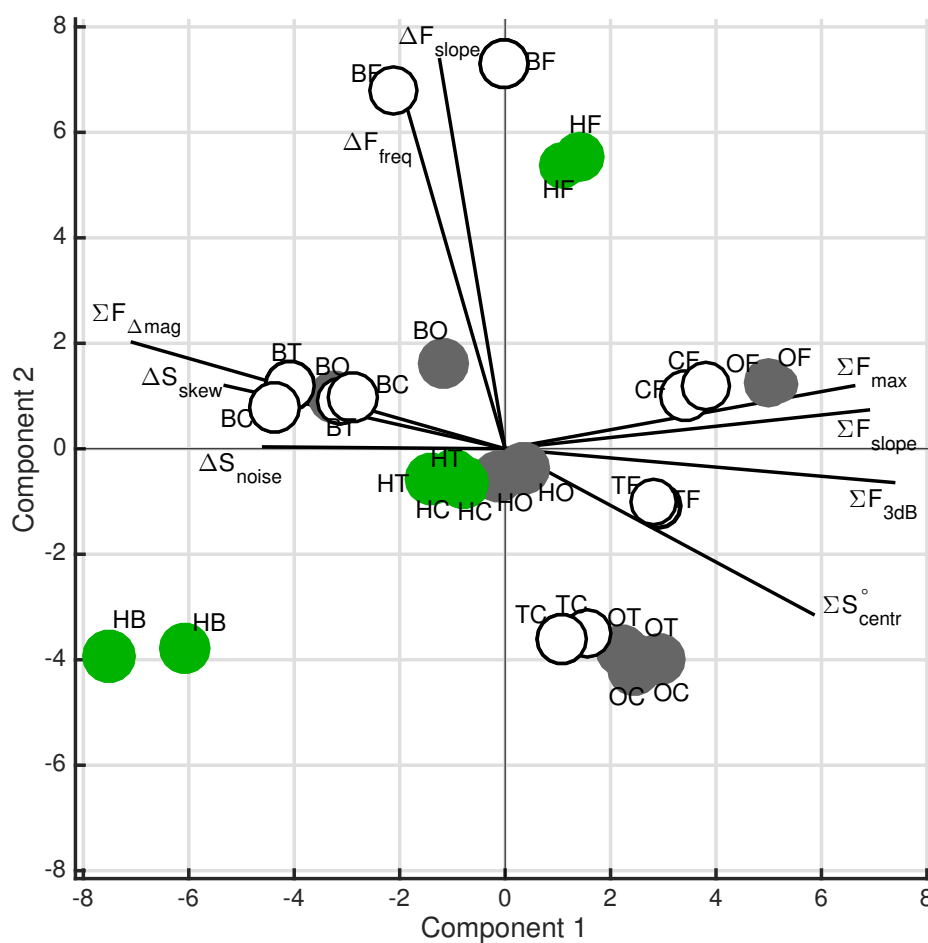


Figure 7. Unison-dyad P_{ortho} and T_{ortho} for PCs 1 and 2. See Figure 4 for legend.

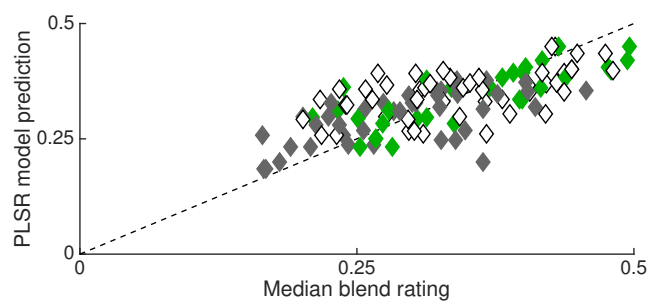


Figure 8. Non-unison-dyad model fit of y variables for X_{ortho} . See Figure 3 for legend.

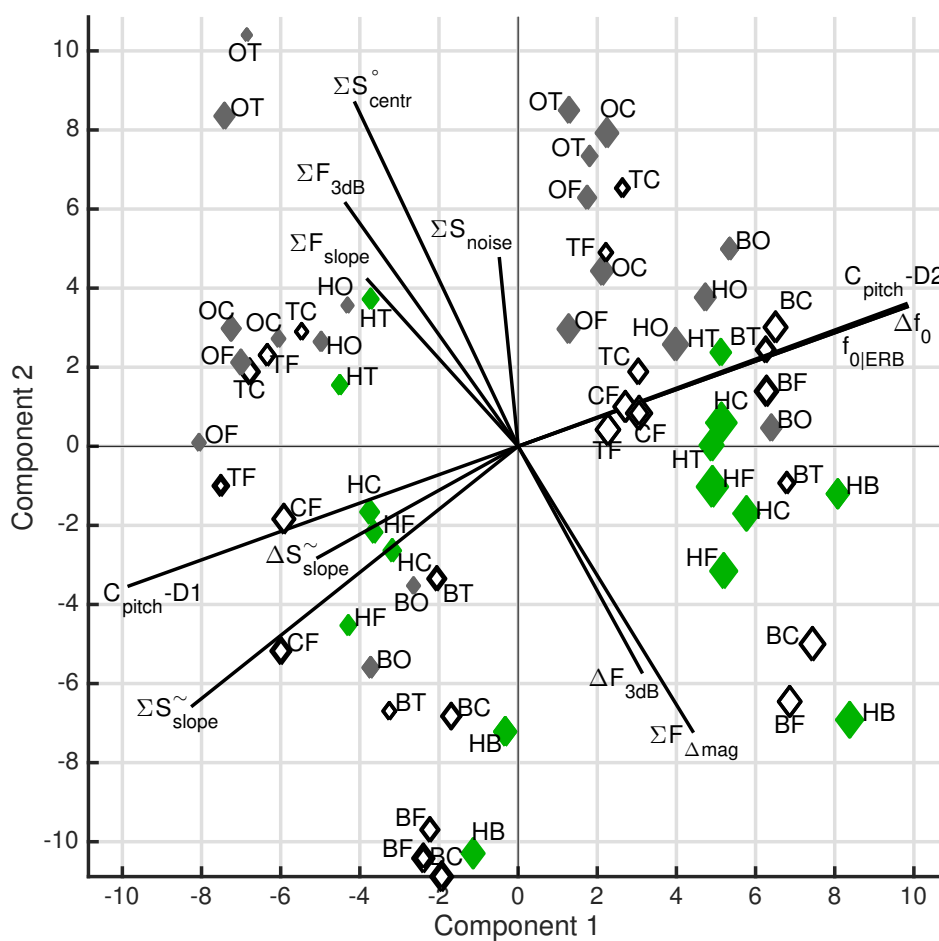


Figure 9. Non-unison-dyad P_{ortho} and T_{ortho} for PCs 1 and 2. See Figure 4 for legend.

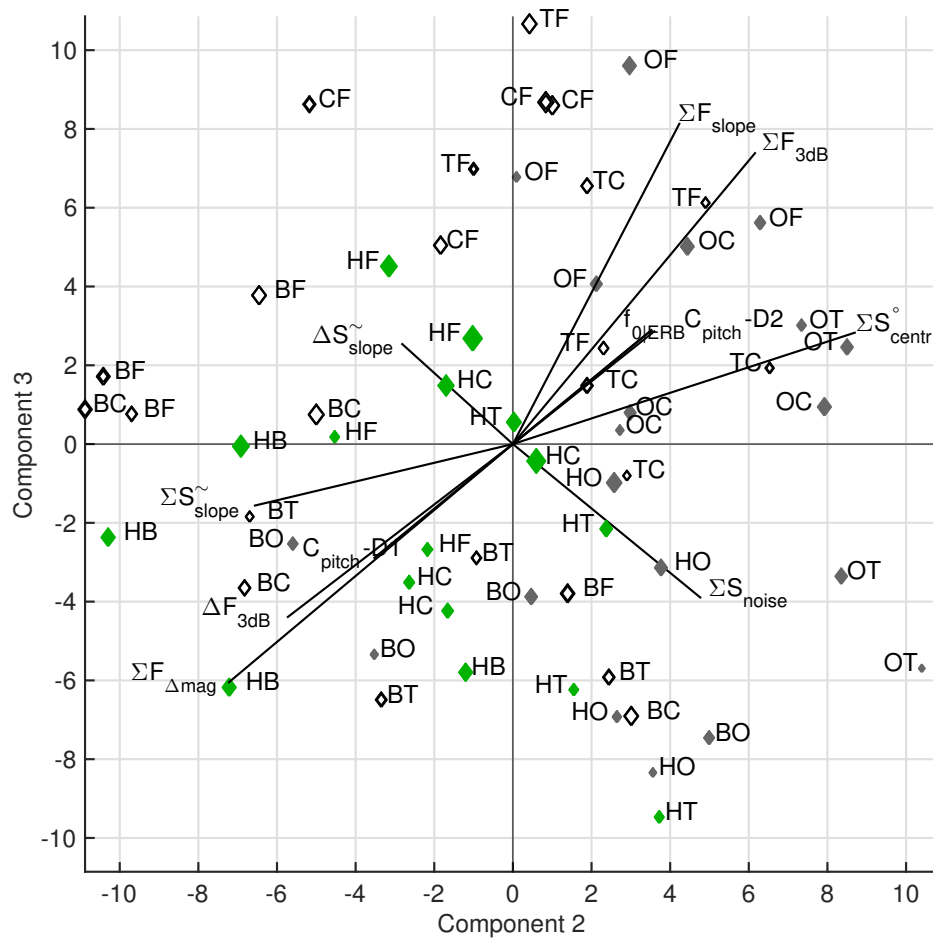


Figure 10. Non-unison-dyad P_{ortho} and T_{ortho} for PCs 2 and 3. See Figure 4 for legend.

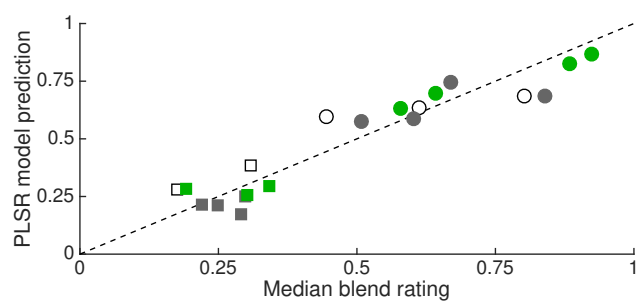


Figure 11. Triad model fit of y variables for X_{ortho} . Legend: squares, including *pizzicati*; circles, excluding *pizzicati*; grey involves oboe; green involves trombone (excl. PTO, TTO, TCO).

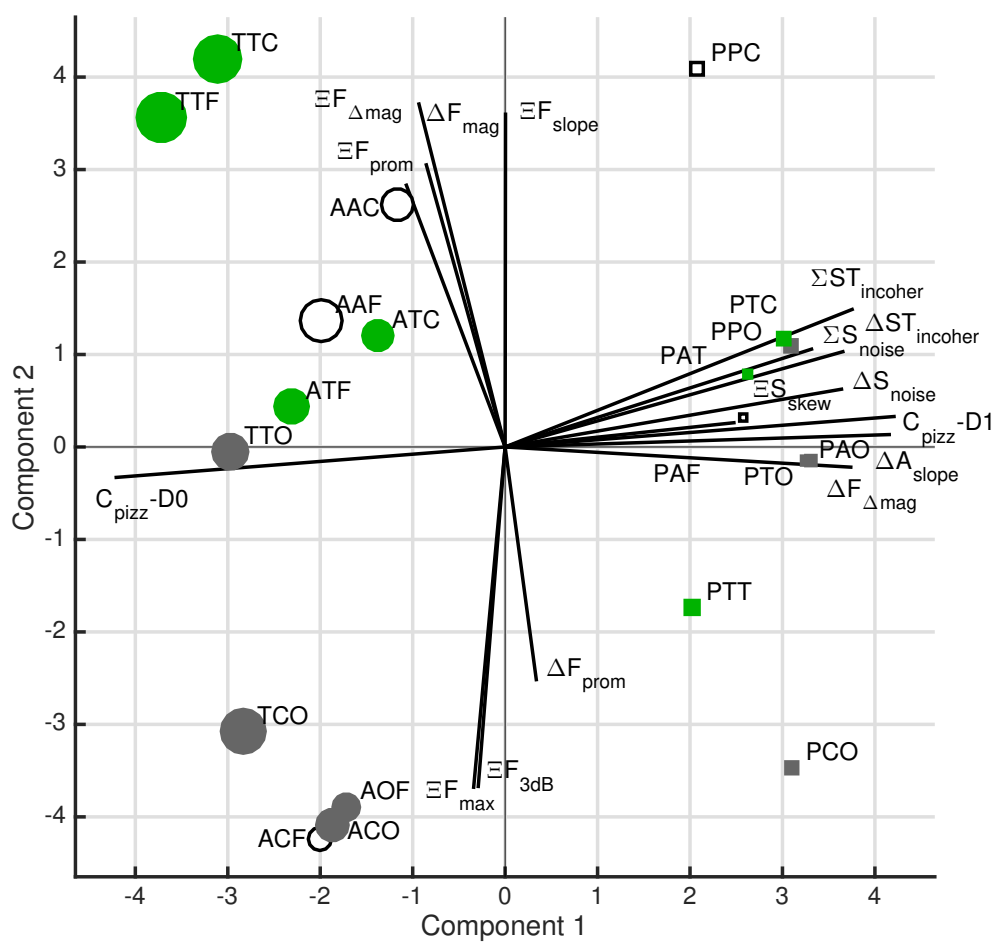


Figure 12. PLSR loadings P_{ortho} (vectors) and scores T_{ortho} (symbols) for PCs 1 and 2 with triads. Legend: squares, including *pizzicati*; circles, excluding *pizzicati*; symbol size represents relative degree of blend; grey involves oboe; green involves trombone (excluding PTO, TTO, TCO).