

1 Running title: Speech rhythm entrainment in a dyadic reading task

2

3

4

5

6 *SPEECH RHYTHM CONVERGENCE IN A DYADIC READING TASK*

7

8 Karina Cerda-Oñate¹

9 Gloria Toledo Vega²

10 Mikhail Ordin^{3,4*}

11

12

13 ¹ Departamento de Lengua Castellana y Literatura, Universidad Católica del Maule,
14 Chile

15 ² Facultad de Letras, Pontificia Universidad Católica de Chile

16 ³ BCBL- Basque Center on Cognition, Brain and Language

17 ⁴ IKERBASQUE- Basque Foundation for Science

18

19 Acknowledgments:

20

21 This research was supported by the Spanish Ministry of Economy and Competitiveness
22 (MINECO), through the "Severo Ochoa" Programme for Centres of Excellence in R&D
23 (SEV-2015-490). KC-O acknowledges financial support from the Program for the
24 Training of Advanced Human Capital (PFCHA) at CONICYT (Scholarship n.
25 21151287/2015). The authors also thank C. García-Meza (UCSH Chile), F. Nocetti
26 (UDEEC Chile), P. Timofeeva (BCBL), J. Aguasvivas (BCBL) and Candice Frances
27 (BCBL) for their help with the data collection and synchrony measurements, and Dr.
28 Magda Altman (BCBL) for language support, proofreading and editing this manuscript.
29 We are also thankful to the editors of this special issue, who highlighted the timely and
30 important questions related to phonetic convergence in interaction.

31

32 *Address for correspondence: Mikhail Ordin, Basque Center on Cognition, Brain and
33 Language, San Sebastián, Mikeletegi 69, 20009, Spain.

34 **Abstract**

35 We tested the effect of co-presence on entrainment to speech rhythm, examining
36 differences in speech rhythm convergence during a reading task, in conditions where the
37 reading partner was present or absent. Speech rhythm was operationalized as a two-
38 level phenomenon. At a lower level, rhythm was operationalized as regularity in the
39 distribution of salient acoustic events (vowel onsets) and regularity in the duration of
40 speech intervals (consonantal and vocalic intervals). At a higher level, rhythm was
41 operationalized in terms of meter, the distribution of salience that can group syllables
42 into metrical structures, based on lexical stress and phrasal prominence. To assess the
43 impact of the presence/absence of a co-speaker on speech entrainment, we asked a
44 *model speaker* and *experimental participants* to sit side-by-side and read two texts
45 aloud at the same time while recording each speaker separately. Later, we requested
46 only *experimental participants* to read in synchrony with the recording obtained from
47 the *recurrent model speaker* during side-by-side reading and, afterwards, we asked
48 *experimental participants* to read in synchrony with a solo recording from the same
49 *recurrent model speaker*. We used a poetic text with strong meter and a narrative text
50 with weak meter as reading materials. To assess the degree of speech rhythm
51 convergence, we measured the degree of durational variability in vocalic and
52 consonantal intervals across conditions and texts and the deviation in vowel onsets
53 between chorusing readers. We found that participants make their speech more regular
54 to facilitate chorusing. Importantly, inter-speaker synchronization was substantially
55 improved during side-by-side reading compared to the condition where the speaker read
56 in synchrony with a recording obtained during side-by-side synchronous reading, even
57 though the acoustic signal received by the experimental participants was the same in
58 both conditions. We also found that the text with strong meter modulated the success of
59 speech rhythm convergence only in the most challenging condition, requiring reading in
60 synchrony with a recording of the *model speaker* reading solo. These results show that
61 co-presence plays a crucial role in inter-speaker entrainment, while meter can provide
62 additional benefits for chorusing in more challenging conditions.

63

64 *Keywords:* phonetic convergence, speech rhythm, dyadic reading, rhythm convergence,
65 rhythm synchronization

66

67 **Introduction**

68 Prosody is an essential aspect of the speech signal that is used to convey post-lexical
69 meaning and is modulated by speaker gender, social status, native language and dialect,
70 discourse, and emotional state, among others (Crystal, 1969; Ladd & Cutler, 1983;
71 Labov, 2006; Szczepek Reed, 2010; Clopper & Smiljanic, 2011). Prosody includes all
72 suprasegmental phenomena, such as stress, tone, intonation, and rhythm, produced by
73 the interplay of pitch, intensity, and duration in stretches of speech larger than a
74 segment (Lehiste, 1976; Ladd & Cutler, 1983; Nooteboom, 1997; Fletcher, 2010).

75 **Rhythm as a suprasegmental speech phenomenon**

76 As a suprasegmental speech phenomenon, speech rhythm is conceptualized and defined
77 in different ways (Moore, 2012; Turk & Shattuck-Hufnagel, 2013; Smith et al., 2014;
78 Nolan & Jeon, 2014; Cummins, 2015). Since we do not focus on establishing a
79 theoretical definition for rhythm, we will not delve into terminological debates and will
80 provide an operational definition of speech rhythm that emphasizes those aspects of the
81 phenomenon that are relevant for our study. In this context, speech rhythm will be
82 operationalized in terms of the duration and variability of sub-syllabic units:
83 *consonantal intervals* (single or multiple consonants that happen consecutively and are
84 not interrupted by vocalic segments) and *vocalic intervals* (vowels uninterrupted by
85 pauses). Lower variability scores thus reflect higher regularity in the duration of vocalic
86 and consonantal intervals, and higher regularity in the temporal distribution of
87 *perceptual centers* – temporal reference points when an acoustic event (e.g., syllable, or
88 vowel onset corresponding to a sharp increase in intensity level) is psychologically
89 perceived to occur (Marcus, 1981) – in the speech signal. These perceptual centers (p-
90 centres) determine the perception of regularity in the acoustic signal (Scott, 1998).

91 **Entrainment and speech rhythm convergence**

92 Tapping, speech-cycling, and dancing tasks reveal the human ability to entrain to a
93 rhythmic sensory input signal (Cummins & Port, 1998; Repp, 2005; Eerola et al., 2006,
94 Merchant et al., 2015). In the case of speech rhythm, it has been found that native
95 speakers can entrain to perceptually salient acoustic events that may correspond to more
96 or less regularly distributed stressed syllables or vowel onsets, (Marcus, 1981; Vos et
97 al., 1995). These events are produced by another speaker when reading aloud. To
98 facilitate synchronization by a partner, a speaker will (intentionally or unintentionally)
99 make the distribution of these salient events more regular, hence more predictable; this
100 adaptation does not require any prior training (Cummins, 2002, 2003, 2007, & 2009;
101 Bowling et al., 2013). It has been suggested that side-by-side synchronous reading

102 allows for convergence of rhythmic patterns in speech (Pardo, 2010; for a full review of
103 research on phonetic convergence in speech imitation tasks, see Pardo et al., 2017).
104 Pardo (2010) and Pardo et al. (2017) consider phonetic convergence to be a general
105 phenomenon, which occurs not only in regard to speech rhythm, but also with the
106 formant structure of vowels, pitch, speech tempo, VOT on plosives, etc. In our study,
107 we focus on the convergence of temporal patterns that a) result in specific rhythmic
108 organization of speech; and b) lead to the rhythmic convergence, and c) allow speakers
109 to use rhythmic predictability to synchronize speech in dyadic reading tasks.

110 Earlier studies have already reported convergence effects during synchronous
111 reading (Cummins 2002; 2009). In Experiment 2, Cummins (2002) asked three
112 participants to synchronize their reading with a recording made either in a synchronous
113 reading condition – when two speakers were reading the same text together, and the
114 speech of one of these speakers was recorded – or with a recording made when the
115 speaker read the text alone. A higher degree of convergence was found for the
116 synchronous reading condition for two out of three participants. In Cummins (2009),
117 this effect was replicated with four participants who had to synchronize with two
118 different model speakers. Three out of four participants synchronized better with the
119 speech originally obtained in dyadic reading conditions.

120 In the current study, we decided to run a hypothesis-driven experiment to
121 confirm this finding in a larger sample of participants (N=30) and generalize the
122 findings over different language populations. Importantly, we sampled our participants
123 from a population with an ambient language (Spanish), which is rhythmically different
124 from English – the native language of participants in Cummins' studies (2002; 2009) –
125 so we could generalize results over a broader range of languages (Spanish exhibits more
126 regular distribution of vowel onsets, stressed syllables, more equal vocalic and
127 consonantal intervals, etc.; White & Mattys, 2007). We believe that rhythm provides
128 cues that facilitate synchronization, and thus the effect should not be limited to
129 languages with a specific type of temporal organization. This hypothesis, however, had
130 not yet been tested on languages which are rhythmically different from English (where
131 this effect was first established). Also, we decided to look at synchronization on
132 different types of text materials – poetic and narrative texts – to understand whether this
133 effect is limited to a certain type of reading material or whether it might increase when
134 the material contains stronger inherent rhythm (in this case, a poetic text) that could
135 potentially provide further cues for synchronization. Most importantly, we introduced

136 another condition: a live speaker. Cummins (2002; 2009) provided some empirical
137 evidence that at least English speakers achieve better rhythmic convergence with
138 recordings obtained in synchronous conditions than with recordings of solo readers. We
139 introduced a live speaker as an additional condition because a recorded model speaker
140 cannot make online adaptations to a partner while the task is being performed; only a
141 live speaker can adapt their speech online to achieve better convergence with their
142 partner. Thus, we posit that in dyadic reading both readers will dynamically change
143 their speech at the same time, (1) providing cues that help their live partner increase
144 synchrony, and (2) trying to converge their own speech with the partner's speech. These
145 dynamics might further enhance convergence of temporal patterns, as we aimed to test
146 in this study. A more technical formulation of our hypothesis is provided below. But
147 first, we introduce two core concepts/terms that are key to our theoretical and
148 methodological framework: beat and meter.

149 **Beat, meter, and their role in speech rhythm convergence**

150 Beats, or perceptually salient acoustic events in the speech stream – p-centres – may be
151 perceived as more or less regularly distributed (Marcus, 1981; Vos et al., 1995). A
152 regular beat licenses predictions for the timing of following beats. This facilitates
153 synchronization with speech as p-centers correspond to more prominent speech
154 segments, for example, vowel onsets (thus each syllable in Spanish would correspond to
155 one p-centre). Beats can be organized into hierarchical patterns based on their relative
156 strengths (London, 2004; Large & Snyder, 2009; Fitch, 2013) to create a hierarchical
157 temporal grid, which we refer to as meter. For example, alternations between stressed
158 and unstressed syllables can be organized into iambic, trochaic, or more complex
159 metrical patterns.

160 Meter has been found to play a significant role in the entrainment of motor
161 output to sensory input when dancing, marching, tapping, clapping, singing, and other
162 types of behavior that require the perception and use of expectancy cues (Hannon &
163 Johnson, 2005; Fitch, 2013). Pitt and Samuel (1990) showed that metrical expectancy
164 cues are weaker in the context of normal narrative speech and that meter only leads to
165 better performance when beat cues are insufficient, e.g., when reading a list of words, or
166 when context encourages the listener to pay closer attention to the meter.

167 The relationship between timing regularity – regular distribution and thus the
168 predictability of when the next syllable will occur – and metrical expectancy in speech
169 perception was also explored by Quené & Port (2005). In one experiment, participants

170 had to listen to a set of words in which the onsets of the stressed syllables were either
171 regularly or irregularly distributed (all words within sets had the same meter). In the
172 second experiment, participants had to listen to a set of words. In one condition, the
173 words had either exclusively trochaic or exclusively iambic meter, making stressed
174 syllables onsets predictable. In the other condition, iambic and trochaic meters were
175 used interchangeably, making the stressed syllable onsets less predictable. Participants
176 in both experiments had to do a phoneme monitoring task, which consisted of listening
177 to a set of words or sentences and pressing a button as soon as the target phoneme was
178 detected. Quené & Port (2005) found a main effect for timing regularity (meaning that
179 reaction time was shorter when the syllable onsets were regularly distributed) but not
180 for metrical expectancy (indicating that reaction time was not modulated by whether
181 words had the same or different meter).

182 Meter can provide expectancy cues for the timing of the following stressed
183 syllable, but the salience of these cues – as tested in phoneme monitoring tasks –
184 remains unclear (Pitt & Samuel, 1990 found a facilitatory effect of meter consistency;
185 Quené & Port, 2005 did not), and their role in the precision and ease of rhythmic
186 convergence has not been established. We therefore decided to include materials with
187 both weak and strong meter in our experiment, choosing a narrative text with weak
188 meter (Text 1) and a poetic text with strong meter (Text 2). We believe that the
189 organization of beats into metrical patterns may provide yet another set of expectancy
190 cues and further facilitate the convergence of rhythmic patterns in speech. These ideas
191 will be further explained in the next section.

192 **Hypothesis and predictions**

193 In this study, we explored speech rhythm as an acoustic phenomenon that is modulated
194 by human co-presence and variation in the degree of meter. To explore human co-
195 presence, we tested the effect of the presence/absence of an interlocutor during a
196 synchronous reading aloud speech task. We posited that the presence of an interlocutor
197 would induce speakers to provide better expectancy cues in an effort to encourage
198 enhanced rhythmic synchronization; this more regular speech rhythm in speech should,
199 in turn, be reflected in lower scores for the rhythm variability metrics and higher %V
200 measures (Ramus et al., 1999; Grabe & Low, 2002; Dellwo, 2006; White & Mattys,
201 2007).

202 We adopted the premise that perception of rhythmicity in speech is based on the
203 extraction of p-centers (Morton, Marcus, & Frankish, 1976), which can serve as

204 expectancy cues, allowing the listener to predict the onsets of vowels, syllables, and
205 higher-level prosodic constituents (Lehiste, 1973). Better expectancy cues, such as the
206 ones provided by a text with strong meter, should allow the listener to better entrain
207 their speech production – motor output – to the speech produced by an interlocutor, that
208 is, the sensory signal emitted by another biological organism. If the timing of p-centers
209 encourages the perception of temporal regularity, then speech rhythm entrainment
210 should also be enhanced (Large & Snyder, 2009; Vuust & Witek, 2014).

211 As for the role of meter, we expected that the text with stronger meter would
212 allow for better timing convergence of stronger beats during speech production, further
213 enhancing synchronization of verbal behavior between partners who read the poetic text
214 compared to the narrative text.

215 To summarize, we hypothesized that co-presence would play a significant role
216 in increasing rhythmic convergence. Bi-directional attempts at speech synchrony (when
217 both speakers dynamically change their speech patterns) result in greater temporal
218 convergence than unidirectional attempts to synchronize speech with a recording. We
219 believe that regular distribution of vowel onsets (and thus syllables) will provide
220 expectancy cues to align timing patterns in dyadic reading. Metrical prominence would
221 affect the degree of speech rhythm convergence, by reinforcing expectancy cues – e.g.,
222 promoting higher regularity between vowel onsets in stressed syllables. We also
223 predicted that, even if the text with stronger meter enhanced the degree of coupling
224 between speakers, the co-presence of reading partners would be the strongest predictor
225 for increased levels of synchronization.

226 To test these hypotheses, we set up an experimental study asking native Spanish
227 speakers to read a narrative text with weak meter (Text 1) and a poetic text with strong
228 meter (Text 2) in synchrony. The *experimental participant* either read together in a live
229 session with the co-present *model speaker*, or read along with a recording of the model
230 speaker obtained when the *model speaker* was reading the texts solo, and when the
231 *model speaker* was reading in synchrony with the *experimental participant*. We
232 measured the degree of speech rhythm convergence by estimating synchrony measures.
233 Synchrony was operationalized in terms of temporal matching of onset starting points
234 for consonantal and vocalic intervals. Vowel and consonants starting points were
235 compared across conditions and texts for the *model speaker* and the *experimental*
236 *participants*. All of these measurements are explained in the Methods and Materials
237 section below.

238 **Methods and materials**

239 This experiment was conducted following the research policies of the Social Sciences
240 and Humanities Ethics Committee at the Pontifical Catholic University of Chile (PUC)
241 and the Chilean National Commission for Scientific and Technological Research.
242 Participants provided written informed consent and were informed that they could
243 withdraw from the experiment at any time without any consequences. All data were
244 stored anonymously.

245 **Participants**

246 Thirty experimental participants were recruited using different social media websites.
247 22 participants were female (mean age = 23.27; SD = 3.80) and 8 participants were
248 male (mean age = 24.5; SD = 3.11). All participants were native, monolingual speakers
249 of Chilean Spanish and reported no previous experience with any activities that might
250 train reading synchronization abilities, such as ensemble vocal music performances or
251 theatrical performances requiring people to speak at the same time and in synchrony.
252 None of the participants had ever resided outside of Chile for longer than six months.
253 This background demographic information was obtained using a background
254 questionnaire that was completed by each participant. The model speaker was a female
255 native Chilean Spanish speaker.

256 **Materials**

257 The experiment included three reading materials in Spanish, which were read in 4
258 different conditions: 1) Warm-up text, Text 0, a list of 16 sentences; 2) A narrative text
259 (text with weak meter), Text 1, which is a Spanish version of *The North Wind and the*
260 *Sun* (Coloma, 2016); 3) A poetry text (text with strong meter), Text 2, which is a poem
261 written in octosyllabic verse, consisting of two stanzas (Parra, 1988). The recordings
262 obtained from participants reading Text 0 were not used as speech data, since Text 0
263 was only used for the warm-up phase. These materials are further explained in Table 1
264 below.

265 INSERT TABLE 1 SOMEWHERE HERE

266 **Design**

267 Data were obtained from the *model speaker* and *experimental participants* reading in 4
268 different conditions. After running the first condition, the order of other conditions was
269 counter-balanced across participants. These four conditions are described in Table 2:

270 INSERT TABLE 2 SOMEWHERE HERE

271 In the first condition or *Synchrony-live condition* (1), participants read aloud and in
272 synchrony with the *recurrent model speaker*, a naïve female monolingual native speaker
273 of Spanish, who had previous experience in acting and was unaware of the objectives of
274 the research. The *recurrent model speaker*, as the name indicates, did not vary across
275 participants; in other words, the same *recurrent model speaker* read along with each one
276 of the 30 *experimental participants*. In the second condition or *Synchrony with*
277 *recording from live condition* (2), *experimental participants* listened to the recording of
278 the *recurrent model speaker* reading with them in the first condition and read in
279 synchrony with that recording. In the third condition or *Synchrony with recording from*
280 *non-live condition* (3), all *experimental participants* read the texts in synchrony with a
281 recording obtained from the model speaker reading the texts alone. The individual
282 recording used in the third condition was the same for all 30 *experimental participants*.
283 Finally, in the fourth condition or *Solo recording condition* (4), *experimental*
284 *participants* read aloud individually. Therefore, three out of the four conditions were
285 synchronic reading conditions.

286 **Procedures**

287 All recordings for each *experimental participant* were made during one session, which
288 lasted around 40 minutes. *Experimental participants* were asked to read the texts in all
289 four conditions. A warm-up material (Text 0) was used which was used to familiarize
290 participants with the task (and was left unanalyzed); a narrative text with weak meter
291 (Text 1); and a poetry text with a strong meter (Text 2). The recording for condition 4
292 (the model speaker reading alone) was made prior to the experiment. The *recurrent*
293 *model speaker* read each of the three texts with the participants in conditions 1, 2 and 3
294 (condition 1 always preceded condition 2; otherwise the order of conditions and texts 1
295 and 2 within conditions was counterbalanced across participants). The only instruction
296 given throughout the experiment to the *experimental participants* was to try to "read in
297 synchrony with their partner and to avoid repeating after the other participant in all
298 synchronous conditions." The model speaker, who was unaware of her experimental
299 role, was never advised on correct pronunciation, timing, or any other aspect related to
300 rhythm, beat, or periodicity. It is also important to note that the *recurrent model speaker*
301 and the *experimental participants* were unfamiliar with any concepts related to speech
302 rhythm, speech convergence, or synchronous speech.

303 As for the *Synchrony-live condition* (1), the *recurrent model speaker* and
304 *experimental participant* sat side by side and read aloud synchronously. In this

305 condition, the reading partners were not able to look at each other because they were
306 both focused on reading the texts that were placed in front of them, positioned at eye
307 level on a book stand. They sat facing in the same direction, which excludes the
308 possibility of seeing something peripherally without simultaneously turning their heads.
309 Thus, we can assume that visual interaction had a minimal impact, if at all, on the
310 results we report. For the *Synchrony with recording from live condition* (2), and
311 *Synchrony with recording from non-live condition* (3), the *experimental participant*
312 wore headphones to listen to the recording from the *recurrent model speaker*. The
313 volume of the recording was adjusted to a comfortable threshold defined by each
314 participant before starting the experimental task. Finally, during the *Solo-recording*
315 *condition* (4), the participant read alone.

316 All recordings took place in a quiet room at PUC. Two head-mounted dynamic
317 uni-directional (cardioid) WH20XLR Shure microphones were used to simultaneously
318 record the speech of both the *recurrent model speaker* and the *experimental participant*;
319 these microphones were chosen to avoid picking up any signal from the co-speaker.
320 Also, to avoid any interference, participants sat side-by-side with a distance of about 60
321 cm between them. Sony MDR-NC8 Headphones (with noise canceling switched off)
322 were used to play the recordings of the *model speaker* during the *Synchrony with*
323 *recording from live condition* (2) and *Synchrony with recording from non-live condition*
324 (3). The noise canceling was switched off. The noise cancelling would have cut
325 participants off from sounds in the ambient auditory environment including their own
326 voice, thus disrupting participants' ability to monitor their own speech production. With
327 noise-cancelling switched off, participants heard their voice not through the
328 headphones, but as a part of the ambient auditory environment. The experimental
329 participants and the model speaker were recorded using the same recorder in a single
330 WAV file as left and right channels, in all conditions; a TASCAM DR-40 recorder
331 (sampling rate = 44.1 kHz, 16 bit) was used for this purpose. It is important to note that
332 all recordings from the experimental participants were extracted during the same
333 session; this helped to control for any unwanted noise or interference in the *Synchrony*
334 *with recording from live condition* (2). We also used windscreens to minimize the
335 possibility that close-up uni-directional microphones would pick up the partner's voice.

336 **Data analysis**

337 The assessment of rhythmic duration took both consonantal and vocalic intervals into
338 consideration. The speech data was labeled using the orthographic transcription of the
339 reading materials (Text 1 or weak meter; Text 2 or strong meter) as a reference. For this,
340 the acoustic software analysis *Praat* (Boersma & Weenink, 2016) and the *Praat* plugin
341 for forced text-to-speech alignment *EasyAlign* (Goldman, 2011) were used. Notably, the
342 *EasyAlign* plugin was employed to segment and label the recordings into phonemic
343 units that were based on the orthographic transcription of the texts the participants read
344 during the tasks. EasyAlign is a Praat plug-in for semi-automatic aligning of segmented
345 words, syllables, and phones to an acoustic signal based on the audio recording and its
346 orthographic transcription. First, orthographic texts were manually segmented into
347 utterances (identical segmentation of texts into utterances was applied to each
348 recording). The second stage was a grapheme-to-phoneme conversion, in which the
349 orthographic transcription for each utterance was converted into a sequence of
350 phonemes (choosing from a set of possible pronunciation variants); optional phonemes
351 were marked by an asterisk. After this, an expert phonetician (the first author) compared
352 the sequence of phonemes against the audible recording and made the necessary
353 adjustments for individual idiosyncrasies. At the third stage, a verified phonetic
354 transcription was aligned to the acoustic signal using HTK-toolkit, an algorithm based
355 on Markov Hidden Chains (<http://htk.eng.cam.ac.uk/>), to produce segmentation of the
356 recording into phones, syllables, and words. The durations of the phones were used for
357 further analysis. Details of the algorithm can be found in Goldman (2011).

358 The labeled data for each participant, *Condition*, and *Text* was saved into a single Praat
359 TextGrid file. This aligned phonetic labels with the acoustic waveform, using different
360 tiers for each transcription. The transcriptions were all in SAMPA for Spanish, a
361 machine-readable phonetic alphabet (Wells, 1996). The annotations were then manually
362 checked to correct the most obvious technical errors. EasyAlign works reliably with
363 Spanish speech produced by native speakers (both Castilian and Latin American
364 Spanish varieties are targeted). The main motivation for using EasyAlign in our study
365 was to automate grapheme-to-phoneme conversion. Automatizing this process improves
366 the reproducibility of the results because the segmentation decisions are not affected by
367 subjective choices, such as those related to the role of liquids and non-nuclear vocoids.
368 Subjective decisions tend to differ across annotators and could be biased by an
369 experimental hypothesis, adding to the inconsistency of the results across studies. Here
370 instead, purely acoustic parameters were used to detect the phoneme boundaries.

371 To assess the temporal patterns of speech, we computed the durational
 372 variability of consonantal (C) and vocalic (V) intervals using *Correlatore* (automatic
 373 rhythm analysis software; Mairano & Romano, 2010). This variability is captured,
 374 among other measures, by the coefficients of variation, VarcoC and VarcoV (White &
 375 Mattys, 2007; Dellwo, 2006); the coefficient is the standard deviation divided by the
 376 mean duration of intervals for each sentence. Fluctuations in standard deviation depend
 377 on the mean duration of intervals, thus standard deviation will be modulated by
 378 individual speech rate or stylistic changes in speech rate even between different
 379 sentences produced by the same speaker (speech delivered at a faster rate will have
 380 shorter overall speech interval durations, thus smaller standard deviations). Varco
 381 measures thus reflect variability, while controlling for between- and within-speaker
 382 differences in speech rates (Dellwo, 2006). We also estimated %V, or percentage of
 383 vocalic material in a sentence (this measure is also reported to be robust to tempo
 384 variation). Varco and %V measures are referred to as global metrics because they
 385 capture durational variability averaged across the whole utterance. Local metrics, on the
 386 other hand, capture variability in the duration of C and V intervals in a pairwise fashion
 387 (respectively, CnPVI and VnPVI), assessing how the second interval differs from the
 388 first, the third from the second, and so on (Grabe & Low, 2002). The difference in
 389 interval durations between each consecutive pair is normalized by the mean duration of
 390 that pair, and the sum of these normalized differences is then divided by the number of
 391 consecutive pairs. The result is often multiplied by 100 to convert it into the percentage
 392 of overall variability within a sentence, in which interval durations fluctuate from one
 393 pair to the next. The formula is given below (1). If a consonantal or vocalic interval was
 394 interrupted by a pause exceeding 150ms, we considered them to be two separate
 395 intervals. These metrics were calculated for each sentence, and then an average was
 396 calculated for each speaker and text. %V, Varco, and nPVI measures have been widely
 397 used in studies on speech rhythm; using these metrics makes it possible to compare and
 398 discuss the results of this research in light of other studies on speech rhythm.

$$399 \quad nPVI = 100 * \left(\sum_1^{m-1} \frac{d_k - d_{k+1}}{(d_k + d_{k+1})/2} \right) / (m - 1) \quad (1)$$

400 where m is the number of intervals, and k is the ordinal number of each consecutive interval.

401 To measure convergence, we estimated the difference between the onset of each speech
 402 segment delivered by the two readers. To assess synchronization between speakers, we
 403 computed synchronization metrics for segment onsets using an in-house Python 3.6

404 script (see Supplementary material) to compare and calculate vowel and consonant
405 onset time differences between the Praat TextGrids of the model speaker and the
406 experimental participant, both recorded by the same recorder in a single WAV file as
407 left and right channels. The TextGrids from the model speaker and experimental
408 participants were compared as two columns. If the segments in the columns were the
409 same, the onset difference was calculated. We only included this data in further analyses
410 if 85% of the compared segments matched, otherwise we considered synchronization
411 had failed (this happened when convergence was too weak and readers could not align
412 with onsets of at least 85% of the segments within a sentence; we discarded 5.5% of
413 speech data per pair).

414 The rhythm metrics were calculated for each utterance and then averaged for
415 each participant across utterances. The averages were used in a repeated-measures
416 ANOVA (in SPSS v.19) to assess the differences between conditions: *Synchrony-live*
417 (1); *Synchrony with recording from live* (2); *Synchrony with recording from non-live*
418 (3); and *Solo-reading* (4). We applied Greenhouse-Geisser correction for sphericity
419 violations and report corrected p-values and degrees of freedom, with the a-priori alpha
420 level set to 0.05. *Condition* and *Text* (weak meter [Text 1] versus strong meter [Text 2])
421 were included as factors in each analysis. The interaction *Condition* by *Text* was
422 included in the models. For pairwise comparisons, paired t-tests (with Bonferroni
423 correction) were used.

424 **Results**

425 **Vowels**

426 **Vowel percentage.**

427 We found a significant effect of *Condition* and *Text* (refer to Table 3 for the statistical
428 parameters). The effect size was substantial for both factors; the interaction was not
429 significant. Further pairwise comparisons (with the Bonferroni correction) indicated that
430 %V was significantly higher in the *Synchrony-live* (1) than in the other conditions. We
431 also observed a significantly higher %V in the text with strong meter (Text 2) than in
432 the text with weak meter (Text 1) across all conditions. Figure 1 displays the differences
433 in %V between conditions and text types.

434 INSERT FIGURE 1 SOMEWHERE HERE

435 **VarcoV.**

436 The analysis showed a significant effect of *Condition* and *Text* with a substantial effect
437 size for both factors; interaction between the factors was not significant. Pairwise

438 comparisons revealed that VarcoV was higher in the *Synchrony with recording from*
 439 *non-live* (3) than in the other conditions (this measure was not different between other
 440 pairs of conditions). The measure scores were lower in the text with strong meter (Text
 441 2) than in the text with weak meter (Text 1). This pattern is displayed in Figure 2.

442 INSERT FIGURE 2 SOMEWHERE HERE

443 **VnPVI.**

444 We observed a significant but moderate effect of *Condition* and a substantial effect of
 445 *Text* on the nPVI scores with no significant interaction between the factors. Pairwise
 446 comparisons showed that the only difference for the factor *Condition* was between *Solo*
 447 *recording* (4) and *Synchrony with recording from non-live* (3) with higher scores in
 448 condition 3 than 4. Again, the nPVI scores were lower in the text with strong meter
 449 (Text 2) compared to the text with weak meter (Text 1). Figure 3 displays this pattern.

450 INSERT FIGURE 3 SOMEWHERE HERE

451 **Consonants**

452 **VarcoC and CnPVI.**

453 We did not observe a significant effect of *Condition* on VarcoC or CnPVI scores (see
 454 Table 3). However, the effect of the factor *Text* on VarcoC was significant and
 455 substantial. No significant interaction between the factors was observed. The effect on
 456 CnPVI was also statistically significant, but unsubstantial ($\eta_p^2=.03$). Durational
 457 variability of consonantal intervals at the timescale of whole utterances – captured by
 458 VarcoC – was higher on the text with strong meter (Text 2, poetic) than the text with
 459 weak meter (Text 1, narrative). Pairwise durational variability (changes in durational
 460 variability between successive pairs of consonantal clusters) was similar in the two texts
 461 (the effect of *Text* on consonantal PVI was negligible). This pattern is displayed in
 462 Figures 4 and 5.

463 INSERT FIGURE 4 SOMEWHERE HERE

464 INSERT FIGURE 5 SOMEWHERE HERE

465 INSERT TABLE 3

466 **Synchrony data**

467 **Onset synchrony**

468 To analyze the differences in vowel onset synchrony between poetic and narrative texts,
 469 and between the three different synchronic conditions (conditions 1, 2 & 3), we ran a
 470 repeated-measures ANOVA with *Condition* and *Text* as factors and onset differences as

471 the dependent variable (onset difference was calculated by subtracting the onset time of
 472 the participant from the onset time of the model speaker). The factor *Condition* had
 473 three levels: *Synchrony live* (1), *Synchrony with recording from live* (2), and *Synchrony*
 474 *with recording from non-live* (3). We report the corrected significance values and
 475 degrees of freedom in Table 4 below and display the pattern of results in Figure 6
 476 (negative onset differences in conditions 2 and 3 confirmed that experimental
 477 participants were following the model). The analysis showed a significant effect of
 478 *Condition*; pairwise comparisons (Bonferroni corrected) showed that vowel onset
 479 differences were smaller in the *Synchrony live condition* (1) than in the other two
 480 conditions, and synchronization was significantly stronger when speakers were
 481 shadowing the recording made in the synchrony reading condition than the recording
 482 made in the solo reading condition. We did not observe a significant effect of *Text*,
 483 however, we did observe a significant if weak interaction between *Text* and *Condition*,
 484 $F(2, 42) = 3,402, p = .04$ (uncorrected), $\eta_p^2 = .018$ (noteworthy, $p = .06$ after applying
 485 Greenhouse-Geisser correction, which signals that further interpretation of this
 486 interaction should be considered with caution and requires verification in further
 487 studies). To understand this interaction, we performed 2-tailed paired t-tests comparing
 488 synchronization on segment onsets for poetic vs. narrative tests in three different
 489 conditions. The tests showed that in the *live* condition ($p=1.0$) and in *Synchrony with*
 490 *recording from live* condition ($p=1.0$), the difference was not significant, while in the
 491 *Synchrony with recording from non-live* condition, synchronization was stronger for the
 492 poetic than the narrative text ($p=.033$, corrected).

493 INSERT TABLE 4 SOMEWHERE HERE

494 INSERT FIGURE 6 SOMEWHERE HERE

495 **Effect of Condition versus model speaker training**

496 As the model speaker performed the task multiple times, it is possible that her
 497 performance improved over the course of the experiment. If so, better synchronization
 498 in the live condition could have been caused by progressively better entrainment by the
 499 model speaker to the experimental participants rather than by co-presence of a live
 500 partner allowing for bi-directional alignment. To ascertain whether this was the case, we
 501 performed ranked correlations between the order of the participant in the experiment
 502 (the higher the order, the more practice could have benefitted the model speaker) and
 503 onset differences. Strong and significant correlations between the ordinal variable and

504 onset differences would signal that the precision of synchronization was associated with
505 the duration of the experiment (i.e., longitudinal effect on the model speaker's skills).
506 The correlations were run separately for poetic ($\rho = -.11$, $p = .566$) and narrative texts
507 ($\rho = .19$, $p = .299$). We found weak and non-significant correlations, showing that the null
508 hypothesis was as likely as the alternative hypothesis (in other words, precision of onset
509 synchronization in the live condition was not associated with the progression of the
510 experiment). To estimate support for the null hypothesis, we calculated Bayes factors
511 for the correlation coefficients ($BF_{10} = .23$ for the poetic text and $BF_{10} = .43$ for the
512 narrative text). This showed moderate but reliable evidence in favour of the null
513 hypothesis and allowed us to exclude the possibility that better synchronization between
514 readers in the live condition occurred because the model speaker progressively honed
515 her speech entrainment skills over the course of the experiment.

516 **Summary of the results**

517 The data showed that the durational variability of vocalic intervals was
518 consistently lower (i.e., VarcoV and nPVI scores were lower) in the text with strong
519 meter (Text 2) than in the text with weak meter (Text 1). At the same time, the
520 durational variability of consonantal intervals was significantly higher in the text with
521 strong meter (Text 2). These patterns suggest that speakers tried to make the vowel
522 durations more or less equal in the text with strong meter (Text 2), possibly for stylistic
523 reasons, and vowels were stretched more naturally. This modulation is also reflected in
524 the higher percentage of vocalic material in the text with strong meter compared to the
525 text with weak meter – %V was lower in the text with weak meter (Text 1). However,
526 this behavior also led to increased variability in the duration of inter-vocalic intervals.
527 Thus, it seems that speakers tried to quasi-equalize the intervals that could form beats,
528 i.e., intervals that typically had a higher level of intensity – vowels – rather than
529 normalizing interbeat intervals.

530 We also observed unequal variance in the measures of vocalic variability, i.e.,
531 significant ANOVAs, between conditions, but we find this result difficult to interpret.
532 We observed lower variability in the duration of vocalic intervals in the *Synchrony-live*
533 condition (1) and the *Synchrony with recording from live condition* (2). It is possible
534 that this result reflects the effort made to provide a partner with an opportunity for
535 efficient synchronization, and from an attempt to synchronize with partners' utterances
536 by changing vocalic durations (vowels can be stretched and compressed easily).

537 However, in the *Solo recording* condition (4), the durational variability of the vowels
538 was lower than in *Synchrony with recording from non-live condition* (3), although we
539 would not expect differences between these two conditions (and if we had expected any
540 difference, it would have been in the opposite direction). The variability in the duration
541 of consonantal clusters between conditions is not significantly different, indicating that
542 people do not try to sync their inter-vowel intervals, but rather take advantage of the
543 stretchability/compressibility of the vowel intervals during synchronization.

544 The segment onset differences did not vary significantly between text types,
545 suggesting that the synchronization strategy is probably the same for texts with weak
546 meter (Text 1) and strong meter (Text 2).

547 As for between-speaker synchrony, we observed substantial enhancement in
548 onset convergence in the *Synchrony-live condition* (1) compared to synchronization
549 with the recordings, irrespective of whether these recordings were made in the live
550 (condition 2) or non-live condition (condition 3). This trend proves that the co-presence
551 of a live partner is important because only in the *Synchrony-live condition* (1) did both
552 partners work in synergy to both provide and perceive synchronization cues and use
553 them to amend their speech production, probably, as suggested above, by manipulating
554 the duration of vocalic intervals rather than consonantal clusters.

555 Teasing apart the significant interaction between *Condition* and *Text* factors on
556 the synchronization measure, we showed that strong meter enhanced synchronization
557 only in the most challenging condition, the condition which required reading in
558 synchrony with a solo recording from the recurrent model speaker, and that it did not
559 prove to be beneficial in other conditions.

560 **Discussion**

561 Our findings show that when speakers read with another speaker or with a recording of
562 another speaker, they regularize the durations of vowels, consonantal intervals, and
563 syllables, thus making the distribution of vowel and syllable onsets – p-centers or beats
564 – more isochronous and thus more predictable. Enhanced expectancy cues at the beat
565 level, in turn, allow for better inter-speaker synchronization at vowel onsets, as reflected
566 in the higher synchrony measures found in our study.

567 We found that the temporal variability of vocalic intervals differed across
568 conditions. This suggests that when people need to provide expectancy cues to enable
569 better synchronization – even if they are unaware of their behavioral patterns – the
570 durational variability of vowels becomes lower. Durational values for both the

571 *Synchrony-live condition* (1) and *Synchrony with recording from live condition* (2) were
572 more or less maintained for both texts. Interestingly, the durational variability of
573 consonantal intervals was not modulated by *Condition*. The fact that experimental
574 participants did not manipulate consonant intervals is not surprising, since the nucleus
575 of the syllable in Spanish is always a vocalic interval. Therefore, keeping the
576 consonantal structure constant could be considered a way to increase the chances of
577 achieving synchrony for vowel onsets.

578 The results indicate that the metric structure of Text 2 made a significant
579 contribution to the regularity of beat distribution in the speech signal, as reflected by the
580 higher values of %V and lower durational variability of vowels, both pairwise (VnPVI)
581 and as measured across whole utterances (VarcoV). Interestingly, for VarcoC, we found
582 the opposite effect. For CnPVI, we did not observe significant differences between the
583 text with weak meter (Text 1) and the text with strong meter (Text 2). This suggests that
584 a prominent meter (that is, regular and predictable distribution of stressed syllables) in a
585 poetic text leads to better inter-speaker coherency in the pronunciation of consonantal
586 clusters and preserves syllabic structure. These results indicate that the mere presence of
587 strong metrical structure makes the distribution of beats more regular. This not only
588 provides additional expectancy cues at the level of meter but also enhances expectancy
589 cues at the level of beat.

590 It is important to note some limitations on generalizing our findings. Spanish is a
591 language that does not have great durational contrasts between stressed and unstressed
592 vowels. In this respect, it differs from languages like English, German, and Russian,
593 whose native speakers likely do not benefit from the effect of meter in the same manner
594 as Spanish speakers do. We believe that in these other languages, different strategies
595 might be used to enhance synchronization. These could include different behaviors
596 related to consonantal intervals or combining synchronization of vocalic and inter-
597 vocalic intervals, something that was not observed in our data. Besides, our conclusions
598 are based on two texts (a narrative and a poem; one text per reading style). Although we
599 had a large sample of participants, which justifies generalization across speakers sharing
600 the same language, the degree to which the results are generalizable across texts has not
601 been tested. The differences that we attributed to the effect of meter might also be due
602 to differences in the phonetic material – idiosyncrasies in the phonotactic composition
603 of the sentences in the narrative and poetic texts. Thus, the effect of meter should be
604 confirmed by future studies and for other languages.

605 Our data demonstrated that vowel and consonant onset time differences were
606 modulated by *Condition*. Inter-speaker alignment of vowel and consonant onsets was
607 the highest in the *Synchrony-live condition* (1). Importantly, in the *Synchrony-live*
608 *condition* (1), the deviation among values for onset timing difference was very small
609 (Figure 6). These findings prove the boosting effect of a live partner on rhythmic
610 synchronization with the speech signal. For example, in the case of the *Synchrony with*
611 *recording from live condition* (2), participants were not able to achieve the same degree
612 of synchrony that they achieved in the *Synchrony-live condition* (1). These results
613 indicate that co-presence provided additional expectancy cues beyond the signal-
614 inherent acoustic correlates of speech rhythm that did not differ between the *Synchrony-*
615 *live condition* (1) and *Synchrony with recording from live condition* (2).

616 The boosting effect of co-presence on speech synchronization can be explained
617 by the mutual efforts of the interactants, who dynamically modified their rhythmic
618 patterns online and adjusted their speech onsets. In the case of the *Synchrony from*
619 *recording live condition* (2), participants had to synchronize their reading with a signal
620 that did not change in real-time, while in the *Synchrony-live condition* (1), the signal of
621 one interlocutor was modulated by the constantly changing signal of their partner. Our
622 analysis of segment onset synchronization showed that, across conditions, when
623 participants were aligning their speech with the recording, the recorded model led,
624 while the experimental participants followed (with a larger onset difference in the
625 condition, in which the recording was made during solo reading). The fact that
626 experimental participants failed to take a leading position in dyadic reading with a
627 recorded model suggests that they find it challenging to use expectancy cues, predict
628 segment onsets, and produce an upcoming segment before it has been initiated by the
629 model. To a large extent, in the conditions with a recorded model, participants relied on
630 shortening an existing lag rather than predicting segment onsets based on expectancy
631 cues and initiating corresponding segments before they actually occurred in the model
632 speech. This result pattern indicates that the optimal use of expectancy cues occurs in a
633 dynamically changing situation (*live* reading), when the cues are modified online, in the
634 course of a reading-aloud task performed by both partners, allowing them to switch
635 between leading and lagging positions interactively. We propose that during live co-
636 presence, participants modify their prosodic timing patterns online and that this mutual
637 effort exerts beneficial effects on task performance. This means convergence is an

638 interactive process, and rhythmic synchronization with a live partner is more precise
639 than unidirectional synchronization, e.g., with a recording of someone's speech.

640 It is interesting to note that the effect found by Cummins (2002, 2009) – that
641 people synch better with recordings made in a synchronous reading condition versus a
642 non-synchronous reading condition – was confirmed in our study only for the narrative
643 text; we did not observe this effect with the poetic text. Meter potentially provides
644 additional expectancy cues for synchronization. Indeed, these cues proved relevant in
645 the most challenging condition (synchronization to a recording made during solo
646 reading) because they allowed participants to predict when the next **stressed** syllable
647 would occur. Thus, meter allows for synchronization to the onset of **stressed** syllables.
648 As attention shifts to the onsets of the **stressed** syllables, convergence of vowel onsets
649 in inter-stress intervals is relaxed. Narrative text does not provide strong expectancy
650 cues for stress onsets, thus reducing synchronization in the *Synchrony with recording*
651 *from non-live condition* for the narrative compared to the poetic text.

652 Our data emphasize the relevance of the co-presence of a live partner who
653 attempts to sync with the other speaker by dynamically changing their speech. We have
654 established that people synchronize better in the *Synchrony-live condition* (1) due to the
655 expectancy cues provided by co-presence. We posit that people achieved better
656 performance when synchronizing with another interlocutor because they were trying to
657 increase cooperation with the model speaker, a point also raised by Bowling et al.
658 (2013). In this sense, the aptitude for joint speech in humans could be linked to the
659 ability for entrainment found in primate communication. Notably, coordination in
660 primate vocal behavior appears to provide a signal of group cohesion and cooperative
661 strength (Ravignani, Bowling, & Fitch, 2014), as also seems to be the case with human
662 beings (Polyanskaya, Samuel, & Ordin, 2019). If we look at the different spheres where
663 joint speech is used – education, religion, sports, among others – we find a clear
664 indication that joint speech is used for group cohesion and cooperative strength.

665 The social and cognitive implications of humans' ability to entrain to each other'
666 behavior (including vocal communication) are enormous, in the sense of both individual
667 and social performance. In other words, there is great potential for exploiting this ability
668 in all spheres of social life and cognition since human beings engage in individual and
669 social entrainment tasks (entrainment of one's motor or vocal output to the signal
670 emitted by another individual) from birth. Furthermore, the effects of joint speech are
671 not only beneficial for achieving cooperation and cohesion. Recent research has found

672 that when two people are in a synchronous speech condition, neural networks that are
673 suppressed during individual speech are activated (Jasmin et al., 2017). For example,
674 activation in the auditory sensory cortices was usually enhanced during speech
675 perception but inhibited during speech production. And yet, during chorusing, these
676 sensory cortices were not suppressed. Moreover, activation of neural substrates in the
677 right hemisphere, which do not constitute part of the classic speech production network,
678 was observed. It seems that our brains detect that joint speech is, in fact, a different type
679 of speech, perhaps due to its highly cooperative and interactive nature. The authors
680 suggested that these neural changes in synchronous speech production overshadow the
681 distinction between one's own and partner's speech, such that the partner's speech is
682 perceived as if it were self-produced. This might facilitate extraction and use of
683 expectancy cues from the speech stream and allow participants to take leading positions
684 in dyadic reading tasks rather than relying solely on shortening the existing lag between
685 the segment they hear and the segment they initiate during their own speech production
686 (as the boundary between one's own speech production and the other's speech
687 production is diffused). Perceiving another person to be more like oneself enhances
688 group affiliation, facilitates joint activities, and encourages cooperation. Enhanced
689 cooperation during joint speech production and inevitable convergence of speech
690 patterns (and thus motor patterns, since speech is a motor activity) may be beneficial for
691 cognitive training, in the case of speech pathologies or any other type of condition
692 involving diminished speech and motor abilities, such as Parkinson's disease (Thaut et
693 al., 2001; deDreu et al., 2012).

694

Conclusions

695 Our data showed that the co-presence of a live speaker has the strongest modulatory
696 effect on rhythmic patterns and on the degree of synchronization between speakers in a
697 chorus reading task. We assume that this facilitatory effect emerges from dynamic
698 online changes resulting from two simultaneous drives: the drive to synchronize with a
699 partner; and the drive to provide cues that allow this partner to synchronize better with
700 oneself. The interaction of these two complementary drives – present in both partners –
701 increases the complexity of cue exchange and speech modification. The exact
702 mechanistic explanation for these interactions should be sought in the mathematics of
703 complexity theory.

704

References

- 705 Boersma, P., & Weenink, D. (2016). *Praat: Doing Phonetics by Computer* [Computer
706 software]. Version 6.0.19.
- 707 Bowling, D. L., Herbst, C. T., & Fitch, W. T. (2013). Social origins of rhythm?
708 Synchrony and temporal regularity in human vocalization. *PLoS One*, 8(11),
709 e80402. doi: 10.1371/journal.pone.0080402
- 710 Clopper, C. G., & Smiljanic, R. (2011). Effects of gender and regional dialect on
711 prosodic patterns in American English. *Journal of Phonetics*, 39(2), 237-245.
712 doi: 10.1016/j.wocn.2011.02.006
- 713 Coloma, G. (2016). Una versión alternativa de “El viento norte y el sol” en español.
714 *Revista de Investigación Lingüística*, 18, 191-212.
- 715 Crystal, D. (1969). *Prosodic systems and Intonation in English* (Vol. 1). Cambridge:
716 C.U.P.
- 717 Cummins, F., & Port, R. (1998). Rhythmic constraints on stress timing in English.
718 *Journal of Phonetics*, 26(2), 145-171. doi: 10.1006/jpho.1998.0070
- 719 Cummins, F. (2002). On synchronous speech. *Acoustic Research Letters Online*, 3(1),
720 7-11. doi: <https://doi.org/10.1121/1.1416672>
- 721 Cummins, F. (2003). Practice and performance in speech produced synchronously.
722 *Journal of Phonetics*, 31(2), 139-148. doi: [https://doi.org/10.1016/s0095-](https://doi.org/10.1016/s0095-4470(02)00082-7)
723 [4470\(02\)00082-7](https://doi.org/10.1016/s0095-4470(02)00082-7)
- 724 Cummins, F. (2007). *The quantitative estimation of asynchrony among concurrent*
725 *speakers*. Technical Report UCD-CSI-2007-2, School of Computer Science and
726 Informatics, University College Dublin.
- 727 Cummins, F. (2009). Rhythm as entrainment: The case of synchronous speech. *Journal*
728 *of Phonetics*, 37(1), 16-28. doi: <https://doi.org/10.1016/j.wocn.2008.08.003>
- 729 Cummins, F. (2015). Rhythm and Speech. In M. Redford (Ed.), *The Handbook of*
730 *Speech Production* (pp. 158-177). New York: Wiley.
- 731 deDreu, M. J., van der Wilk, A. S., Poppe, E., Kwakkel, G., and van Wegen, E. E.
732 (2012). Rehabilitation, exercise therapy and music in patients with Parkinson's
733 disease: a meta-analysis of the effects of music-based movement therapy on
734 walking ability, balance, and quality of life. *Parkinsonism Relat. Disord*, 18,
735 114–119. doi: 10.1016/S1353-8020(11)70036-0

- 736 Dellwo, V. (2006). Rhythm and speech rate: a variation coefficient for deltaC. In P.
737 Karnowski & I. Szigeti (Eds.), *Language and Language Processing* (pp. 231-
738 241). Frankfurt am Main: Peter Lang.
- 739 Eerola, T., Luck, G., & Toiviainen, P. (2006, August). An investigation of pre-
740 schoolers' corporeal synchronization with music. In *Proceedings of the 9th*
741 *International Conference on Music Perception and Cognition* (pp. 472-476).
742 The Society for Music Perception and Cognition and European Society for the
743 Cognitive Sciences of Music Bologna.
- 744 Fletcher, J. (2010). The prosody of speech: Timing and rhythm. In W. J. Hardcastle, J.
745 Laver, & F. E. Gibbon (Eds.), *The Handbook of Phonetic Sciences* (pp. 521-
746 602). John Wiley & Sons.
- 747 Fitch, W. (2013). Rhythmic cognition in humans and animals: distinguishing meter and
748 pulse perception. *Frontiers in Systems Neuroscience*, 7, 68. doi:
749 10.3389/fnsys.2013.00068
- 750 Goldman, J. (2011). EasyAlign: an automatic phonetic alignment tool under Praat. In
751 *Interspeech '11, Proceedings of the 12th Annual Conference of the International*
752 *Speech Communication Association* (pp. 3233-3236). Florence, Italy:
753 International Speech Communication Association.
- 754 Grabe, E., & Low, E. L. (2002). Durational variability in speech and the rhythm class
755 hypothesis. *Papers in Laboratory Phonology 7* (pp. 515-546). Berlin: Mouton de
756 Gruyter.
- 757 Hannon, E. E., & Johnson, S. P. (2005). Infants use meter to categorize rhythms and
758 melodies: Implications for musical structure learning. *Cognitive Psychology*,
759 50(4), 354-377. doi: 10.1016/j.cogpsych.2004.09.003
- 760 Jasmin, K., McGettigan, C., Agnew, Z., Lavan, N., Josephs, O., Cummins, F., Scott, S.
761 (2017). Cohesion and Joint Speech: Right Hemisphere Contributions to
762 Synchronized Vocal Production. *Journal of Neuroscience*, 36(17), 4669–4680.
763 doi: 10.1523/JNEUROSCI.4075-15.2016
- 764 Labov, W. (2006). *The social stratification of English in New York city* (3rd ed.).
765 Cambridge University Press.

- 766 Ladd, D. R., & Cutler, A. (1983). Models and measurements in the study of prosody. In
767 A. Cutler, & D. R. Ladd (Eds.), *Prosody: Models and measurements* (pp. 1-10).
768 Heidelberg: Springer.
- 769 Large, E., & Snyder, J. (2009). Pulse and meter as neural resonances. *Annals of the New*
770 *York Academy of Sciences*, 1169(1), 46–57. doi: 10.1111/j.1749-
771 6632.2009.04550.x
- 772 Lehiste, I. (1973). Rhythmic units and syntactic units in production and perception. *The*
773 *Journal of the Acoustical Society of America*, 54(5), 1228-1234. doi:
774 10.1121/1.1914379
- 775 Lehiste, I. (1976). Suprasegmental features of speech. In N. Lass (Ed.), *Contemporary*
776 *issues in Experimental Phonetics: Perspectives in Neurolinguistics and*
777 *Psycholinguistics* (pp. 225-242). New York: Academic Press.
- 778 Low E., Grabe E., & Nolan, F. (2000). Quantitative Characterizations of Speech
779 Rhythm: Syllable-Timing in Singapore English. *Language and Speech*, 43(4),
780 377-401. doi: 10.1177/00238309000430040301.
- 781 London, J. (2004). *Hearing in Time: Psychological Aspects of Musical Meter*. Oxford
782 University Press. Oxford, U.K.
- 783 Mairano, P. & Romano, A. (2010). *Correlatore*. [Computer software]. Retrieved May
784 22, 2016 from <http://www.lfsag.unito.it/correlatore/index.html>
- 785 Marcus, S. M. (1981). Acoustic determinants of perceptual center (P-center) location.
786 *Perception & Psychophysics*, 30(3), 247-256. doi: 10.3758/BF03214280
- 787 Merchant, H., Grahn, J., Trainor, L., Rohrmeier, M., & Fitch, W. T. (2015). Finding the
788 beat: a neural perspective across humans and non-human primates. *Phil. Trans.*
789 *R. Soc. B*, 370(1664), 20140093. doi: 10.1098/rstb.2014.0093
- 790 Moore, R. (2012). Finding Rhythm in Speech: A response to Cummins. *Empirical*
791 *Musicology Review*, 7(1-2), 36-44. doi: 10.18061/1811/52977
- 792 Morton, J., Marcus, S., & Frankish, C. (1976). Perceptual centers (P-centers).
793 *Psychological Review*, 83(5), 405–408. doi: 10.1037/0033-295x.83.5.405
- 794 Nolan, F., & Jeon, H.S. (2014). Speech rhythm: a metaphor? *Phil. Trans. R. Soc. B*,
795 369(1658), 20130396. doi: 10.1098/rstb.2013.0396

- 796 Nooteboom, S. (1997). The prosody of speech: melody and rhythm. In W.J. Hardcastle
797 & J. Laver (Eds.), *The Handbook of Phonetic Sciences* (pp. 640-673). Oxford:
798 Blackwell.
- 799 Pardo, J. S. (2010). Expressing oneself in conversational interaction. In E. Morsella
800 (Ed.), *Expressing oneself/expressing one's self: Communication, cognition,*
801 *language, and identity* (pp.183-196). Hove, England: Psychology Press/Taylor &
802 Francis. doi: 10.4324/9780203848708
- 803 Pardo, J. S., Urmanche, A., Wilman, S., & Wiener, J. (2017). Phonetic convergence
804 across multiple measures and model talkers. *Attention, Perception, &*
805 *Psychophysics*, 79(2), 637-659. doi: 10.3758/s13414-016-1226-0
- 806 Parra, V. (1988). La Muerte. In V. Parra, *Décimas. Autobiografía en verso*. Buenos
807 Aires: Editorial Sudamericana.
- 808 Pitt, M. A., & Samuel, A. G. (1990). The use of rhythm in attending to speech. *Journal*
809 *of Experimental Psychology: Human Perception and Performance*, 16(3), 564-
810 573. doi: 10.1037/0096-1523.16.3.564
- 811 Polyanskaya, L., Samuel, A.G., & Ordin, M. (2019). Speech rhythm convergence as a
812 social coalition signal. *Evolutionary Psychology* 17(3), 1-11. Doi:
813 10.1177/1474704919879335.
- 814 Quené, H. & Port, R. F. (2005) Effects of timing regularity and metrical expectancy on
815 spoken word perception. *Phonetica*, 62(1), 1-13. doi: 10.1159/000087222
- 816 Ramus, F., Nespors, M., & Mehler, J. (1999). Correlates of linguistic rhythm in the
817 speech signal. *Cognition*, 73(3), 265-292. doi: 10.1016/S0010-0277(99)00058-X
- 818 Ravnani, A., Bowling, D. & Fitch, T. (2014). Chorusing, synchrony and the
819 evolutionary functions of rhythm. *Frontiers in Psychology*, 5, 1118. doi:
820 10.3389/fpsyg.2014.01118
- 821 Repp, B. H. (2005). Sensorimotor synchronization: a review of the tapping literature.
822 *Psychonomic Bulletin & Review*, 12(6), 969-992. doi: 10.3758/BF03206433
- 823 Scott, S. (1998). The point of P-centres. *Psychological Research* 61, 4–11.
- 824 Smith, R., Rathcke, T., Cummins, F., Overy, K. & Scott, S. (2014). Communicative
825 rhythms in brain and behaviour. *Phil. Trans. R. Soc. B*, 369(1658), 20130389.
826 doi: 10.1098/rstb.2013.0389

- 827 Szczeppek Reed, B. (2010). Speech rhythm across turn transitions in cross-cultural talk-
828 in-interaction. *Journal of Pragmatics*, 42(4), 1037-1059. doi:
829 10.1016/j.pragma.2009.09.002
- 830 Thaut, M. H., McIntosh, K. W., McIntosh, G. C., & Hoemberg, V. (2001). Auditory
831 rhythmicity enhances movement and speech motor control in patients with
832 Parkinson's disease. *Functional Neurology*, 16(2), 163-172.
- 833 Turk, A., & Shattuck-Hufnagel, S. (2013). What is speech rhythm? A commentary on
834 Arvaniti and Rodriquez, Krivokapić, and Goswami and Leong. *Laboratory*
835 *Phonology*, 4(1), 93-118. doi:10.1515/lp-2013-0005
- 836 Vos, P. G., Mates, J., & van Kruysbergen, N. W. (1995). The perceptual centre of a
837 stimulus as the cue for synchronization to a metronome: Evidence from
838 asynchronies. *The Quarterly Journal of Experimental Psychology Section A*,
839 48(4), 1024-1040. doi: /10.1080/14640749508401427
- 840 Vuust, P., & Witek, M. A. (2014). Rhythmic complexity and predictive coding: a novel
841 approach to modeling rhythm and meter perception in music. *Frontiers in*
842 *Psychology*, 5, 1111. doi: 10.3389/fpsyg.2014.01111
- 843 Wells, J. (1996). *SAMPA for Spanish*. London: Division of Psychology and Language
844 Sciences, University College London.
- 845 White, L., and Mattys, S. (2007). Calibrating rhythm: First language and second
846 language studies. *Journal of Phonetics* 35(4), 501–522.
- 847

848 **FIGURE CAPTIONS**

849 *Figure 1.* Vowel % across conditions and texts. Here and in subsequent figures, error
850 bars show 2SE around the mean. Asterisks mark significant contrasts, * - $p < .05$, ** -
851 $p < .005$, *** - $p < .001$.

852 *Figure 2.* VarcoV across conditions and texts

853 *Figure 3.* VnPVI across conditions and texts

854 *Figure 4.* VarcoC across conditions and texts

855 *Figure 5.* CnPVI across conditions and texts

856 *Figure 6.* Mean differences for onsets across synchrony conditions and texts
857

858 Table 1 *Experimental materials*

Text number	Long name	Short name	Description	Objective
Text 0	Warm-up text	Warm-up text	16 sentences	Warm-up: 16 sentences were created to habituate experimental participants to the synchrony-reading task. These sentences also introduced specific vocabulary found in Texts 1 and 2, so as to prevent any semantic or phonological interference during the experiment.
Text 1	Narrative text with weak meter (Text 1)	Text with weak meter (Text 1)	A Spanish version of <i>The North Wind and the Sun</i> (Coloma, 2016)	Experimental: A phonetically balanced text, with weak meter, which includes every segmental sound in Spanish.
Text 2	Poetic text with strong meter (Text 2)	Text with strong meter (Text 2)	A poem written in octosyllabic verse, consisting of two stanzas (Parra, 1988).	Experimental: A text with strong meter.

859
860

861

Table 2 *Description of Conditions*

Condition	Name	Description	Synchronic condition
(1)	Synchrony-live	Experimental participants were recorded reading aloud and in synchrony with a female native speaker of Spanish (henceforth, the model speaker)	Yes
(2)	Synchrony with a recording from live condition	Experimental participants were recorded reading aloud and in synchrony with the recording of the model speaker obtained from the <i>Synchrony-live condition</i> (1)	Yes
(3)	Synchrony with a recording from non-live condition	Experimental participants were recorded reading aloud and in synchrony with a recording of the model speaker obtained from a non-synchronous condition	Yes
(4)	Solo recording condition	Experimental participants were individually recorded reading aloud	No

862

863

864
865Table 3 *Vocalic and consonantal interval measurements (none of the Condition*Text interactions were significant, and thus interactions are not reported).*

Measurement	Factor	Intra-subject effects	Pairwise Contrasts (²)
V%	Conditions (4)	F(2,375, 64,116) = 14,556; p < 0.0005; $\eta_p^2 = 0.350$ (¹)	Conditions 1 & 2 (p = .027); Conditions 1 & 3 (p < 0.0005); Conditions 1 & 4 (p = 0.001); Conditions 2 & 3 (p = .001)
	Text (2)	F(1, 27) = 281.132; p < 0.0005; $\eta_p^2 = 0.912$	Text 1 and Text 2 (p < 0.0005).
VarcoV	Conditions (4)	F(3, 75) = 18.076; p < 0.0005; $\eta_p^2 = 0.420$	Conditions 1 & 3 (p < 0.0005); Conditions 2 & 3 (p < 0.0005); Conditions 3 & 4 (p = .002)
	Text (2)	F(1, 25) = 84.173; p < 0.0005; $\eta_p^2 = 0.771$	Text 1 and Text 2 (p < 0.0005).
VnPVI	Conditions (4)	F(3, 66) = 4.355; p = 0.007; $\eta_p^2 = 0.165$	Conditions 3 & 4 (p = 0.004)
	Text (2)	F(1, 22) = 176.295; p < 0.0005; $\eta_p^2 = 0.889$	Text 1 and Text 2 (p < 0.0005).
VarcoC	Conditions (4)	No significant differences	No significant differences
	Text (2)	F(1, 26) = 23.102;	Text 1 and Text 2 (p < 0.0005).

$p < 0.0005$;

$\eta_p^2 = 0.470$

CnPVI	Conditions (4)	No significant differences	No significant differences
	Text (2)	F(1, 26) = .792; p = 0.025; $\eta_p^2 = 0.030$.	Text 1 and Text 2 (p = 0.025)

866 *Notes.* ¹Greenhouse-Geisser correction was applied to all repeated-measures ANOVAs in this study if the assumption
867 of sphericity was violated; ²Bonferroni correction was applied to all pairwise contrasts in the study.

868

869

Table 4 *Measurements for Speech Rhythm Synchrony*

Measurement	Factor	Intra-subject effects	Pairwise contrasts
Onset synchrony	Conditions (3)	F(1.415, 28.302) = 21.933; p < 0.0005; $\eta_p^2 = 0.523$	Conditions 1 & 2 (p < 0.0005); Conditions 1 & 3 (p < 0.0005); Conditions 2 and 3 (p = .025)
	Text (2)	No significant differences	No significant differences

870
871