



HAL
open science

Covert variations of a musician's loudness during collective improvisation capture other musicians' attention and impact their interactions

Armand Schwarz, Arthur Faraco, Coralie Vincent, Patrick Susini, Emmanuel Ponsot, Clément Canonne

► To cite this version:

Armand Schwarz, Arthur Faraco, Coralie Vincent, Patrick Susini, Emmanuel Ponsot, et al.. Covert variations of a musician's loudness during collective improvisation capture other musicians' attention and impact their interactions. *Proceedings of the Royal Society B: Biological Sciences*, 2025, 292 (2039), pp.20242623. 10.1098/rspb.2024.2623 . hal-04905196

HAL Id: hal-04905196

<https://hal.science/hal-04905196v1>

Submitted on 22 Jan 2025

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

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



Distributed under a Creative Commons Attribution 4.0 International License

1 Covert variations of a musician's loudness during collective improvisation capture other mu-
2 sicians' attention and impact their interactions

3
4

5 Armand Schwarz^{1,2}, Arthur Faraco³, Coralie Vincent¹, Patrick Susini¹, Emmanuel Ponsot¹ and
6 Clément Canonne^{1*}

7
8

1- STMS UMR 9912 (CNRS/IRCAM/Sorbonne Université), Paris, France

9
10

2- Institut Pasteur (Université Paris Cité/INSERM/Institut de l'Audition/IHU
Reconnect), Paris, France

11
12

3- Universidade de Sao Paulo, Sao Paulo, Brazil

13
14

* corresponding author: clementcanonne@hotmail.com

15 **Funding.** This work was supported by the INSEAD-Sorbonne Université Behavioural Lab
16 and by the Collegium Musicae.

17
18

18 **Acknowledgments.** The authors thank all the wonderful improvisers who participated in this
19 study.

20
21

21 **Data availability statement.** All data collected for the study and code used in the analyses are
22 available at <https://osf.io/bsgx8/>.

23
24

24 **Authors contributions.** All authors designed the study. A.S. prepared the experimental mate-
25 rial and programmed the software used in the experiment. C.C., A.S., A.F., and C.V. collected
26 the data. C.C., E.P., and A.S. prepared the analysis plan. A.S. analyzed the data. A.F. ran the
27 cross-recurrence quantification analyses. C.C. wrote the paper with edits and comments from
28 E.P.

29
30

30 **Ethical approval and consent.** Ethical approval for this study was obtained at INSEAD/ Sor-
31 bonne University Center for Behavioural Science, Paris, France. All methods were carried out
32 in accordance with their guidelines and regulations. All participants signed an informed con-
33 sent.

34

35 Covert variations of a musician's loudness during collective improvisation capture other mu-
36 sicians' attention and impact their interactions

37

38 0. Abstract

39 While research on auditory attention in complex acoustical environment is a thriving field,
40 experimental studies thus far have typically treated participants as passive listeners. The pre-
41 sent study – which combined real-time covert loudness manipulations and online probe detec-
42 tion –investigates for the first time the effects of acoustic salience on auditory attention during
43 live interactions, using musical improvisation as an experimental paradigm. We found that
44 musicians were more likely to pay attention to a given co-performer when this performer was
45 made sounding louder or softer; that such salient effect was not due to the local variations
46 introduced by our manipulations but rather likely to be driven by the more long-term context;
47 and that improvisers tended to be more strongly and more stably coupled when a musician
48 was made more salient. Our results thus demonstrate that a meaningful change of the acousti-
49 cal context not only captured attention but also impacted the ongoing musical interaction it-
50 self, highlighting the tight relationship between attentional selection and interaction in such
51 social scenario, and opening novel perspectives to address whether similar processes are at
52 play in human linguistic interactions.

53

54 1. Introduction

55

56 Research on auditory attention in complex acoustical environment is a thriving field. Starting
57 with Cherry's seminal paper (1953) – which coined the expression “cocktail party effect” to
58 describe our ability to focus on selected aspects of a given acoustic scene while blocking out
59 non-relevant sonic streams – an impressive amount of psychophysical as well as neuroimag-

60 ing research has been conducted on both sides of the “cocktail party problem” (McDermott,
61 2009): how do we segregate concurrent sonic streams that are somehow mixed?; and how do
62 we direct our attention to a source of interest while ignoring other sources?
63 Strikingly, the whole field was thus started by the description of an imaginary situation – the
64 cocktail party – in which agents *participate* to the acoustic scene in which they are immersed
65 in. They are not merely passive listeners, observing the scene from afar; rather, they direct
66 their attention from one source to another as part of broader dialogic interests – namely, iden-
67 tifying the conversation(s) they are going to engage with, thus adding their own voice to the
68 overall acoustic scene. However, this interactional aspect has all but disappeared from extent
69 psychophysical studies addressing the cocktail party problem from the perspective of speech
70 intelligibility with multiple concurrent talkers. Baring a few recent exceptions of a rather ob-
71 servational nature (e.g., Miles et al., 2023 or Ryan et al., 2023), it is fair to say that most stud-
72 ies approach speech-on-speech situations from an experimental perspective in which the par-
73 ticipants are treated as external listeners of a scene rather than as agents *involved* in that scene
74 (see Bidelman & Yoo, 2020 for a recent example). While there are obvious methodological
75 reasons for this state of affair (studying participants’ auditory attention in sound-proof booths
76 using precisely calibrated sounds presented through headphones or spatialized speakers al-
77 lows for maximum experimental control), disconnecting attention from interaction might also
78 come at a cost. This is particularly clear when it comes to the problem of selective auditory
79 attention – whose underlying mechanisms are arguably less studied than the mechanisms al-
80 lowing for smooth sound segregation (McDermott, 2009). Selective auditory attention is in-
81 deed often seen as being strongly impacted by bottom-up factors – often described by the ge-
82 neric term “salience” to refer to acoustical events that draw auditory attention independently
83 of volition, either because they possess acoustic characteristics that are known to catch audito-
84 ry attention such as high levels of sound intensity or roughness (Arnal et al., 2019) or because

85 they stand out from the local sonic context (Dalton & Lavie, 2004). However, one might
86 wonder whether stimulus-driven factors such as acoustic salience still play a significant role
87 in shaping auditory attention over time in interactional contexts in which volitional attentional
88 selection processes might seem to be predominant (Koch et al., 2011). Moreover, it could be
89 the case that agents are typically less sensitive to local variations in the acoustic context in
90 settings which are more cognitively demanding (as interactional settings typically are, com-
91 paratively to more passive settings), in line with already established effects of cognitive load
92 on sound evaluation or auditory susceptibility (see Steffens et al., 2020; Van der Heiden et al.,
93 2023). This would constitute yet another case of a finding on social attention obtained from
94 studies in which participants merely observe an interaction between agents (for example by
95 relying on pictures or videos) that does not straightforwardly generalize to participants' be-
96 haviors during live interactions (Dawson & Foulsham, 2022).

97 The present paper precisely seeks to investigate the impact of such acoustic salience on audi-
98 tory attention during live interactions, by introducing a novel experimental paradigm involv-
99 ing musicians freely improvising together in a trio setting. Music provides a highly ecological
100 setting – and yet still largely unexplored (but see Faraco et al., 2024, for a recent example) –
101 to study selective auditory attention in interactive contexts, as music is, most often than not, a
102 collective affair (Blacking, 1973), presenting simultaneously different instrumental or vocal
103 parts. In performances, this means that musicians must navigate within a complex acoustic
104 environment, focusing their attention on such or such aspect of the sonic tapestry collectively
105 produced to fulfill various performance goals – which typically means focusing their auditory
106 attention on such or such musician (or sub-group of musicians) (Keller 2001). Amongst the
107 wide variety of collective musical practices, free improvisation is a particularly interesting
108 case for our purpose, if only because joint improvisation plays such a central role in most of
109 our social interactions (Noy et al., 2011). In collective free improvisation (CFI from now on),

110 musicians aim at spontaneously creating music without relying on pre-defined plans or pre-
111 existing musical structures (Saint-Germier & Canonne 2020). CFI is thus in stark contrast
112 with score-based, well-rehearsed performances, in which the distribution of musicians' audi-
113 tory attention is likely to be mostly guided by the structural information contained within the
114 score (e.g., where the melodic part lies) or the "ideal sound" and other performance goals that
115 are slowly built through rehearsals. In CFI, on the contrary, how musicians distribute their
116 auditory attention is an integral part of the interactional dynamics, and plays a crucial role in
117 shaping how the performance will unfold (Clarke, 2005).

118 But while CFI might provide a highly relevant setting to assess whether acoustic salience play
119 a significant role in driving musicians' auditory attention, a crucial problem is to find a way to
120 track such auditory attention over the course of the performance. Recent advances in auditory
121 attention decoding, based on the modeling of neural information collected through mobiles
122 EEG (Straetmans et al., 2022) offer promising perspectives, but we are still far from being
123 able to use such methods in a setting as sonically and interactionally complex as collective
124 free improvisation. Relying on post-hoc verbalizations through interviews with musicians
125 (Seddon, 2005) is of course also an option but such verbalizations tend to focus on conscious
126 decisions made by the musicians (Canonne & Garnier, 2012), which might make the method-
127 ology ill-suited to study non-intentional shifts of musicians' auditory attention. In the follow-
128 ing study, we thus relied on a behavioral approach, using a probe-based method which al-
129 lowed us to assess, each time a probe is sent, whether or not a given musician in a trio was
130 paying close auditory attention to one of the two other musicians. This probe-based method
131 was then combined with real-time manipulations of the salience of the acoustic signals pro-
132 duced by the musicians. As they were playing together, and unknowingly to them, the musi-
133 cians' individual signals were selectively made, for short periods of time of a few seconds,
134 either a bit louder or a bit softer – a manipulation which presents the advantage of being ap-

135 plicable in the same way to every musical instrument, contrary to manipulations that would
136 involve, e.g., real-time manipulation of the spectral information. This novel experimental par-
137 adigm thus made it possible to systematically study, for the first time, whether and how a
138 change in loudness – which has been shown in previous studies to act as a salient cue (Dalton
139 & Lavie, 2004) – could shape selective auditory attention in a context in which participants
140 are able to freely interact with one another. On the one hand, contrasting the effects of both
141 “louder” and “softer” manipulations on musicians’ attention allowed us to assess whether this
142 loudness manipulation was more a matter of intrinsic acoustic cues (which would predict that
143 louder events are more salient than softer events) or of contextual cues (which would predict
144 that both louder and softer events can act as salient cues when they stand out of given context,
145 as was found in Dalton & Lavie, 2004). On the other hand, our repeated loudness manipula-
146 tions combined with our probe-based methodology – with probes being sent to participants
147 both *during* the manipulations and at various points *after* the manipulations – allowed us to
148 assess whether such potential salience effect would be rather driven by the short-term devia-
149 tion in loudness or more long-term contextual changes in musicians’ typical loudness.
150 But beyond studying the impact of loudness changes on musicians’ attention, our paradigm
151 also allowed us to directly assess whether such events would impact the ongoing interaction.
152 Attention and action have indeed been found to be closely related (Humphreys et al., 2010), to
153 the point that a recent influential philosophical account of attention has explicitly defined at-
154 tention in terms of action selection (Wu, 2014). More specifically, studies have shown that
155 acoustic salience within musical stimuli was associated with spontaneous muscle activity
156 from the participants exposed to such stimuli (Schultz et al., 2021), and that joint attention
157 increased feelings of connectedness between participants (Wolf et al., 2015). All this makes
158 plausible that musicians would be more likely to interact more strongly with their co-
159 improvisers when they are more acoustically salient, and our experimental design made it

160 possible to test for this hypothesis. In other words, by assessing its impact on both attention
161 and interaction, our study investigates how acoustic salience can act as a full-fledged commu-
162 nicational strategy in social settings.

163

164 2. Methods

165

166 *2.1. Participants*

167 Fifteen musicians participated in the experiment (mean age = 40.5 years, SD = 9.5, 9 male, 2
168 female, 1 non-binary and 3 that did not provide the information), divided into five trios. They
169 were highly trained musicians (with a mean of 28.8 years of musical practice, SD = 8.6 years)
170 and had a significant experience with collective free improvisation (with a mean of 19 years
171 of practice, SD = 9.5), which they all practiced in a professional proficiency (invitation to
172 important festivals; numerous live and studio recordings published on well-regarded labels;
173 etc.). They were recruited from the Parisian Free Improvisation scene, that comprises a high
174 number of musicians with a wide variety of different stylistic backgrounds (Roueff, 2006).

175 The overall instrumentation was saxophone ($N = 4$), guitar ($N = 3$), trumpet ($N = 2$), drums,
176 piano, clarinet, contrabass clarinet, double bass and electronics ($N = 1$). One participant also
177 used their voice during the improvisations. In each trio, the three musicians always played
178 different instruments in order to minimize timbral – and thus source – confusion for the
179 participants.

180 Trios were intentionally composed to minimize potential effects of familiarity on selective
181 attention (i.e., avoiding that participants paid more attention to the musicians they were used
182 to play with). Consequently, the majority of musicians within each trio were unacquainted or
183 had not previously collaborated together. Participants assessed their prior familiarity with the

184 other members of their trio using a 7-point Likert scale. As expected, mean familiarity was
185 low ($M = 2.06$; $SD = 0.93$).

186 All participants gave their informed written consent and were compensated at the standard
187 rate for the employment of professional musicians in France.

188

189 2.2. Procedure

190 The experiment was held in a professional recording studio. Each musician of the trio was
191 allocated to an individual booth, and was equipped with headphones (Beyerdynamics DT 770
192 pro, 80 ohms) in order to listen to each other. Musicians could thus not see each other, so that
193 communication between them would only be based on acoustic information. Importantly,
194 musicians' headphones were panned in such a way that they heard one improviser completely
195 on the right side and one completely on the left side, while hearing themselves in the center.
196 The panning (i.e., which musician is heard on the right side, which is heard on the left side)
197 was made randomly.

198 Each trio performed between four and six improvisations so as to improvise for a total of
199 approximately 40 minutes (mean = 42 minutes, $SD = 2.5$ minutes). They had a 30-minute
200 break after having performed the first two or three improvisations, depending on the duration
201 of those first improvisations. The recording was made on *Protools* (version 2022.9) by a
202 professional sound engineer.

203 While the musicians were playing, we introduced at various points real-time manipulations of
204 the RMS level (i.e., average loudness) of a given musician during approximately a 10-
205 second window. There were two distinct patterns of variation. Taking the musician's actual
206 loudness as the baseline, the *louder* pattern featured an +5 dB increase in the signal's
207 amplitude for a duration comprised between 1.76s and 3.13s (randomly drawn from a
208 gaussian distribution with mean=2.5s and $SD=0.25$), sustained the heightened level for a

209 duration comprised between 4.29s and 5.69s (randomly drawn from a gaussian distribution
210 with mean=5s and SD=0.25), and then decreased for a duration comprised between 1.85s and
211 3.27s (randomly drawn from a gaussian distribution with mean=2.5s and SD=0.25) back to
212 the baseline. Conversely, the *softer* pattern started with a -5 dB decrease from the baseline for
213 approximately 2.5s (drawn from the same distribution as for the louder pattern), maintained
214 the reduced level for approximately 5s, and finally increased for approximately 2.5s back
215 to the baseline.

216 As these manipulations were repeated numerous times over a given performance, the
217 variations of durations described above were meant to make the manipulations more
218 unpredictable, and thus less recognizable. Moreover, since the musicians did not see each
219 other while playing, and that their interactions was entirely mediated by the sounds perceived
220 in their headphones, they had no way to know that the level heard in their headphones was not
221 the one at which the musicians were actually playing. The 5 dB increase/decrease was pre-
222 tested by the authors of the present paper, in order to identify a threshold value that would be
223 both perceivable but would not sound too “unnatural” (see Supplementary Material for the
224 presentation of a perceptual follow-up study which confirmed, in a post-hoc fashion, the
225 naturalness of our manipulations).

226 In order to implement these manipulations, six different 15-minute automation tracks were
227 created for each trio, each of them defining the moments in which these RMS manipulations
228 would occur for each musician. For each musician A in a given trio, 2 automation tracks were
229 created: one applied to B’s signal as received in A’s headphones and one applied to C’s signal
230 as received in A’s headphones. Note that musicians always heard themselves as they actually
231 played, i.e., without any RMS manipulation. A *Max/MSP* patch was designed to synchronize
232 in real time those automation tracks with the ongoing *Protools* recording session.

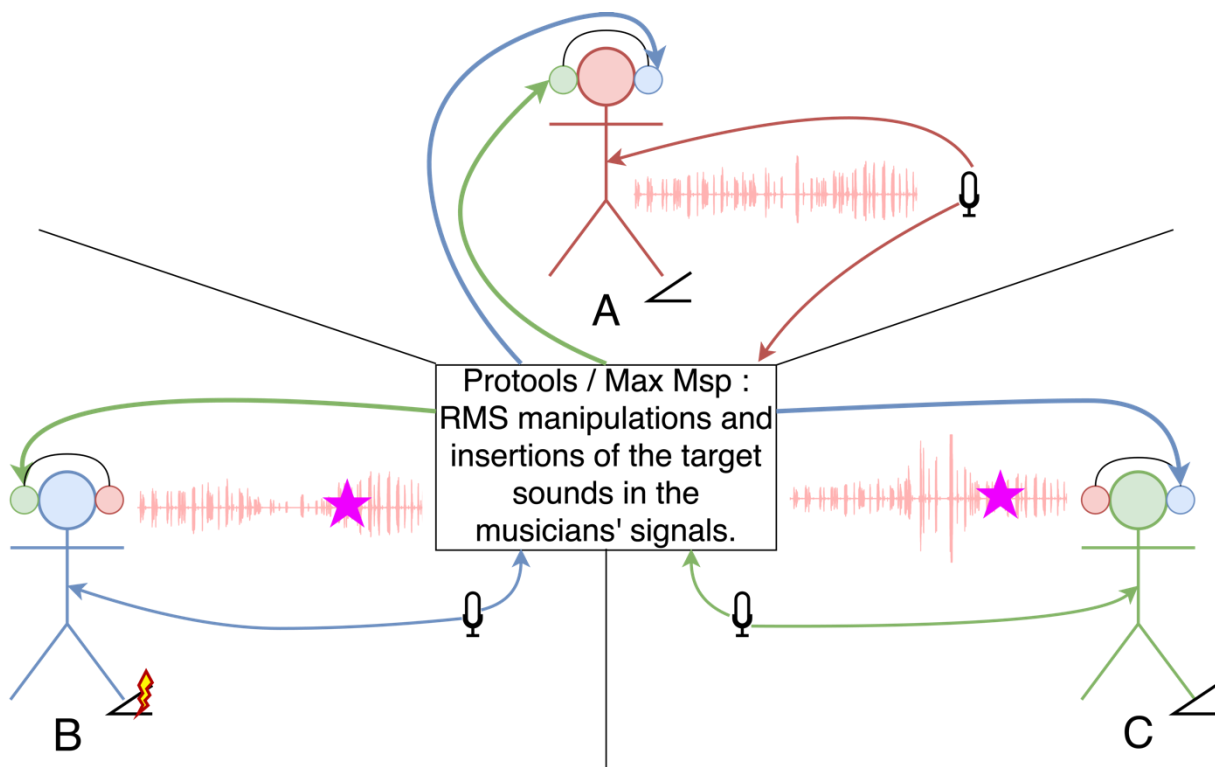
233 The following four constraints were used in the random generation of the automation tracks
234 for a given trio to minimize the impact of our manipulations on the ongoing group interaction,
235 while still collecting as many data points as possible. First, a given musician never heard both
236 of their co-improvisers modified at the same time. Second, for a given musician, RMS
237 manipulations were never applied to a same co-improviser twice in a row. Third, if a given
238 musician's RMS was being modified for their two co-improvisers, then one co-improviser
239 heard that musician through a louder pattern, and the other through a softer pattern. Fourth,
240 and finally, when considering the 6 automation tracks of a given trio as a whole, RMS
241 manipulations occurred continuously, with a 1-second buffer period between two successive
242 manipulations.

243 To track the auditory attention of the participant in real time, participants had to perform a
244 detection task while improvising. Target sounds were indeed added at different times to the
245 music improvised by the musicians and the musicians were asked to press a MIDI pedal
246 (ROLAND DP-10 Piano style Sustain Pedal) each time they noticed one of those target
247 sounds – thus triggering a signal recorded as a separate MIDI track in the Protools session.
248 Musicians were free to place their pedal wherever they felt it would be more convenient for
249 them to press it while playing. At the beginning of each session, a 3-minute training period
250 was organized so that musicians could familiarize themselves with the target sounds as well
251 as with the pedal pressing.

252 Each target sound was associated to a given musician, meaning that the two other musicians
253 heard the target sound on the channel associated to that musician. Musicians also never heard
254 the target sounds associated with themselves. For example, every time a target sound was
255 associated to musician A, the two other musicians B and C heard the target sound on the side
256 attributed to musician A in their headphones panning (i.e., either the left or the right side)
257 while A did not hear it (see Figure 1).

258 The target sound was composed of two bursts of noise of 60ms each, with 40ms of attack and
 259 20ms of release and without any inter-onset interval, so as to make detection harder. The first
 260 burst was filtered between 400 and 1600 Hz and the second one between 320 Hz and 1280
 261 Hz, with a strong decrease outside of the passing band (~ -300 dB/octave). Importantly, the
 262 level of the target sound, was always 1.5 dB louder than the level of the musician's track at
 263 the moment at which it was embedded. This was implemented through a Max/MSP patch
 264 which performed a real-time analysis of the RMS level of all the recorded tracks during the
 265 performance. The choice of this value for the emergence level of the target sound was made
 266 on the basis of observations from the experimenters during a pilot session in which several
 267 values were tested, trying to ensure sufficient audibility (which was also assessed in a post-
 268 hoc fashion, see section 2.3 below) while limiting salience, and thus avoiding both floor and
 269 ceiling effects regarding target detection performance.

270



271

272 Figure 1. Experimental procedure. Musicians A, B and C are improvising together, without seeing each other
 273 and hearing each other through a panned mixed in their headphones. B hears A's signal through a "softer

274 pattern” real-time manipulation (see the smaller amplitude) and C hears A’s signal through a “louder pattern”
275 real-time manipulation (see the larger amplitude). A target sound (represented here by a purple star) is associated
276 to musician A during these manipulations. B and C hear the target sound on the side on which they hear
277 musician A (i.e., the left channel for musician B and the right channel for musician C). B detected the target
278 sounds and pressed the pedal while C did not detect the target sound and did not press the pedal.

279

280 Target sounds appeared every 7 to 14 seconds (randomly drawn from a uniform distribution).
281 Since musicians did not hear the target sounds associated to their own musical track, this
282 means that, on average, a musician heard a target sound every 15.75 seconds. For each trio,
283 three 15-minute sound files were created beforehand with all the target sounds associated to
284 each of the three musicians. Each sound file was then superposed to the two different
285 automation tracks associated with that musician (one for each of their co-improvisers within
286 the trio) so that roughly 75% of the target sounds would be heard at moments where the level
287 of the musician associated with the target sound was manipulated. The remaining target
288 sounds were heard when no RMS manipulation was happening and were thus used as a
289 baseline. On average, each participant heard 137 target sounds ($SD = 20$) over the duration of
290 the whole experiment. This resulted in a total of 2170 target sounds to be detected by the
291 musicians.

292 Importantly, in order to decrease potential learning effects, target sounds and RMS
293 manipulations were not systematically associated and the distribution of the time intervals
294 between two successive target sounds heard by a same participant was quite wide (see
295 Supplementary Material for further analyses on potential learning effects). Furthermore, the
296 first 30 seconds of each improvisation were without target sounds nor RMS modification, so
297 that the musical interaction between the three musicians would begin in a way as natural as
298 possible.

299 Two video examples of the improvisations produced during this experiment by two different
300 trios, and edited in such a way as to visually illustrate our experimental procedure, can be
301 seen in the Supplementary Material.

302 Ethical approval for this study was obtained at [Anonymized] (Protocol ID: 2023-16). All
303 methods were carried out in accordance with their guidelines and regulations.

304

305 *2.3. Estimating the target sounds' perceptual emergence*

306 Since the music was entirely improvised, the difficulty to detect a target sound was likely to
307 vary depending on the moment at which it appeared, regarding not only the acoustical charac-
308 teristics of the track in which it was embedded, but also that of the other tracks. We assumed
309 that the latter aspect was minimal in the present context, since the tracks were separated in the
310 mix and therefore only subject to contralateral masking effects. However, regarding the for-
311 mer aspect, the masking effects related to the track associated with the target sound could
312 have influenced detection. Indeed, even though target sounds were always inserted in a track
313 using a constant level offset of +1.5 dB, their perceptual emergence is likely to have depended
314 on more complex acoustical parameters of that specific track (e.g., average level, spectral con-
315 tent of the track, etc.). Thus, to ensure that the observed performance to detect these target
316 sounds would mainly reflect the degree of attention to the track with which it was associated,
317 and not simply their perceptual emergence at a pure psychoacoustical level, we estimated
318 such perceptual emergence based on a well-adopted computational model of the human audi-
319 tory system (King et al., 2019). This model allowed us to compute a perceptual emergence
320 score for each one of the target sounds triggered during the experiment (see Supplementary
321 Material for more details), and thus to assess whether our experimental manipulation had an
322 effect on the detection of target sounds, over and beyond their perceptual emergence.

323

324 2.4. Acoustical features

325 As the musicians were recorded in separate studio booths, we had access to each musician's
326 individual tracks. This allowed us to calculate audio descriptors of interest and explore
327 potential relationships between these descriptors and our salience manipulations. Specifically,
328 we computed two main audio descriptors often used to account for interaction between the
329 musicians in collective free improvisation: Root Mean Square (RMS) and spectral centroid
330 (Goupil et al., 2021). RMS is indicative of the loudness in each musician's signal, while the
331 spectral centroid provides information on the signal's timbre (more specifically on its
332 brightness). These two audio descriptors were computed using the python library *Librosa*
333 (McFee et al., 2015) from each musician's individual WAV files, with a 100ms window size.

334

335 2.5. Variables

336 The target sounds presented to the performers were divided into three categories:

- 337 - "louder" target sounds: the target sounds associated to a musician heard through a
338 louder pattern;
- 339 - "softer" target sounds: the target sounds associated to a musician heard through a
340 softer pattern;
- 341 - baseline target sounds: the target sounds associated to a musician heard without any
342 RMS manipulation.

343 This provided us with our main independent variable: RMS manipulations (base-
344 line/softer/louder).

345 To account for the effects of RMS manipulations on attention, a target sound was considered
346 to be detected by a given musician if the musician's pedal was pressed in the 7 seconds fol-
347 lowing the triggering of the target sound (7 seconds was the minimal temporal distance be-

348 tween two successive target sounds).¹ This provided us with our first dependent variable: De-
349 tection (yes/no). Note that the false alarms (i.e., pressing the pedal in the absence of a target
350 sound) were rare (mean number of false alarms per participant = 5.133, SD = 3.631), amount-
351 ing to less than 3% of the total number of target sounds participants had to detect. Since there
352 was moreover no clear method for assigning them to a given experimental category (i.e., base-
353 line, “louder pattern”, or “softer pattern”), they were ultimately not taken into account in our
354 analysis.

355 In order to assess whether the potential effects of RMS manipulations on attention were due to
356 a change in the local acoustic context, we were also interested in investigating whether target
357 sounds that would occur immediately after a RMS manipulation (and thus, immediately after
358 a change back to the baseline level) would be better detected than target sounds that would
359 occur farther away from a RMS manipulation. Baseline target sounds were thus further divid-
360 ed into two sub-categories: Baseline 1 – containing the target sounds occurring less than 10
361 seconds after the end of a RMS manipulation (i.e., during a time window of the same length
362 as our patterns); and Baseline 2 – containing the target sounds occurring more than 10 se-
363 conds after the end of a RMS manipulation.

364 We were finally interested in assessing whether RMS manipulations would impact the
365 musicians’ interactions. Since all of our target sounds were associated with one of our three
366 experimental conditions, we relied on the timings of such target sounds to estimate the
367 interaction between a given pair of musicians (the musician who had to detect the target sound
368 and the musician with who the target sound was associated). For each musician, we thus
369 extracted the time series corresponding to our two audio descriptors (RMS and spectral
370 centroid) over a 7-second window (the minimum temporal distance between two target
371 sounds, so that there would never be any overlap in successive measurements) centred around

¹ See Supplementary Material for a replication of our main analysis using a much shorter time window but yielding similar results.

372 the target sound under consideration. Note that 14 target sounds appeared less than 3.5s
373 before the end of a performance and were thus discarded from our analyses as they did not
374 allow us to extract complete time series following the procedure described above. To assess
375 the interaction between any two improvisers, we then used cross-recurrence quantification
376 analysis (CRQA) – a method that provides insights into various correlational characteristics
377 between two or more time series and that has been widely used in joint action studies (Wallot
378 & Leonardi, 2018).

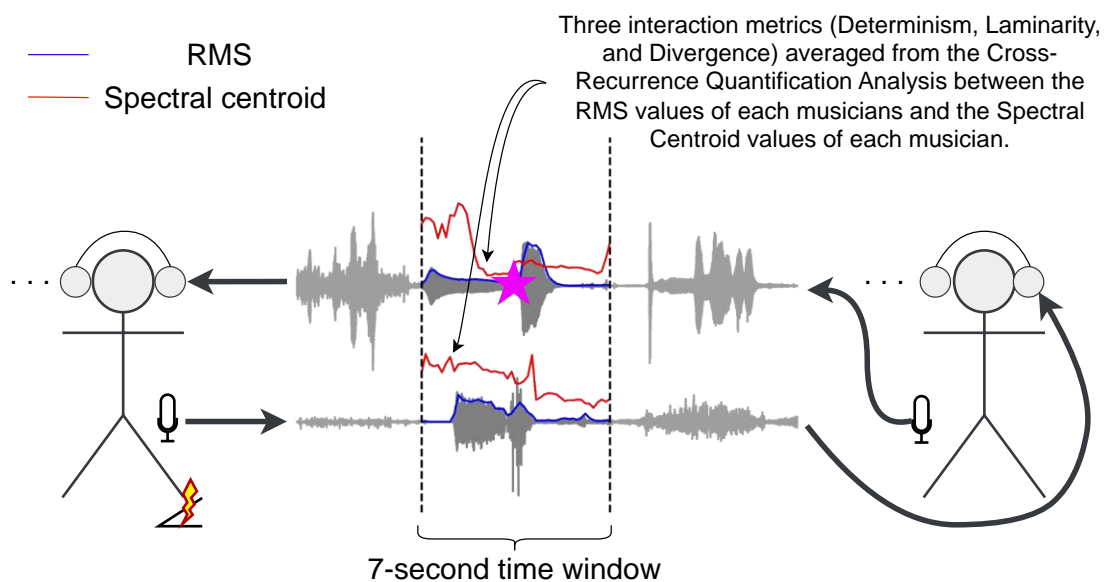
379 After identifying the relevant hyperparameters (see Supplementary Material for more details),
380 we performed the CRQA using the PyRQA library in Python (Rawald, Sips & Marwan, 2017)
381 on the extracted pairs of RMS time series and pairs of Spectral Centroid time series. Three
382 metrics were used as variables to describe the interaction between the musicians, namely:

383 - *Determinism* (DET), which corresponds to the proportion of recurrent points that
384 forms diagonal lines in the cross-recurrence plot. Diagonal lines represent periods
385 where the two time-series follow similar trajectories for an extended period, meaning
386 that there are patterns in one time-series that are also found in the other time-series.
387 This provides an estimation of the degree to which the two musicians had similar
388 behaviors.

389 - *Laminarity* (LAM), which measures the proportion of recurrent points that form
390 vertical lines in the recurrence plot. Vertical lines indicate periods of time where the
391 system remains in the same state or a similar state for an extended period, reflecting
392 phases of low variability. This provides an estimation of the degree of stability of the
393 interaction between the two musicians.

394 - *Divergence* (DIV), which is calculated as $1/L_{max}$, L_{max} being the longest diagonal
395 line. The higher this metric, the more chaotic is the relation between the two time
396 series. This provides an estimation of the degree of independence of the two musicians

397 As mentioned above, the CRQA was performed for both RMS and spectral centroid. For each
398 moment of RMS manipulation, and for each one of our three metrics (Determinism,
399 Laminarity, Divergence), we thus obtained two values (one for RMS and one for spectral
400 centroid). Those two values were then averaged, resulting in a total of three Interaction
401 variables (Average Determinism, Average Laminarity, Average Divergence) which, taken
402 together, provides a rich picture of the interaction between the musician who had to detect the
403 target sound and the musician with who the target sound was associated (see Figure 2).
404



405

406 Figure 2. Computation of the three interaction variables.

407

408 2.6. Statistical analyses

409 As the duration of the improvisations were different each time, each participant ended up
410 being presented with a different number of targets. To best account for this aspect and the fact
411 that the overall detection performance greatly varied not only across individuals, but also
412 within participants, depending on the musician (i.e., left or right) the target sounds were
413 associated with, we relied on general linear mixed models for our statistical analyses, and

414 used the grouping between Detecting Participant (i.e., the participant that has to press the
415 pedal) and Detected Participant (i.e., the participant with who the target sound is associated)
416 as our random intercept.

417 To test the impact of our RMS manipulations on detection, we did a four-scale hierarchical
418 regression by comparing nested models, starting with a null model (m0), then adding the
419 emergence score as a predictor (m1), and finally adding the RMS category to which the target
420 sound belongs (i.e., baseline, louder or softer) as a second predictor (m2). A last model was
421 fitted which included the interaction term between our two predictors (m3). The models were
422 fitted with the function *glmer* from the R package *lme4* (using a binomial family), and
423 compared using a likelihood ratio test. “Baseline” was used as a base level whenever models
424 included RMS manipulations as a predictor.

425 To test whether Baseline1-target sounds were better detected than Baseline2-target sounds,
426 we did a two-scale hierarchical regression by comparing nested models starting with a model
427 with only the emergence score as a predictor (m1), and then adding “Baseline type”
428 (Baseline1 or Baseline2) as a second predictor. The models were fitted with the function
429 *glmer* from the R package *lme4* (using a binomial family), and compared using a likelihood
430 ratio test. “Baseline2” was used as a base level.

431 To test the impact of our RMS manipulation on the musicians’ interactions, we did a series of
432 two-scale hierarchical regressions by comparing nested models, with each one of our three
433 Interaction variables as our dependent variables. We started with a null model (resp. m4, m6,
434 and m8) and then added the RMS category (i.e., baseline, louder or softer) to which the
435 interaction score under consideration was associated as our predictor (resp. m5, m7 and m9),
436 using “Baseline” as the base level. The models were fitted with the function *lmer* from the R
437 package *lme4*, and compared using a likelihood ratio test.

438 We also performed pairwise post-hoc tests using estimated marginal means, with Kenward-
439 Roger degrees of freedom method and Benjamin-Hochberg method for p-values adjustment,
440 using the *emmeans* package in R.

441

442 3. Results

443 *3.1. RMS manipulations impact musicians' attention*

444 The likelihood ratio test for model comparison between the null model (m0) and the model
445 with perceptual emergence as a predictor (m1) was significant ($X^2=103.93$, $p<0.001$). The
446 likelihood ratio test for model comparison between m1 and the model with both perceptual
447 emergence and RMS manipulations as predictors (m2) was also significant ($X^2=8.343$,
448 $p=0.015$), but adding the interaction term (m3) did not result in a better fit ($X^2=1.964$,
449 $p=0.375$). As expected, our model m2 showed a highly significant effect of perceptual emer-
450 gence on detection (Estimate=2.496, SE=0.256, $z=9.739$ and $p<0.001$): target sounds which
451 emerged more strongly from the musician's acoustic signal were better detected. But more
452 importantly, the model also showed a significant effect of the "louder pattern" condition on
453 detection (Estimate=0.343, SE=0.122, $z=2.813$, $p=0.005$) as well as a significant effect of the
454 "softer pattern" condition on detection (Estimate=0.239, SE=0.119, $z=2.008$, $p=0.045$), as
455 compared to the base level ("Baseline"). In other words, musicians were more likely to detect
456 the target sounds when such sounds happened during a loudness manipulation, and such ma-
457 nipulations had an effect on our participants over and beyond the effects of the perceptual
458 emergence of target sounds. Strikingly, a post-hoc pairwise test also showed that there was no
459 significant difference in the participants' detection rate between the "louder pattern" and
460 "softer pattern" conditions (Estimate=0.104, SE=0.119, $z=0.873$, $p=0.382$), suggesting that
461 the effects of our manipulations on attention was not a matter of intrinsic acoustic features

462 (i.e., the louder, the more salient), but rather of dynamic contrast – in line with most current
463 theories of salience (Kaya et al., 2020).

464 Note however that the effect size of the manipulations remained small (Cohen’s $d=0.097$ for
465 louder patterns and 0.114 for softer patterns). This might be due to the subtleness of our ma-
466 nipulations, which was precisely designed to be as undistruptive as possible.

467

468 *3.2. The impact of RMS manipulations on musicians’ attention is not driven by short-term*
469 *contrasts*

470 While we established that both louder and softer patterns had an effect on our participants’
471 auditory attention, the reason why they did so remains unclear. In line with the many studies
472 highlighting the link between surprise and attention (see Hortsman, 2015, for a review),
473 prevalent accounts of acoustic salience generally highlight the role of contrasts on both short-
474 term and long-term contexts (Huang & Elhilali, 2017; Kaya et al., 2020), but it is not clear
475 that they necessarily have the same weight, particularly in an interactional setting.

476 If the salience of our RMS manipulations was mainly driven by contrastive variations within
477 the local sonic context, this would imply that the target sounds that occur immediately *after* a
478 change back to the baseline (Baseline 1) should also be better detected than the target sounds
479 that occur farther away from the end of the manipulation (Baseline 2), during periods of over-
480 all RMS stability. However, the likelihood ratio test for model comparison between m_1 and
481 our model with Baseline Type as an additional predictor was not significant ($X^2=0.130$,
482 $p=0.718$), suggesting that whether a target sound immediately followed or not a loudness
483 change did not make a difference in the participants’ detection score.

484 It thus seems plausible that the effect of our manipulations on the participants’ attention was
485 driven by the more long-term context, although the precise temporal scale of such long-term
486 context remains unknown. The change back to the baseline could be less salient than the de-

487 parture from the baseline simply in virtue of being the second loudness change within a 10-
488 second interval, thus making it less surprising. Alternatively, it could be less salient because it
489 simply signals a return to a “normal” state of affair. To further substantiate this latter interpre-
490 tation, we ran two additional analyses. First, we computed the amount of time a musician was
491 hearing a given co-improviser under each one of our 3 conditions (baseline, louder pattern,
492 softer pattern). We found that, on average, a participant heard a given co-improviser in the
493 baseline condition 78% of the time, 10.8% of the time in the louder pattern condition, and
494 11.2% of the time in the softer pattern. Second, we computed, for each participant, the mean
495 RMS level they received from their two co-improvisers’ across our 3 conditions, confirming
496 that, on average, co-improvisers were indeed heard significantly louder in the “louder pattern”
497 condition, and significantly softer in the “softer pattern” condition, as compared to the base-
498 line (see Supplementary Material for more details). In other words, participants were exposed
499 on a long-term basis to the typical sound of their co-improvisers as heard in the baseline con-
500 dition, so that in particular, the RMS range perceived during the baseline condition was much
501 more likely to occur than the RMS range perceived during the experimental manipulations.
502 This makes it highly plausible that the improvisers formed internal models of the standard
503 sound of their co-improvisers based on how they typically sounded during the baseline condi-
504 tion, and that the experimental manipulations disrupted the broader expectations associated
505 with these models. The ensuing prediction error is likely to have resulted in a heightened at-
506 tention to what the manipulated musician was playing (Den Ouden et al., 2012), and thus in a
507 better detection score when a target sound occurred during that same period.

508

509 *3.3. RMS manipulations impact musicians’ interactions in a targeted manner*

510 The likelihood ratio tests for model comparisons between the null models (m4, m6 and m8)
511 and the models with RMS manipulation as a predictor (m5, m7 and m9) were significant (X^2

512 = 13.747, $p = 0.001$; $X^2 = 23.326$, $p < 0.001$; and $X^2 = 7.569$, $p = 0.022$, respectively). Model
513 m5 showed a significant effect of the “louder pattern” condition on the average determinism
514 (Estimate=0.03, SE=0.014, $t=2.172$, $p=0.03$) as well as a significant effect of the “softer pat-
515 tern” condition (Estimate=0.05, SE=0.013, $t=3.70$, $p < 0.001$), as compared to the base level
516 (“Baseline”). Model m7 also showed significant effects of both the louder pattern (Esti-
517 mate=0.049, SE=0.014, $t=3.382$, $p < 0.001$) and softer pattern (Estimate=0.067, SE=0.014,
518 $t=4.725$, $p < 0.001$) on average laminarity. Finally, m9 showed significant effects of both the
519 louder pattern (Estimate=-0.03, SE= 0.012, $t=-2.440$, $p= 0.015$) and softer pattern (Estimate=-
520 0.029, SE=0.012, $t=-2.366$, $p= 0.018$) on average divergence.

521 In other words, our results suggest that the detecting musician and the detected musician had
522 more similar musical behaviors (i.e., similar recurring patterns were more likely to be found
523 for the two musicians), had a more stable interaction (i.e., the two musicians were more likely
524 to reach a state that persisted over time with little variation), and were less independent from
525 one another (i.e., the two musicians’ behaviors were more likely to be predictive of one an-
526 other) during louder and softer patterns as compared to the baseline.

527 Once again, the effect sizes were small (Cohen’s $d=0.058$ for louder patterns and 0.143 for
528 softer patterns on average determinism; Cohen’s $d=0.124$ for louder patterns and 0.213 for
529 softer patterns on average laminarity; Cohen’s $d=0.106$ for louder patterns and 0.106 for soft-
530 er patterns on average divergence). The highly complex and ever-changing nature of the mu-
531 sic produced in a CFI context means that results obtained through algorithmic tools primarily
532 designed to extract acoustical features from much simpler musical signals (e.g., pop music)
533 are bound to be very noisy. What is remarkable here is not so much the effect sizes them-
534 selves but rather that a significant effect on all three of our interaction metrics was found for
535 both patterns, despite the noisy nature of the acoustical data, suggesting a clear and consistent
536 effect of salient cues on the musicians’ interactional dynamics.

537

538 4. Discussion

539

540 Our study – which combined real-time covert loudness manipulations and online probe detec-
541 tion – made it possible to systematically investigate for the first time the effects of acoustic
542 salience on auditory attention during live and complex interactions, using musical improvisa-
543 tion as an experimental paradigm. We found, first, that musicians were more likely to pay
544 attention to a given musician when this musician had been made sounding louder or softer;
545 second, that such salience cue appeared to be primarily driven by the long-term acoustic con-
546 text, as no significant salience-like effect was observed when loudness returned back (increas-
547 ing or decreasing) to the baseline value after our manipulation; and third, that salient cues also
548 had an effect on the ongoing musical interaction – with improvisers tending to interact more
549 strongly and in a more stable way with a musician that had been made sounding louder or
550 softer. Taken together, these results shed new light on the role of salience in social-attentional
551 processes and extend previous experimental observations to a full-fledge, complex interac-
552 tional setting.

553 Attention can take many forms and objects, and attentional processes in musical contexts are
554 no exception: musicians' attention can be alternatively mostly focused on themselves or on
555 the other performers, depending on one's own musical part – i.e., how demanding or virtuosic
556 it is – or on coordination requirements – i.e., how the sonic and temporal unfolding of your
557 own part is dependent of that of the other performers (Keller, 2001; Faraco et al., 2024;
558 Abalde et al., 2024); it can also shift from deep absorption to something more akin to mind-
559 wandering (Høffding, 2019). But – as hinted by the observed inter-participant variability in
560 detection score (see Figure S3 in Supplementary Material) – improvisers' attentional profiles
561 can also be widely different from one another, for reasons that might have to do with the in-

562 strument they play and the musical functions that are traditionally associated with such in-
563 strument (e.g., in a jazz context, whether one own's instrument is usually part of the frontline
564 or of the rhythm session, see Monson, 1996), their musical background, and their own repre-
565 sentations of the practice of improvisation itself. Our paradigm opens new avenues to address
566 these questions empirically.

567 In contexts as indeterminate and open-ended as collective improvisations, it can become cru-
568 cial for agents to come up with strategies that can capture their co-agents' attention as a way
569 to signal a given intention (e.g., the intention to end the performance, see Goupil et al., 2021).
570 Previous observational works on musical collective improvisation (Canonne & Garnier, 2012)
571 have suggested that modifying one's own production in such a way as to make it more salient
572 within the group precisely is one of these strategies. Our results provide considerable addi-
573 tional ground to this suggestion by demonstrating that the interactional impact of our loudness
574 manipulations goes over and beyond simple behavioral adaptations such as the Lombard ef-
575 fect (which was indeed also present in our study, see Supplementary Material): on the one
576 hand, improvisers' interactions were also impacted when someone in the group was made to
577 be perceived as playing *softer* (whereas the Lombard effect is only observed when the ambi-
578 ent noise *increase*); and on the other hand, salient events also impacted the very dynamics of
579 the interaction between improvisers, creating points of stability within the musical flux. In
580 other words, we should not think of acoustic salience as only a set property of a given audito-
581 ry stream; it is also something that can be actively manipulated by the agents themselves for
582 communicative purposes. Acoustic salience thus acts as a kind of ostensive gesture which
583 guide not only attention but also the various inferential processes that are bound to emerge in
584 most forms of social communication (Frith, 2008).

585 The salience effect of our loudness manipulations can be further interpreted in the light of
586 studies on attentional selection. Most studies on salience have considered the salience of local

587 events in the auditory signal in a stimulus-driven, purely bottom-up way, and thus have de-
588 veloped models that mostly approach events' salience with respect to the emergence of their
589 instantaneous acoustical statistics in a short-term context (see Kaya & Elhilali, 2014) and
590 more recently to their semantics (Kothinti & Elhilali, 2023). But it might also be relevant to
591 frame acoustic salience within a more general predictive-coding/Bayesian theory (Friston et
592 al., 2012), and approach it not through the event/context relationship as this is typically done
593 but rather through the short-term/long-term statistical relationship. In that respect, it is inter-
594 esting to mention three factors that have been found to lead to an enhancement of attentional
595 selection beyond perceptual emergence (see Theeuwes et al., 2019, for a review): (i) history-
596 driven selection – which relates to statistical learning, with stimuli that belong to high-
597 probability categories leading to less attentional capture compared to those associated to low-
598 probability categories; (ii) value-driven selection – where a reward associated to a stimulus
599 can enhance its attentional capture; and (iii) goal-driven selection – where a change in stimu-
600 lus feature that is useful for a particular task can enhance attentional capture. Beyond the role
601 of the statistical structure of the stimuli in explaining the effect of loudness manipulations on
602 attentional capture observed here, one might wonder whether a change in a specific musi-
603 cian's overall loudness carries an intrinsic musical value in this context, e.g., by providing
604 opportunities for the group to renegotiate the interactional dynamics. Yet, since disentangling
605 those various mechanisms was not the initial purpose of our research, new studies with specif-
606 ic experimental designs would be needed to further explore the different factors shaping
607 attentional selection in collective musical improvisation, beyond the too simple bottom-
608 up/top-down dichotomy (Awh et al., 2012). It would also be interesting to study whether other
609 acoustic dimensions would play a comparable role in such a musical context, e.g., whether
610 experimentally-induced variations in timbre or average pitch would lead to similar effects on
611 the other musician's attention and interaction. We could hypothesize that as soon as these

612 changes are meaningful musically-speaking, that they would induce comparable effects. This
613 remains to be tested empirically.

614 The preeminence of the long-term context over the short-term context observed in our study
615 might also precisely be due to its interactional nature. The sonic information perceived by the
616 participants was not treated as mere acoustical cues but also as full-fledged behavioral cues,
617 ultimately providing them with a sense of their (musical) personalities. This certainly led the
618 improvisers to pay particular attention to the acoustic cues that were highly informative about
619 their co-performers. In that perspective, a sonic event that strongly stands out from the local
620 acoustic context might be less relevant than a less spectacular sonic event that does not match
621 with the expected sonic and musical personality of a given co-improviser – as typically in-
622 ferred over a longer time scale. The complex relations between short-term and long-term con-
623 texts, as well as between acoustic cues and interactional cues, should provide an exciting topic
624 for further experimental explorations.

625 It is also important to acknowledge that the effects of our loudness manipulations (on both
626 attention and interaction) were quite small. Our main motivation for selecting such a subtle
627 loudness variation as our experimental manipulation was to make the manipulation as unno-
628 ticeable as possible, so as to not artificially disrupt the interactions between the musicians. As
629 such, it remains an open question whether our small effect size is due to a slightly too subtle
630 manipulation, or whether more overt acoustic manipulations would have yield similarly small
631 effects, because, in this kind of context, musicians' attention is merely driven by top-down
632 factors (associated to specific interactional goals or idiosyncratic musical intentions) rather
633 than salient cues. Further studies could shed light on this issue by manipulating top-down at-
634 tention at the same time (e.g., by providing musicians with dynamical attentional scripts that
635 they must follow as they improvise) to assess the relative importance of both factors. On the
636 other hand, our diotic-based target detection paradigm could be directly transferred to further

637 address such questions in lab-based experiments with isolated participants, extending previous
638 psychophysical paradigms to assess auditory attention (Huang & Elhilali, 2017). Participants
639 would not be interacting with the music played in their headphones but rather performing an
640 additional unrelated cognitively demanding task, while we would measure how their perfor-
641 mance for detecting a target appearing in one ear or the other of a diotic auditory scene varies
642 with both bottom-up and top-down factors. Results from such in-lab experiments could pro-
643 vide insightful benchmarks for interpreting the size of the effects observed in online experi-
644 ments, where the active task of the listener is arguably the most demanding.

645 Another interesting issue is the extent to which our results could generalize to other forms of
646 sonic interactions between humans. In particular, it is an open question whether loudness ma-
647 nipulations could play a similar role in impacting interactions in multi-speaker conversations
648 in natural languages. CFI is mainly a hedonistic practice, devoid of any practical concerns,
649 making it possible for the improvisers to simply follow whatever catch their attention, and to
650 let the music emerge from the ongoing interactions. Moreover, the lack of a clear musical
651 semantics also means that salience is bound to play an even greater role in making the com-
652 munication between the musicians possible, as it is often the case in semantically undeter-
653 mined contexts (Kecskes, 2013). In contrast, it might well be the case that the role of auditory
654 salience is lessened in the regulation of verbal interactions in which participants have more prac-
655 tical concerns (which might strongly constrain how the interaction will unfold) and can all
656 rely on clear semantic resources. Here, we believe that our experimental procedure could be,
657 in a principled-way, easily and fruitfully transposed to the study of complex multi-speaker
658 conversations, allowing us to contrast effectively the role of auditory salience in both verbal
659 interactions and musical interactions. Further studies addressing the respective impact of
660 short-term feature emergence vs. changes in long-term stimulus statistics will thus be key for

661 a better understanding of how meaningful dynamics in speech prosody are integrated and con-
662 sequently shape conversational interactions.

663 Finally, although our study suggests a link between attention and interaction, its exact nature
664 remains unclear. Two concurrent explanations could account for this connection. On the one
665 hand, attention might be seen as a condition for interaction: paying more attention to a fellow
666 agent might make it more likely to engage in an interaction with them. But on the other hand,
667 it is also possible that attention (to someone) emerges from the interaction itself – considered
668 as an autonomous system made of mutual couplings (De Jaegher et al., 2010): what is felt or
669 noticed first is the interaction, which then increases the attention towards the other parts of the
670 interactive network. This hints at a variety of communicational strategies that are mediated
671 either through interaction or through attention, and which could be more sharply teased apart
672 in further studies. Moreover, while our study shed light on the dyadic musical interactions
673 between the manipulated musician and the musician that listened to the manipulated musician,
674 further experimental protocols have to be set up in order to fully understand how these in-
675 duced changes in salience on a single musician shape the interactional dynamics at the broad-
676 er group level. In particular, one could assume that if musicians rely on a constant amount of
677 interactional resources (Keller, 2001), the strengthening of the interaction with the salient
678 musician would necessarily come at the cost of a weaker interaction with the non-salient mu-
679 sicians. Addressing such questions of course requires further conceptual elaboration as well as
680 the development of novel experimental protocols that can also provide information about the
681 overall amount of attention deployed by each musician, but doing so might produce key data
682 to further develop theoretical models of musical interaction.

683 Systematically studying attention in the wild – and moving from individual lab booths to
684 more complex interactional settings – is arguably one of the most pressing challenges for so-
685 cial cognition research. Music – and in particular CFI – might provide a relevant experimental

686 paradigm for such research program, allowing us to balance experimental control with ecolog-
687 ical validity in a satisfying way (D'Ausilio et al., 2015). Our study illustrates once again the
688 value of musical settings for investigating social cognition, by providing a first attempt at
689 tracking auditory attention in real-time during complex interactions. While music is arguably
690 a much more harmonious environment than the standard noisy cocktail party situation, there
691 is no doubt that it has still a lot to reveal on the inner workings of auditory attentional pro-
692 cesses in highly social contexts.

693

694 References

695

696 Abalde, S. F., Rigby, A., Keller, P. E., & Novembre, G. (2024). A framework for joint music
697 making: behavioral findings, neural processes, and computational models. *Neuroscience &*
698 *Biobehavioral Reviews*, 105816.

699

700 Arnal, L. H., Kleinschmidt, A., Spinelli, L., Giraud, A. L., & Mégevand, P. (2019). The rough
701 sound of salience enhances aversion through neural synchronisation. *Nature communica-*
702 *tions*, 10(1), 3671.

703

704 Awh, E., Belopolsky, A. V., & Theeuwes, J. (2012). Top-down versus bottom-up attentional
705 control: A failed theoretical dichotomy. *Trends in cognitive sciences*, 16(8), 437-443.

706

707 Bidelman, G. M., & Yoo, J. (2020). Musicians show improved speech segregation in competi-
708 tive, multi-talker cocktail party scenarios. *Frontiers in Psychology*, 11, 1927.

709

710 Blacking, J. (1973). *How musical is man?*. University of Washington Press.

711

712 Canonne, C., Garnier, N. (2012). Cognition and Segmentation in Collective Free Improvisa-
713 tion: An Exploratory Study. *Proceedings of the 12th International Conference on Music Per-*
714 *ception and Cognition 8th Triennial Conference of the European Society for the Cognitive*
715 *Sciences of Music*, Thessaloniki, Greece.

716

717 Cherry, E. C. (1953). Some experiments on the recognition of speech, with one and with two
718 ears. *The Journal of the acoustical society of America*, 25(5), 975-979.
719

720 Clarke, E. (2005). *Ways of listening: An ecological approach to the perception of musical*
721 *meaning*. Oxford University Press.
722

723 Dalton, P., & Lavie, N. (2004). Auditory attentional capture: Effects of singleton distractor
724 sounds. *Journal of Experimental Psychology: Human Perception and Performance*, 30(1),
725 180. <https://doi.org/10.1037/0096-1523.30.1.180>
726

727 D'Ausilio, A., Novembre, G., Fadiga, L., & Keller, P. E. (2015). What can music tell us about
728 social interaction?. *Trends in cognitive sciences*, 19(3), 111-114.
729

730 Dawson, J., & Foulsham, T. (2022). Your turn to speak? Audiovisual social attention in the
731 lab and in the wild. *Visual Cognition*, 30(1-2), 116-134.
732

733 De Jaegher, H., Di Paolo, E., & Gallagher, S. (2010). Can social interaction constitute social
734 cognition?. *Trends in cognitive sciences*, 14(10), 441-447.
735

736 Den Ouden, H. E., Kok, P., & De Lange, F. P. (2012). How prediction errors shape percep-
737 tion, attention, and motivation. *Frontiers in psychology*, 3, 548.
738

739 Faraco, A., Schwarz, A., Vincent, C., Susini, P., Ponsot, E., & Canonne, C. (2024). Listening
740 Behaviors and Musical Coordination in Collective Free Improvisation. *Music & Science*, 7,
741 20592043241257023.

742

743 Friston, K., Adams, R. A., Perrinet, L., & Breakspear, M. (2012). Perceptions as hypotheses:
744 saccades as experiments. *Frontiers in psychology*, 3, 151

745

746 Frith, C. D. (2008). Social cognition. *Philosophical Transactions of the Royal Society B: Bio-*
747 *logical Sciences*, 363(1499), 2033-2039.

748

749 Goupil, L., Wolf, T., Saint- Germier, P., Aucouturier, J. J., & Canonne, C. (2021a). Emergent
750 shared intentions support coordination during collective musical improvisations. *Cognitive*
751 *Science*, 45(1), e12932.

752

753 Høffding, S. (2019). *A phenomenology of musical absorption*. Springer.

754

755 Horstmann, G. (2015). The surprise–attention link: A review. *Annals of the New York Acade-*
756 *my of Sciences*, 1339(1), 106-115.

757

758 Huang, N., & Elhilali, M. (2017). Auditory salience using natural soundscapes. *The Journal*
759 *of the Acoustical Society of America*, 141(3), 2163-2176.

760

761 Humphreys, G. W., Yoon, E. Y., Kumar, S., Lestou, V., Kitadono, K., Roberts, K. L., &
762 Riddoch, M. J. (2010). The interaction of attention and action: From seeing action to acting on
763 perception. *British Journal of Psychology*, 101(2), 185-206.

764

765 Kaya, E. M., & Elhilali, M. (2014). Investigating bottom-up auditory attention. *Frontiers in*
766 *human neuroscience*, 8, 327.

767

768 Kaya, E. M., & Elhilali, M. (2017). Modelling auditory attention. *Philosophical Transactions*
769 *of the Royal Society B: Biological Sciences*, 372(1714), 20160101.

770

771 Kaya, E. M., Huang, N., & Elhilali, M. (2020). Pitch, timbre and intensity interdependently
772 modulate neural responses to salient sounds. *Neuroscience*, 440, 1-14.

773

774 Kecskes, I. (2013). Cross-cultural and intercultural pragmatics. Oxford University Press.

775

776 Keller, P. E. (2001). Attentional resource allocation in musical ensemble perfor-
777 mance. *Psychology of Music*, 29(1), 20-38.

778

779 King, A., Varnet, L., & Lorenzi, C. (2019). Accounting for masking of frequency modulation
780 by amplitude modulation with the modulation filter-bank concept. *The Journal of the Acousti-*
781 *cal Society of America*, 145(4), 2277-2293.

782

783 Koch, I., Lawo, V., Fels, J., & Vorländer, M. (2011). Switching in the cocktail party: explor-
784 ing intentional control of auditory selective attention. *Journal of Experimental Psychology:*
785 *Human Perception and Performance*, 37(4), 1140.

786

787 Kothinti, S. R., & Elhilali, M. (2023). Are acoustics enough? Semantic effects on auditory
788 salience in natural scenes. *Frontiers in Psychology*, 14, 1276237.

789

790 McDermott, J. H. (2009). The cocktail party problem. *Current Biology*, 19(22), R1024-
791 R1027.

792

793 McFee, B., Raffel, D., Dawen L., Ellis, D., McVicar, M., Battenberg, E., Nieto, O. (2015).

794 "librosa: Audio and music signal analysis in python." *Proceedings of the 14th python in sci-*

795 *ence conference*, pp. 18-25.

796

797 Miles, K., Weisser, A., Kallen, R. W., Varlet, M., Richardson, M. J., & Buchholz, J. M.

798 (2023). Behavioral dynamics of conversation, (mis) communication and coordination in noisy

799 environments. *Scientific Reports*, 13(1), 20271.

800

801 Monson, I. (1996). *Saying something: Jazz improvisation and interaction*. University of Chi-

802 cago Press.

803

804 Noy, L., Dekel, E., & Alon, U. (2011). The mirror game as a paradigm for studying the dy-

805 namics of two people improvising motion together. *Proceedings of the National Academy of*

806 *Sciences*, 108(52), 20947-20952.

807

808 Rawald, T., Sips, M., Marwan, N. (2017). PyRQA. Conducting Recurrence Quantification

809 Analysis on Very Long Time Series Efficiently. *Computers and Geosciences*, 104, pp. 101-

810 108. doi: <https://doi.org/10.1016/j.cageo.2016.11.016>.

811

812 Roueff, O. (2006). L'invention d'une «scène» musicale, ou le travail du réseau: La program-

813 mation d'un club de musiques improvisées entre radicalisation et consécration (1991-

814 2001). *Sociologie de l'Art*, (1), 43-76.

815

816 Ryan, F., Jiang, H., Shukla, A., Rehg, J. M., & Ithapu, V. K. (2023). Egocentric auditory at-
817 tention localization in conversations. In *Proceedings of the IEEE/CVF Conference on Com-*
818 *puter Vision and Pattern Recognition* (pp. 14663-14674).

819

820 Saint-Germier, P., Canonne, C. (2020). Coordinating free improvisation: An integrative
821 framework for the study of collective improvisation. *Musicae Scientiae*, 26(3), 455–475,
822 <https://doi.org/10.1177/1029864920976182>.

823

824 Schultz, B. G., Brown, R. M., & Kotz, S. A. (2021). Dynamic acoustic salience evokes motor
825 responses. *Cortex*, 134, 320-332.

826

827 Seddon, F. A. (2005). Modes of communication during jazz improvisation. *British Journal of*
828 *Music Education*, 22(1), 47-61.

829

830 Steffens, J., Müller, F., Schulz, M., & Gibson, S. (2020). The effect of inattention and cogni-
831 tive load on unpleasantness judgments of environmental sounds. *Applied Acoustics*, 164,
832 107278.

833

834 Straetmans, L., Holtze, B., Debener, S., Jaeger, M., & Mirkovic, B. (2022). Neural tracking to
835 go: auditory attention decoding and saliency detection with mobile EEG. *Journal of neural*
836 *engineering*, 18(6), 066054.

837

838 Theeuwes, J. (2019). Goal-driven, stimulus-driven, and history-driven selection. *Current*
839 *opinion in psychology*, 29, 97-101.

840

841 Van der Heiden, R. M., Kenemans, J. L., Donker, S. F., & Janssen, C. P. (2022). The effect of
842 cognitive load on auditory susceptibility during automated driving. *Human factors*, *64*(7),
843 1195-1209.

844

845 Wallot, S. & Leonardi, G. (2018). Analyzing Multivariate Dynamics Using Cross-Recurrence
846 Quantification Analysis (CRQA), Diagonal-Cross-Recurrence Profiles (DCRP) and Multidi-
847 mensional Recurrence Quantification Analysis (MdRQA) – A Tutorial in R. *Frontiers in Psy-*
848 *chology*, *9*(2232). DOI: doi: 10.3389/fpsyg.2018.02232.

849

850 Wolf, W., Launay, J., & Dunbar, R. I. (2016). Joint attention, shared goals, and social bond-
851 ing. *British Journal of Psychology*, *107*(2), 322-337.

852

853 Wu, W. (2014). *Attention*. Routledge.

854

855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904

Supplementary information

1. Video examples

The two video examples provided show musicians from two different trios improvising during our experimental protocol. The videos are displayed in such a way as to illustrate our experimental procedure. Each line is meant to represent the subjective perspective of the musician displayed in the center. The musician shown on the left was the musician heard on the left channel of their headphones, and the musician shown on the right was the musician heard on the right channel of their headphones. RMS manipulations applied to these musicians (i.e., louder and softer patterns) are shown through changes in brightness (the video of a musician gets darker when this musician is heard through a softer pattern, and the video of a musician gets brighter when this musician is heard through a louder pattern). Finally, sudden flashes in the video happen at times a target sound associated with that musician was presented. Those videos were filmed and edited by Elsa Laurent.

2. Additional information on the participants

We also tested our participants for hearing losses. Audiometric thresholds were measured between 0.25 and 8 kHz in all participants, using an Echoscan device. Importantly, we ensured that participants had clinically-normal thresholds (<15 dB HL) in the frequency range (from 0.25 to 2 kHz) of the probe used for the detection task. Most of them also had a normal hearing at high frequencies, except for two participants who had a moderate hearing loss: one at 4 kHz and one at 8 kHz. For both of them, the loss was predominant in the left ear (~ 40 dB HL) in regard of the right ear (25 dB HL). Both participants were aware of this. As their results were similar to the others, and because the target sounds were lower in frequency and never presented at threshold, we decided to keep them in the analysis.

3. Assessing the naturalness of the RMS manipulations.

This follow-up experiment was designed to estimate, in a post-hoc fashion, the degree of naturality of the RMS manipulations introduced during our recording session. In particular, we were wondering whether external listeners would be more likely to perceive the overall musical output as less natural when one of the musicians was subjected to a RMS manipulation, as compared to the same output but without the RMS manipulation.

3.1. *Participants*

27 participants took part in study 2 (mean age = 26 years, SD = 5, 14 women and 13 men). They all had practiced music for at least 5 years (mean number of years = 12.5 years, SD = 5.5). All participants gave their informed written consent and were paid at a standard rate.

3.2. *Stimuli*

18 audio excerpts were chosen randomly at moments where one musician's RMS level was modified (either with a high-level pattern or a low-level pattern) for another musician. Those excerpts were directly taken from the audio tracks heard by the improvisers in their headphones. Each excerpt was 15 seconds long, corresponding to the 10-second period of the RMS manipulations, with the addition of two 2.5-second buffer zones at the start and at the end of the excerpt. For each excerpt, we then created two alternative versions: one containing

905 the opposite pattern from the one originally applied (i.e., a louder pattern instead of a softer
906 pattern); and one with no RMS manipulation whatsoever. This resulted in a total of 54 ex-
907 cerpts. Finally, 50-ms fade-ins/fade-outs were imposed at the beginning and the end of each
908 excerpt to avoid clipping.

909 3.3. Procedure

911 The following cover story was provided to the participants: “We are testing an artificial intel-
912 ligence designed to help sound engineers in mixing various instruments playing live together.
913 For each rough recording given to the AI, the AI offers various mixes. In this study, you will
914 have to compare, for several different excerpts, two of the mixes that were created by the AI.
915 We ask you to tell us which version you think sounds the most natural”. Our design thus fol-
916 lowed a two alternative forced choice (2-AFC) procedure. For each trial, two versions of the
917 same musical excerpt (e.g., the same music but once with a high-level pattern applied to one
918 of the musicians of the trio, and the other time with a low-level pattern applied to the same
919 musician) were successively presented to the participants (with a 2-second gap in between),
920 and participants had to select the version that felt the most “natural” for them.

921 The participants heard the stimuli through headphones (Beyerdynamics DT 770 pro, 80 ohms)
922 and answered via a custom-made interface designed with Max/MSP.

923 Over the course of the experiment, participants had to compare each one of the three versions
924 of each musical excerpt to the other two (i.e., louder pattern vs softer pattern; louder pattern
925 vs no-manipulation; softer pattern vs no-manipulation). This resulted in a total of 54 trials,
926 which were divided in three successive blocks of 18 trials. The order of the stimuli was pseu-
927 do-randomized so that, over all trials and participants, each stimulus was presented an equal
928 number of times as the first excerpt and as the second excerpt. Participants were invited to
929 take a break in between blocks.

930 3.4. Statistics

932 A two-scale hierarchical regression was conducted by comparing nested models, starting with
933 a null model and adding our experimental condition (i.e., whether the stimulus belong to the
934 no-pattern, softer pattern or louder pattern category). The effect of our experimental condition
935 on whether a stimulus was chosen to sound more natural than the comparison stimulus was
936 tested through a mixed binomial regression model with Condition as predictor and partici-
937 pants’ choice (natural: yes/no) to be predicted. “No-pattern” was used as our model base lev-
938 el. In the random structure we included the random slope for each listener. The models were
939 fitted with the function *glmer* from the R package *lme4* and compared using a likelihood ratio
940 test.

941 3.5. Results

943 Figure S1 shows the results of this follow-up experiment. The likelihood ratio test for model
944 comparison was not significant ($X^2=3.730$, $p=0.155$), suggesting that whether the RMS of a
945 musician had been modified or not did not make a clear difference in third-party listeners’
946 naturalness judgments. In particular, it does not seem that listeners were more likely to find
947 the excerpt’s sound more natural when we presented them with the musicians’ signals as ac-
948 tually played by them – without any RMS manipulations. This shows that our RMS manipula-
949 tions were indeed subtle enough to not be clearly noticeable by external listeners, suggesting
950 that they could have passed as real musical intentions (i.e., playing a bit louder or a bit softer)
951 during the improvisers’ interactions.

953

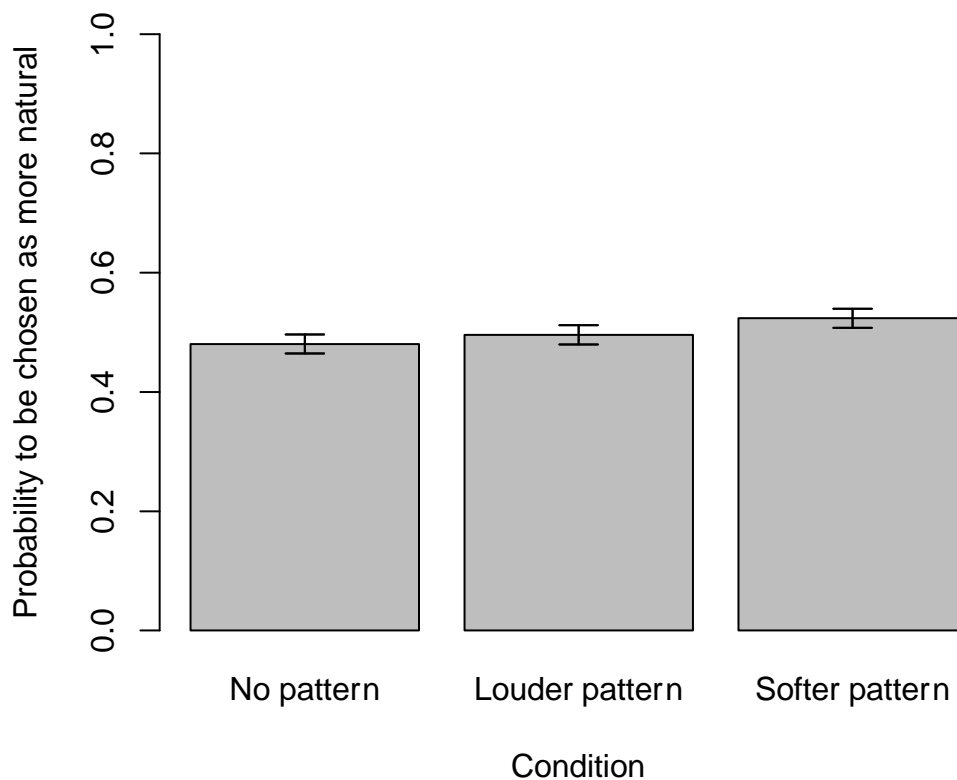


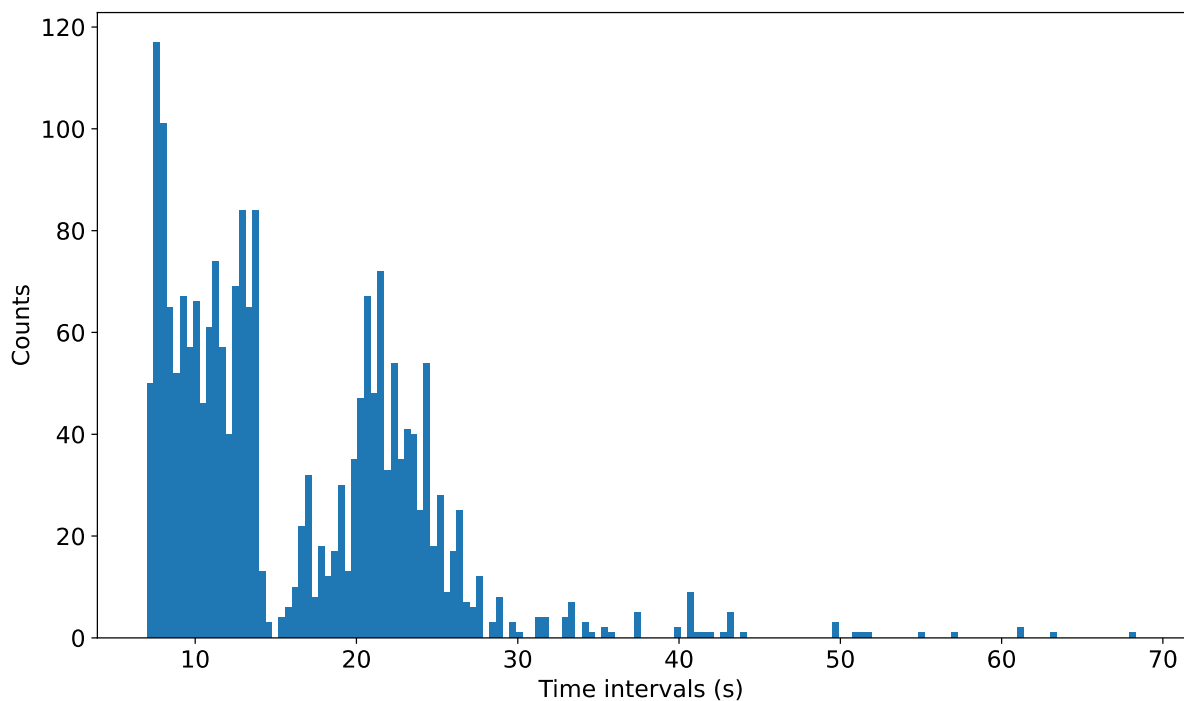
Figure S1. Bars represent the probability of a snippet from a certain Condition (No-pattern, Louder pattern, or Softer pattern) to be chosen as more natural than the comparison snippet. There were 36 datapoints for each participant and condition ($36 * 3 = 108$ datapoints per participant), which were averaged per participant ($N = 27$) and Condition before they were used in the plot. Error bars show standard error.

4. Analyses suggest that participants did not merely learn to associate RMS manipulations and target sound occurrences.

In order to decrease potential learning effects, target sounds and RMS manipulations were not systematically associated: on the one hand, roughly 31% of the target sounds occurred in the “baseline” condition, during which no RMS manipulation was applied; on the other hand, roughly 6.5% of the RMS manipulations were not associated to any target sound. Similarly, the distribution of the time intervals between two successive target sounds heard by a same participant was quite wide (see Figure S2 below), making it unlikely that the participants would press their pedals based on their learning a typical interval between two consecutive target sounds.

Post-hoc analyses of our results were performed to ensure that the effects of RMS manipulations on target sound detection was not the mere result of a learning effect. As our participants were asked to perform several improvisations during the experiment, we were able to compare detection scores for their first and last improvisations. Our reasoning was that if that there was any RMS-driven learning association issue, the effect of loudness patterns on target detection would be larger in the last improvisation (i.e., after learning) compared to the first improvisation. To do so, a hierarchical regression was conducted by comparing nested

980 models starting with a model with only the emergence score as a predictor (m1), second,
 981 adding “Order” (First improvisation or Last improvisation) as a predictor (m2), and third, the
 982 interaction term between Order and RMS manipulations (baseline, louder pattern, softer
 983 pattern) as a predictor (m3). The likelihood ratio test for model comparison between m1 and
 984 m2 was significant ($X^2=13.763$, $p<0.001$). Our model m2 showed a highly significant effect
 985 of Order on detection (Estimate=0.059, SE=0.160, $z=3.692$ and $p<0.001$): target sounds were
 986 indeed better detected in the last improvisation than in the first performance. However, the
 987 likelihood ratio test for model comparison between m2 and m3 was not significant
 988 ($X^2=8.553$, $p=0.073$), showing that participants were not significantly better to detect the
 989 target sounds associated with higher or louder patterns in the last improvisation than in the
 990 first improvisation. Hence, these analyses support the view that the effect of patterns on
 991 detection could not be simply explained by a RMS-driven learning association.
 992



993
 994 Figure S2. Distribution of the time intervals between two consecutive target sounds heard by a same participant

995
 996 5. Computing the perceptual emergence of target sounds
 997

998 Although target sounds were embedded in real time in the musical mix at a constant signal-to-
 999 noise ratio, this did not account for masking effects and therefore did not guarantee equal
 1000 “audibility” of all target sounds. To account for such effects, we *a posteriori* estimated the
 1001 “perceptual emergence” associated to each embedded target using the auditory modulation
 1002 filterbank model (MFB), a computational model² of the human auditory system that mimics
 1003 the main peripheral and central processing stages and can therefore account for basic percep-
 1004 tual masking effects and target-in-noise psychometric characteristics.

1005 We used the MFB implementation of King et al. (2019) using the code provided by the *Audi-*
 1006 *tory Modelling Toolbox* (Majdak et al., 2022; AMToolbox v. 1.5.0). A detailed description of

² The model is built from a cumulative body of works since the 90’s (Dau et al., 1997) and is now used to account for a large range of human perceptual capacities ranging from complex targets detection in noise to speech recognition in noise (e.g., Relación-Iborra et al., 2016).

1007 the MFB model and its different processing stages can be found in Ponsot et al. (2021). The
1008 main parameters of the MFB model were chosen to best reflect auditory processing of nor-
1009 mal-hearing listeners (frequency range of analysis was limited to that of the target, i.e., 100-
1010 2000Hz; N=10 cochlear filters with ERB=1; Q-value of modulation filters = 1; modulation
1011 phase cut-off at 5 Hz). Importantly, since we the overall sound pressure level of the mix re-
1012 ceived by the musicians in their headphones remained unknown (each musician could indeed
1013 freely adjust the overall level of the mix that they received), we used a linear version of the
1014 model with respect to sound level and entered a value of 1 for the compression parameter (i.e.,
1015 no compression), meaning that the model did not attempt to account for level-dependence of
1016 masking and other more complex non-linear effects with respect to sound level.

1017 We used the 3-D internal representation {frequency channels x modulation channels x time}
1018 of sounds provided at the output of the MFB model. We began by computing the internal
1019 model representation of the target sound alone (i.e., without any additional musical signal) to
1020 create a ‘clean template’. Next, we extracted from each individual track 1-second temporal
1021 segments corresponding to the temporal windows of the overall mix at which target sounds
1022 were inserted, and computed their model representations, which yielded ‘noisy templates’.
1023 The similarity between the clean and each noisy template (matrix were transformed into 1-D
1024 vectors) was then evaluated using Pearson correlation. This provided, for each extract, a sca-
1025 lar value of “perceptual emergence” of the target sound, which is assumed to reflect the
1026 amount of perceptual evidence that the target sound was present in the acoustical signal from
1027 the viewpoint of a human auditory observer. By controlling for the degree of “perceptual
1028 emergence” of each target sound on detection performance in the regression models, we
1029 aimed to more precisely assess the degree of attention to the track with which it was associ-
1030 ated, so as to optimally probe the attentional dynamics of the musicians during the perform-
1031 ance.

1032 6. Identifying the hyperparameters for the cross-recurrence quantification analysis

1033 CRQA has been described as “a class of multivariate and generalized correlational analyses,
1034 that are suited for joint action data because they make very few assumptions and are
1035 particularly robust in case of non-linearities, non-stationary dynamics and time-series with
1036 extreme outliers” (Wallot & Leonardi, 2018, p. 2). As such, CRQA has been widely used in
1037 joint action studies, such as ones involving synchronized arousal between performers in
1038 collective rituals (Konvalinka et al., 2011), conversations between medical doctors and
1039 patients (Angus et al, 2012), communicative interaction between children (Lira-Palma et al,
1040 2018), and interaction between parents and children (Dale & Spivey, 2006; Reddy et al.,
1041 2013).

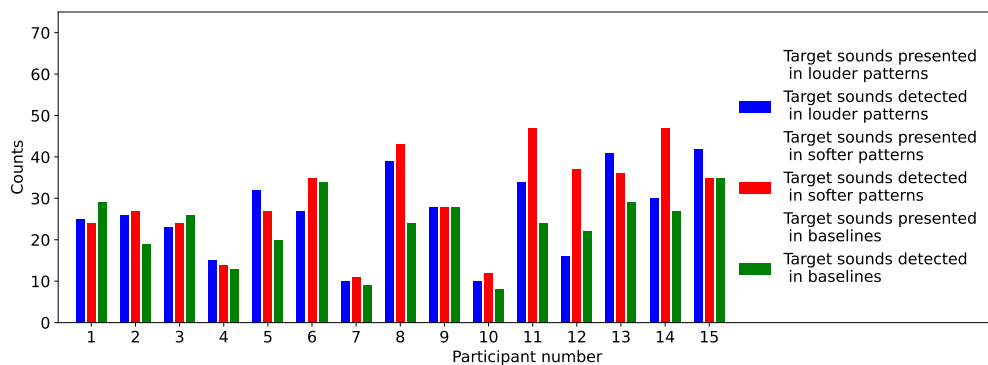
1042 CRQA involves reconstructing two or multiple time series in a phase space, identifying
1043 recurrent states within the system. In other words, the analysis detects moments when the time
1044 series return to similar states, focusing on patterns of recurrence that can reveal the
1045 deterministic or stable nature of the dynamics between the two time-series. In order to
1046 perform the analysis, we followed the procedures described in Wallot & Leonardi (2018) and
1047 Wallot & Monster (2018). In order to do so, we needed to estimate the hyperparameters for
1048 the CRQA, namely the *embedding dimension*, the *time delay* and the *radius*. The time delay
1049 defines the number of time steps between consecutive points in the phase space
1050 reconstruction, while the embedding dimension refers to how many consecutive points of the
1051 time series are considered together when reconstructing the phase space. As for the radius, it
1052 sets the threshold for considering two points in the phase space as “close” or recurrent,
1053 determining the “sensitivity” of the recurrence. To achieve this, we first normalized our time
1054 series data (RMS and spectral centroid) by rescaling them. This step was important for two
1055 reasons: (1) as we used the average of the found optimal radii as a fixed parameter (see
1056

1057 below), both time series needed be on the same scale for the radius to be meaningful; and (2)
 1058 since our interaction measurements are based on the average of the CRQA metrics, both time
 1059 series needed to be standardized to the same scale for accurate comparison. Second, we
 1060 looped over all pairs of segments (comprised of the detecting participant and the detected
 1061 participant) for both RMS and spectral centroid and found the optimal time delay by using the
 1062 average mutual information function, and the embedding dimension by using the false
 1063 nearest-neighbors algorithm (both functions are taken from Wallot & Monster (2018) and
 1064 were adapted from MATLAB to Python). Third, we used the average time delay and
 1065 maximum embedding dimension to calculate the optimal radius based on an iterative process
 1066 that finds the radius in which the percentage of recurrence (%REC) is between 1-10% (Wallot
 1067 & Leonardi, 2018). Fourth, we averaged all hyperparameters thus identified.

1068
 1069 7. Participants' detection scores and response times

1070 Figure S3 shows, for each participant and each experimental condition, the number of target
 1071 sounds they had to detect as well of the number of target sounds they actually detected. Figure
 1072 S4 shows the distribution of the participants' response time for the detection of target sounds,
 1073 revealing a fair amount of inter-participant variability, which is probably due to the large vari-
 1074 ety of instruments involved (e.g., it might have been easier to press a pedal for the saxophone
 1075 players than for the drummers, who already have to use their feet to play their instrument).
 1076 That being said, most pedal presses occurred quite shortly after the target sound occurrence
 1077 ($M = 2.224$; $SD = 0.541$), suggesting that, when a target sound was detected, participants were
 1078 rather fast at pressing their pedal, despite the very demanding setting they were placed in.

1079



1080
 1081 Figure S3. Target sounds presented and detected for each participant in each condition
 1082

1080
 1081
 1082

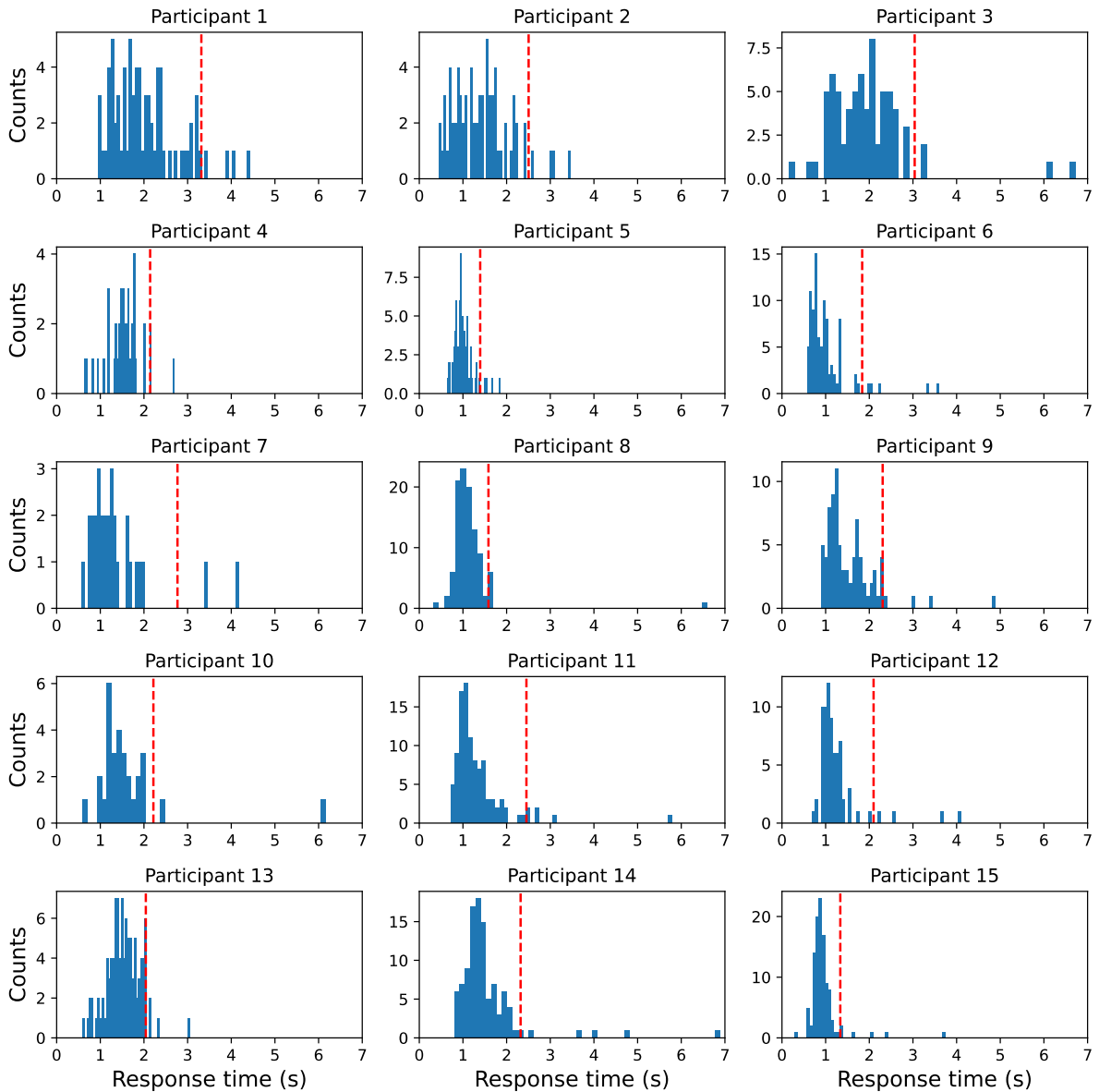


Figure S4. Occurrences of pedal presses as a function of response time for individual participants. The vertical dashed red line shows the cutoff for the 95% percentiles of these distributions.

1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101

8. Replicating our analyses using a much shorter time window for detection yields similar results

To ensure that the pattern of results reported in our main analysis did not depend on the time window we chose to compute the participants' detection score (a target sound was considered to be detected by a participant if the participant pressed their pedal in the 7 seconds following the target sound), we ran the same analyses using a much shorter time window, different for each participant, corresponding to the 95% percentile shown in Figure S4 (which was below 3 seconds for every participant). Strikingly, changing the size of the time window in such a way did not modify our pattern of results. The model still revealed a significant effect of the "louder pattern" condition (Estimate = 0.343, SE = 0.119, $z = 2.872$, $p = 0.004$) and a significant effect of the "softer pattern" condition (Estimate=0.258, SE=0.116, $z=2.216$, $p=0.027$).

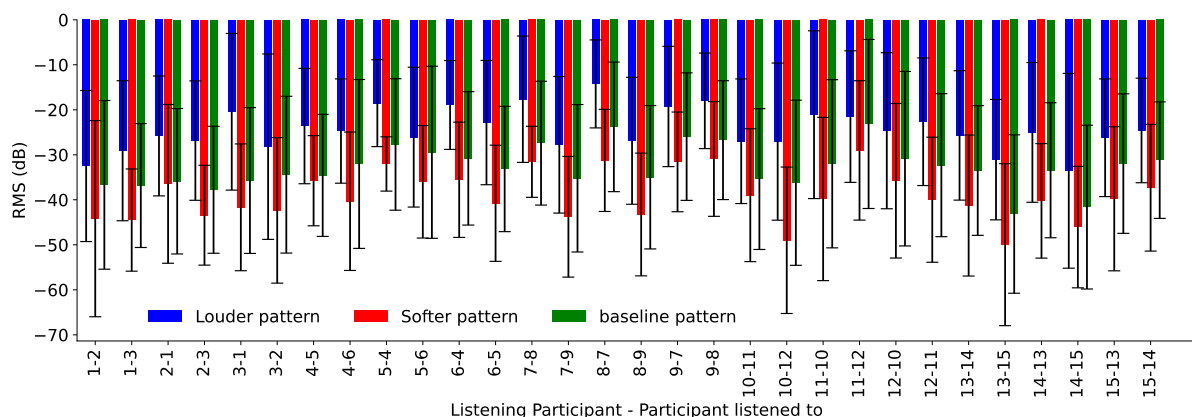
1102 9. Computing the RMS level received by each participant from their two co-improvisers
1103 across the three experimental conditions

1104

1105 For each pair of detecting participant/detected participant, we segmented the audio track of
1106 the detected participant in 10-s windows. Each window was ascribed to one of our 3 condi-
1107 tions (baseline, louder pattern, and softer pattern) depending on the manipulation (or absence
1108 of manipulation) applied to the detected participant during that window. The RMS level was
1109 then averaged over each window. On average, the mean RMS received by a given performer
1110 from a co-performer in the baseline condition was -31.493 dB; the mean RMS received by a
1111 given performer from a co-performer in the louder pattern condition was -23.097 dB; and the
1112 mean RMS received by a given performer from a co-performer in the softer pattern was -
1113 37.502 dB (see Figure S5).

1114 We then conducted a statistical analysis. As the RMS data did not follow a normal distribu-
1115 tion, we did a ranked transformation of our RMS data using the R library bestNormalize. This
1116 allowed us to run a two-scale hierarchical regression by comparing nested models, starting
1117 with a null model and then adding our experimental condition (i.e., whether the window be-
1118 long to the baseline, softer pattern or louder pattern category) as a predictor. The likelihood
1119 ratio test for model comparison was significant ($X^2=537.22$, $p<0.001$). Our model with the
1120 experimental condition as a predictor showed that the RMS level received from musicians
1121 exposed to a louder pattern was significantly higher as compared to the baseline (Estimate =
1122 0.688, $SE=0.034$, $t= 20.087$, $p<0.001$), and that the RMS level received from musicians ex-
1123 posed to a softer pattern was significantly lower as compared to the baseline (Estimate = -
1124 0.321, $SE=0.033$, $t=-9.713$, $p<0.001$).

1125



1126

1127 Figure S5. Mean RMS level received by each participant from their two co-improvisers across the three experi-
1128 mental conditions. Error bars show the standard deviation.

1129

1130

1131

1132 10. Improvisers were susceptible to a Lombard-like effect

1133

1134 We were interested to explore whether our participants would tend to play louder when they
1135 heard one of their co-improvisers under a louder pattern (and thus, heard this co-improviser
1136 playing louder than usual). This would be in line with the many studies on the so-called Lom-
1137 bard effect in the speech literature (Kunc et al., 2022), which showed that humans tend to
1138 speak louder as the ambient noise increase – a phenomenon that has also be identified in cho-
1139 ral singing (Tonkinson, 1994).

1140 First, we segmented the recorded audio track of each participant in 10-s windows. Each win-
1141 dow was ascribed to one of our 3 conditions (baseline, louder pattern, and softer pattern) de-
1142 pending on the manipulation (or absence of manipulation) applied to their co-improvisers
1143 during that window. Second, we eliminated the windows in which the mean RMS was < -60
1144 dB (meaning that the musician was simply not playing at that time). Averaging RMS values
1145 over each type of window revealed that musicians played +1.58 dB louder (SE = 0.43) when
1146 exposed to a louder pattern (as compared to the baseline) and +0.6 dB (SE = 0.41) louder
1147 when exposed to a softer pattern (as compared to the baseline).
1148 We then conducted a statistical analysis. As the RMS data did not follow a normal distribu-
1149 tion, we did a ranked transformation of our RMS data using the R library bestNormalize. This
1150 allowed us to run a two-scale hierarchical regression by comparing nested models, starting
1151 with a null model and adding our experimental condition (i.e., whether the window belong to
1152 the baseline, softer pattern or louder pattern category). The likelihood ratio test for model
1153 comparison was significant ($X^2=8.768$, $p=0.012$). Our model with the experimental condition
1154 as a predictor showed that musicians indeed tended to play significantly louder when exposed
1155 to a louder pattern (Estimate = 0.126, SE=0.023, $t=2.954$, $p=0.003$). Conversely, no signifi-
1156 cant effect was found for the softer pattern (Estimate = 0.048, SE=0.042, $t=1.141$, $p=0.254$).

1157

1158 Additional References

1159

1160 Angus, D., Watson, B., Smith, A., Gallois, C., & Wiles, J. (2012). Visualizing conversation
1161 structure across time: insights into effective doctor-patient consultations. *PLoS One*,
1162 7:e38014. DOI: doi: 10.1371/journal.pone.0038014

1163 Dale, R. & Spivey, M. J. (2006). Unraveling the dyad: using recurrence analysis to explore
1164 patterns of syntactic coordination between children and caregivers in conversation. *Lang.*
1165 *Learn.* (56). DOI: doi: 10.1111/j.1467-9922.2006.00372.x.

1166 Dau, T., Kollmeier, B., & Kohlrausch, A. (1997). Modeling auditory processing of amplitude
1167 modulation. I. Detection and masking with narrow-band carriers. *The Journal of the Acousti-*
1168 *cal Society of America*, 102(5), 2892-2905.

1169 King, A., Varnet, L., & Lorenzi, C. (2019). Accounting for masking of frequency modulation
1170 by amplitude modulation with the modulation filter-bank concept. *The Journal of the Acousti-*
1171 *cal Society of America*, 145(4), 2277-2293.

1172 Kunc, H. P., Morrison, K., & Schmidt, R. (2022). A meta-analysis on the evolution of the
1173 Lombard effect reveals that amplitude adjustments are a widespread vertebrate mecha-
1174 nism. *Proceedings of the National Academy of Sciences*, 119(30), e2117809119.

1175 Konvalinka, I., Gálatas, D., Bulbulia, J., Schjodt, U., Jegindo, E., Wallot, S., Van Orden, G.,
1176 & Roepstorff, A. (2011). Synchronized arousal between performers and related spectators in a
1177 fire-walking ritual. *Proceedings of the National Academy of Sciences*, 108(20). DOI:
1178 <http://www.pnas.org/cgi/doi/10.1073/pnas.1016955108>.

1179 Lira-Palma, D., González-Rosales, K., Castillo, R. D., Spencer, R., & Fresno, A. (2018). Cat-
1180 egorical Cross- Recurrence Quantification Analysis Applied to Communicative Interaction
1181 during Ainsworth's Strange Situation. *Complexity*, 2018(1), 4547029.

1182 Ponsot, E., Varnet, L., Wallaert, N., Daoud, E., Shamma, S. A., Lorenzi, C., & Neri, P.
1183 (2021). Mechanisms of spectrotemporal modulation detection for normal-and hearing-
1184 impaired listeners. *Trends in hearing*, 25, 2331216520978029.

1185 Reddy, V., Markova, G., & Wallot, S. (2013). Anticipatory adjustments to being picked up in
1186 infancy. *PLoS One*, 8:e65289. DOI: doi: 10.1371/journal.pone.0065289.

1187

1188 Relaño-Iborra, H., May, T., Zaar, J., Scheidiger, C., & Dau, T. (2016). Predicting speech in-
1189 telligibility based on a correlation metric in the envelope power spectrum domain. *The Jour-*
1190 *nal of the Acoustical Society of America*, 140(4), 2670-2679.

1191 Tonkinson, S. (1994). The Lombard effect in choral singing. *Journal of Voice*, 8(1), 24-29.

1192 Wallot, S. & Leonardi, G. (2018). Analyzing Multivariate Dynamics Using Cross-Recurrence
1193 Quantification Analysis (CRQA), Diagonal-Cross-Recurrence Profiles (DCRP) and Multidi-
1194 mensional Recurrence Quantification Analysis (MdRQA) – A Tutorial in R. *Frontiers in Psy-*
1195 *chology*, 9(2232). DOI: doi: 10.3389/fpsyg.2018.02232.

1196 Wallot, S., Monster, D. (2018). Calculation of Average Mutual Information (AMI) and False
1197 Nearest-Neighbors (FNN) for the Estimation of Embedding Parameters of Multidimensional
1198 Time Series in MATLAB. *Frontiers in Psychology*, 9. DOI:
1199 <https://doi.org/10.3389/fpsyg.2018.01679>.

1200
1201
1202
1203