

Regularity in speech rhythm as a social coalition signal

Leona Polyanskaya,¹ Arthur G. Samuel^{1,2,3} and Mikhail Ordin^{1,2}

¹ BCBL – Basque Centre on Cognition, Brain and Language, Donostia, Spain

² IKERBASQUE – Basque Foundation for Science, Bilbao, Spain

³ Department of Psychology, Stony Brook University, Stony Brook, NY

Address for correspondence: Leona Polyanskaya, BCBL – Basque Centre on Cognition, Brain and Language, Mikeletegi 69, Donostia, Gipuzkoa, Spain. l.polyanskaya@bcbl.eu

Graphical abstract

Humans treat rhythmicity in speech as a signal of cooperation between talkers. This rhythm-to-cooperation mapping is a biological faculty supported by the universality of rhythm cognition and social interaction. The evolutionary origin of this faculty is possibly the need to transmit and perceive coalition information and to evaluate cooperativeness in social groups of human ancestors. A link between social bonding, isochronous rhythm in vocalization, and synching of motor and vocal rhythm between individuals is a possible precursor of the speech faculty in humans.

Short title: Rhythmicity Manifests Cooperation

Abstract

Regular rhythm facilitates audiomotor entrainment and synchronization in motor behavior and vocalizations between individuals. As rhythm entrainment between interacting agents is correlated with higher levels of cooperation and prosocial affiliative behavior, humans can potentially map regular speech rhythm onto higher cooperation and friendliness between interacting individuals. We tested this hypothesis at two rhythmic levels: pulse (recurrent acoustic events) and meter (hierarchical structuring of pulses based on their relative salience). We asked the listeners to make judgments of the hostile or collaborative attitude of two interacting agents who exhibit either regular or irregular pulse (Experiment 1) or meter (Experiment 2). The results confirmed a link between the perception of social affiliation and rhythmicity: Evenly distributed pulses (vowel onsets) and consistent grouping of pulses into recurrent hierarchical patterns are more likely to be perceived as cooperation signals. People are more sensitive to regularity at the level of pulse than at the level of meter, and they are more confident when they associate cooperation with isochrony in pulse. The evolutionary origin of this faculty is possibly the need to transmit and perceive coalition information in social groups of human ancestors. We discuss the implications of these findings for the emergence of speech in humans.

Keywords: cooperation signal; speech rhythm; speech evolution; isochrony

Introduction

Patterns of verbal behavior of interlocutors become more similar as communication progresses. This widely studied phenomenon is referred to as speech convergence, interspeaker synchronization, accommodation, spontaneous phonetic imitation, to name just a few of the most frequently used terms. Alignment in the course of communication is a general behavioral phenomenon¹⁻³. In human-to-human communication via the medium of speech, convergence happens at multiple linguistic levels simultaneously: phonetic, lexical, syntactic, etc. In this investigation, the focus will be on the perception of rhythmic interspeaker synchronization. Rhythmic synchronization, once established, is maintained because it facilitates speech processing through routinization (allowing interlocutors to develop and use routine expressions) and through the enhancement of monitoring speech processing⁴.

Rhythm entrainment happens at multiple timescales and includes entrainment of speech gestures, movements, vocalizations, etc⁵. This happens when the rhythmic output of the speaker becomes input for the listener's perception systems and allows the listener to predict speech events and maintain a stable representation of recurrent acoustic patterns, which are later reproduced, forming a loop between rhythm production and perception^{6,7}.

Rhythm affects motor coordination, including articulation, interpersonal coordination, and social behavior⁸⁻¹¹. Rhythmic coordination increases social bonding^{12,13}, group identity^{14,15}, and is positively correlated with mutual attractiveness of interacting agents³ and with communication success^{7,16,17}. A mismatch in rhythm of speakers is a correlate of communication failure or mutual dislike. The current project is designed to determine whether regularity of speech rhythm is used by a listener to make pragmatic inferences regarding interpersonal relations between interacting agents. We also assess the degree of awareness in making pragmatic inferences based on rhythmic properties of the vocal signal. This might give us insight into whether listeners are aware of their ability to assign mental states to the interacting agents based on vocal signals.

There are differing views regarding the concept of rhythm in speech research¹⁸⁻²⁰. In the current study, speech rhythm will be understood as a succession of durations that define acoustic events in time. This approach agrees with McAuley's definition of rhythm as the structure that determines signal organization in time for music²¹, and with Clopper and Smiljanic's definition of rhythm as prosodic timing controlling temporal organization of speech²². However, we acknowledge that the collective opinion of researchers on the phenomenon of rhythm will not converge to a single definition.

Operationally, we will separate pulse and meter²³. Pulse is the occurrence of salient acoustic events, and rhythmic differences in pulse are captured by the durational ratios between salient acoustic events²⁴⁻²⁷. For the auditory system, the most salient acoustic event is a sharp amplitude rise²⁸, and in continuous speech, these rises are caused by the vowel onsets. This has been shown for Spanish, French, English, Chinese, Finnish, and Dutch²⁹. Vowel onsets generate recurring physiological responses²⁸ at the frequency of the syllable rate³⁰. These responses are used to extract the syllable as a distinguishable constituent³¹.

There are cross-language differences in the temporal distribution of vowel onsets^{26,27,32}. These differences are primarily determined by language-specific phonotactic constraints. Specifically, in languages that allow complex consonantal clusters some utterances may exhibit higher irregularity in the distribution of the vowel onsets than other utterances. In languages that primarily allow simple syllables consisting of alternating single consonants and vowels (due to strict phonotactic constraints), vowel onsets are distributed more evenly

across all utterances. Other linguistic factors include phonetic manifestations of lexical and phrasal stress, stress-induced vowel lengthening, phonetic reduction in unstressed vowels, the presence of long and short vowels or geminates (double consonants), etc. These factors can also affect durational variability between vowel onsets, thus contributing to cross-language differences in pulse.

Although Dauer suggested that such rhythmic differences are merely automatic consequences of linguistic circumstances³³, further studies showed that even when linguistic factors are controlled, rhythmic differences between languages still persist^{26,34}. The temporal distribution of vowel onsets may vary within languages in order to enhance temporal prediction and coordinate joint action³⁵. A regular distribution of pulses should then facilitate temporal prediction of an upcoming event, preparation of the response, and enhance attention to the sensory input^{36,37}.

In speech production, rhythmic convergence can be achieved via temporal regularity of the vocalizations by adjustment of the phases of vowel onsets. Although speech rhythm is often characterized by irregularity¹⁹, through an interaction between pulse extraction and rhythmic social entrainment, the percept of isochrony is likely to emerge³⁸. A regular distribution of syllables allows easier adjustment of phases and periods of movements, as well as synchronization of motor outputs⁵. Speech production is a kind of motor activity, and the main articulatory movement is the mandibular oscillation that builds the syllable frame, or vocalization frame³⁹. Temporal regularity of mandibular cycles enables predicting the onset of the following syllable and facilitates joint actions like chorusing, moving together or adjusting motor output to the acoustic input and thus interindividual motor actions^{38,40–43}. The tendency towards isochrony at the syllabic timescale can be driven by the urge to coordinate vocalizations and movements between individuals^{38,40,44}. As rhythmic synchronization promotes prosocial behavior, social affiliation and cooperation, and isochrony facilitates synchronization, we hypothesize that a regular temporal distribution of pulses, which results in easier entrainment between communicating agents, will be perceived by a third-party observer as a signal of cooperation and social affiliation between interacting individuals.

We manipulated durational ratios of the vowels to induce the changes in pulse. The durational ratios produce the differences in syllabic rhythm between individual utterances. The stimuli were then paired to imitate the exchange of utterances between two interacting agents (see Methods for details). In Experiment 1, we tested the prediction that a third-party observer (listener) can make pragmatic inferences regarding the cooperation level between interlocutors based on the regularity of pulse distribution in their utterances.

Pulses can be grouped into hierarchical patterns based on their relative perceptual salience. Structuring the pulses into hierarchical patterns is referred to as meter⁴⁵. As every syllable in the speech stream represents a pulse, different distributions of stressed syllables (i.e., stronger pulses) may result in different ways to structure the pulses. In bisyllabic words, these metrical patterns can be iambic (weak–strong) or trochaic (strong–weak). In Experiment 2, we tested the prediction that the interaction between individuals producing a regular rhythm (i.e., the same metrical pattern) will be perceived as more cooperative than the interaction between individuals producing more varied metrical patterns. Thus, we implemented only trochaic patterns in one stimulus type, and 50% trochaic and 50% iambic patterns in another stimulus type. At the same time, we reduced the overall temporal variability of syllabic durations, thus eliminating the differences in pulse (this allowed us to focus exclusively on the meter effect).

Experiment 1

Method.

Participants. Twenty-five Spanish-Basque bilinguals (age range: 18–30) without speech or hearing problems were recruited. The participants were Spanish dominant, fully functional in both languages, living in a bilingual environment. They had Spanish-speaking parents and their age of acquisition of Basque as their second language was 2–3 years.

Material. To construct an artificial language, we used five consonants (/s,m,n,l,f/) and five vowels (/a,u,e,o,i/). Ten bisyllabic nonsense words (samu, nelo, noma, namo, fenu, fale, lufe, mesu, sofu, and sela) were constructed by concatenating these phonemes into simple CV (consonant-vowel) syllables. None of these words is a real word in either Spanish or Basque. The vowel /i/ was only used in ‘*filler*’ syllables (fi, si, mi, ni, li) that were interspersed with the nonsense words. In a stream like: **FIMESUMISELALISAMUMIFALESINELO...**, concatenation of fillers with statistical words resulted in transitional probabilities (TPs) between syllables within words equal to 1.0, and between syllables straddling the word boundaries equal to 0.2, thus providing a reliable statistical cue for detecting the edges of statistically defined words. Filler syllables were more frequent than the syllables that comprised the statistical words. Frequent syllables were meant to model function words (prepositions, articles), while less frequent syllables were meant to model content words.

We imposed a prosodic hierarchy on the syllable strings to encourage filtering the acoustic input through speech processing mechanisms as opposed to just low-level auditory mechanisms. As shown in Figure 1, the prosodic constituents we included were phonological words (PWs), phonological phrases (PPs), and intonational phrases (IPs). Each stimulus (6.6 s in duration) was a single intonational phrase consisting of four phonological phrases, and each phonological phrase was composed of either two or three phonological words.

A phonological word also referred to as a clitic word or a prosodic word, is the domain for assigning lexical stress, i.e., only one of the syllables that compose a PW is stressed. Unlike lexical words, PWs include clitics, prepositions, or articles, i.e., frequent and often one-syllable function words that make a single speech unit. For example, “in the house” could be a phonological word, versus the lexical word “house.” In the experimental material, PWs consisted of a frequent syllable followed by a bisyllabic statistical word, e.g., **FIMESU** (see Fig. 1).

Stress was acoustically manifested by syllable lengthening of the word-initial syllables of the statistical words (durational parameters are given in Table 1). Thus, each PW was defined by transitional probabilities, the distribution of stressed syllables, and the distribution of frequent and infrequent syllables. A sequence of either two or three PWs produced as a single intonational unit made a phonological phrase (PP). The first PP in the sample stimulus in Figure 1 is **FIMESUMISELA**. The right edge of a PP was marked by lengthening the final vowel by 20% (to model the phrase-final lengthening phenomenon) and by pitch rising during the last vowel of the PP (high boundary tone). The left edge of a PP was marked by a higher pitch (160Hz) compared to the pitch on the preceding vowel. Thus, each PP was defined by the initial and final boundary tones and phrase-final lengthening. In addition, there was an overall declination contour going down to 120 Hz (for PPs with two PWs) or 110 Hz (for PPs with three PWs). A sequence of four PPs yielded an intonational phrase (IP). The right edge of the intonational phrase was marked by a falling pitch to 75 Hz (low boundary tone), replacing the high boundary tone of the last PP in each IP. IPs were defined by low (i.e., falling) boundary tones only. In general, we followed the prosodic hierarchy presented in Nespor and Vogel⁴⁶ to implement the prosodic structure. The constituents of lower hierarchical levels were embedded into the constituents of a higher hierarchical level,

with the edges of higher-level constituents fully aligned with the edges of lower-level constituents.

The acoustic stimuli were synthesized with the MBROLA speech synthesizer (ES2 voice). Syllable duration for stressed syllables was 280 ms, and syllable duration for unstressed syllables, including filler syllables, was 180 ms. Consonantal durations were calculated as the difference between syllable and vowel durations. We jittered the durations of vowels (for the sake of naturalness) and implemented the differences in vowel ratios between regular and irregular stimuli and between stressed and unstressed vowels by drawing the durational values with discrete steps from uniform distributions, varying according to rhythm type and syllable prominence (Table 1). These values reflected the fact that the proportion of vocalic to consonantal material is larger in utterances produced with regular rhythm in natural languages²⁷. One hundred twenty stimuli with regular rhythm and 120 stimuli with irregular rhythm were constructed.

In a pilot study, we tested whether the stimuli were perceived as language-like by naïve listeners. We played 30 IPs of the artificial language, 30 sentences in Estonian (a language not known by the listeners) resynthesized with the same method used to create the experimental stimuli, and 30 synthetic stimuli composed of random concatenations of the same syllables used in the experimental stimuli, with the same pitch parameters, but with randomly dispersed pitch accents and boundary tones. Twenty listeners listened to the 90 sound files (played in random order) and indicated for each whether they thought it was a sentence in a real language or not (on a 7-point scale, from 7—“yes, it is definitely a real language” to 1—“it is definitely not a language”). The “naturalness ratings” assigned to resynthesized Estonian sentences ($M = 5.6$) and to the stimuli used in the experiment ($M = 5.2$) were significantly higher ($P < 0.005$) than the rating given to the random concatenation of syllables ($M = 3.9$). The prosodic hierarchy and statistical regularities of the artificial language make the continuous acoustic stream segmentable into discrete embedded constituents (words, PWs, PPs, and IPs) and encourage filtering the acoustic input through speech processing mechanisms. Using a real language would limit options for manipulating the durations of vowels and rhythmic patterns. Given that the resynthesized sentences in a real language (Estonian) did not produce naturalness ratings significantly greater than the planned experimental stimuli ($P > 0.1$), the latter are well-suited to the purposes of the present study.

Procedure. We set up an AX discrimination experiment. Each trial included a pair of stimuli, separated by a 1-s pause. Each stimulus had either a regular or irregular rhythm, thus creating four types of pairs: regular–regular, irregular–irregular, regular–irregular, and irregular–regular. Each stimulus was used only once in one of the pairs, yielding 120 pairs of each type.

To understand whether regularity and rhythmic isochrony can be used by the listeners to make pragmatic inferences, we devised the following experimental cover story. Participants were told that they were going to hear two robots with artificial intelligence, talking to each other in an extraterrestrial language. The robots had been created by many different alien races sharing the common language of the galaxy. However, some races are currently at war, while others are currently in partnership. The aliens created their robots to speak the common galactic language, but the creators did not care to give them individual voices. The robots, which communicate on behalf of their inventors, convey either mutually hostile or mutually friendly messages. For each pair of robots, participants had to decide if their inventors are in a good relationship with each other, or if they are enemies: “Listen carefully and decide on the status of the diplomatic relationships: hostile to each other or collaborating with each other.”

This story was invented to be able to explain why all the utterances, produced by different interacting agents, nevertheless had the same voice. In designing the experiment, we did not want to reduce statistical power by introducing additional factors that might affect the judgments regarding the hostility or friendliness of the interacting agents (one voice, for example, might be potentially perceived as intrinsically more friendly than the other; there could be an interaction between the voice and the rhythm type and between the presence of rhythmic synchronization and the order of the voices in the stimuli pairs, etc.). Using the same voice forced the listeners to focus attention on the pronunciation rather than on paralinguistic factors like the spectral peculiarities of individual voices.

Task sequencing

First, the listeners were familiarized with the artificial language. We played them utterances (not used in the main experiment) with different rhythms in the artificial language. These utterances were continuously concatenated into a single 3-min acoustic stream. This is sufficient to initiate segmentation (extraction of the statistical words based on statistical and prosodic cues), thus engaging speech processing mechanisms. After familiarization, the AX discrimination experiment was launched. After hearing each pair of stimuli, the listeners had to press the button—“hostility” or “partnership”—corresponding to what they thought the relationships of the robots’ creators were. The participants had to listen to the whole trial to the end (the 6.6-s first stimulus followed by a 1-s pause and then the 6.6-s second stimulus) before they could register their response. We expected the listeners to think that the collaborative interacting agents would produce utterances characterized by a higher degree of syllabic isochrony, whereas mutually hostile agents would produce irregular utterances. Immediately after answering the first question, the listeners had to press the button “sure” or “not sure” to say whether they were certain in their evaluation of the diplomatic status.

A short training period with four stimulus pairs—not used in the experiment—preceded the main experiment to familiarize the participants with the interface and to make sure that the participants understood the instructions. The experiment was programmed in PsychoPy (v.1.83.04), the stimuli were presented via professional quality headphones, and the experiment was conducted in a soundproof cabin.

Results and discussion

One participant was excluded because he did not follow the instructions. The results are reported based on the responses of the remaining 24 participants. For each *pair type*, we performed one-sample *t*-tests comparing the percentage of “partnership” responses with the percentage of trials on which participants could have given this response by chance (50%). As shown in Table 2, participants responded at chance levels when stimuli with different syllabic rhythms were paired (regular–irregular and irregular–regular pair types), responded significantly above chance when the stimuli with regular rhythms were paired and below chance when the stimuli with irregular rhythm were paired.

A repeated-measures ANOVA was used to test the differences in the percentage of the trials among the different pair types on which participants responded “partnership.” The normality assumption for each reported test was verified by means of the Shapiro-Wilks test, and the Greenhouse-Geiser correction for *df* ($\epsilon = 0.635$) was applied to control for a sphericity violation (revealed by the significant Mauchly’s test, $P < 0.0005$). The ANOVA revealed a significant and substantial effect of *pair type*, $F(3, 69) = 20.57$, $P < 0.0005$, $\eta_p^2 = 0.472$. Pairwise comparisons showed that the percentage of “partnership” responses was significantly higher for regular–regular pairs than for irregular–irregular, regular–irregular, and irregular–regular pairs (at $P < 0.0005$ for each *pair type*, after Bonferroni correction). No difference was observed between regular–irregular and irregular–regular pairs ($P = 1.0$). For

irregular–irregular pairs, the percentage of “partnership” responses was lower than for the pairs pitting different rhythms (at $P < 0.05$, after Bonferroni correction). Taken together, this result pattern suggests that people tend to associate the irregular rhythm with hostility and regular rhythm with cooperative behavior. If both interacting agents produce speech with temporal regularity or irregularity, the listeners respond that the agents express correspondingly cooperative or hostile attitudes to each other. If the rhythm in utterances delivered by the interacting agents is different, the listeners are at a loss and respond at the chance level.

Importantly, there was considerable variation in the degree of irregularity for individual stimuli, esp. with irregular rhythm. As a result, some irregular stimuli may be perceived as less friendly than others, which might affect the overall percentage of “partnership” responses. To remove the potential effect of the variance in the item-based irregularity, we subjected the data to a mixed model binary logistic regression with *subject* and *stimulus* as random factors. For this analysis, we only included regular–regular and irregular–irregular trials, excluding pairs pitting utterances with different rhythms, which pose a dilemma for the listeners regarding the presence of cooperation between interacting agents. The remaining pair types were used as two levels of the fixed factor in the regression model. The corrected model was significant, $F(1, 1438) = 158.870$, $P < 0.0005$, classification accuracy 70%. The effect of the *pair type* (whether the response will be “partnership” or “hostile” depending on the rhythmic patterns in the auditory stimuli) was significant at $P < 0.0005$, $B = 1.43$ (SE = 0.114). Regular–regular pairs are 4.2 (95%CI: 3.35–5.23) times more likely to be perceived as signaling cooperation between interacting agents than irregular–irregular pairs (Fig. 2).

Experiment 2

The second experiment was designed to test the effect of regularity at the level of meter, when isochrony at the level of pulse is kept constant across stimuli.

Method

Participants. Twenty-seven Spanish-Basque native speakers with the same profile as in Experiment 1 were recruited.

Material. For Experiment 2, we used the 120 stimuli from Experiment 1 with regular rhythm. The second set of 120 (irregular) stimuli was created by making the location of the stress within words variable: half of the statistical words had stress on the word-initial syllable, and a half on the word-final syllable. Thus, regular stimuli in Experiment 2 exhibited stress on every third syllable, while in irregular stimuli the distribution of stressed syllables was random. Again, we paired regular and irregular stimuli in a 2 x 2 design.

Procedure: Identical to Experiment 1

Results and discussion

Two participants were excluded because they did not follow the instructions. The results are reported based on the responses of the remaining 25 participants. Table 3 shows the outcomes of one-sample *t*-tests, for each pair type, comparing the percentage of “partnership” responses with the percentage of trials on which participants could have given this response by chance (50%). The data show that participants were inclined to report that the interacting agents are cooperating when both agents produced utterances with regular meter. For the other pair types, the percentage of “partnership” responses did not differ significantly from what could be expected by chance.

A repeated-measures ANOVA was used to test for any difference in the percentage of the trials on which participants responded “partnership” among the four different pair types. The Greenhouse–Geiser correction ($\epsilon = 0.732$) was applied to control for a sphericity violation (revealed by a significant Mauchly’s test, $P = 0.014$).

The ANOVA revealed a significant effect of *pair type*, $F(3, 72) = 3.627$, $P = 0.03$, $\eta^2 = 0.131$. Further pairwise comparisons revealed that only the differences in the percentage of “partnership” responses for regular–regular versus regular–irregular ($P = 0.027$) and versus irregular–regular ($P = 0.017$) pairs were significant. The other comparisons did not exceed the significance threshold. This result suggests that listeners tend to say that interacting partners are cooperating when both elicit utterances with regular meter.

The stimuli with irregular, i.e., variable meter could contain a sequence of iambic or a sequence of trochaic meter, while other irregular stimuli could exhibit alternating iambic-trochaic metric patterns. Potentially, this could affect the perception of rhythmicity and thus of the cooperative urge between interacting agents. To control for the variability in the irregular stimuli, we selected only regular–regular and irregular–irregular pairs and subjected the responses to mixed models binary logistic regression with *subject* and *stimuli* as random effects, *pair type* as a fixed effect. The corrected model was significant, $F(1, 1498) = 14.211$, $P < 0.0005$, classification accuracy 62%. The effect of the *pair type* (whether the response will be “partnership” or “hostile” depending on the rhythmic patterns in the auditory stimuli) was significant at $P < 0.0005$, $B = 0.403$ ($SE = 0.107$). Figure 3 shows that participants tend to perceive the interacting agents as cooperating in all pair types in Experiment 2. However, in regular–regular pairs they are 1.5 times more likely to respond “partnership” than in irregular–irregular pairs (95%CI: 1.21–1.85).

It should be noted, however, that meter consistency could be perceived as a cooperation signal because trochaic metrical patterns are more typical in Spanish—one of the native languages of the participants. However, the consistency in metrical patterns is less rigid in Basque, the other native language of our participants. Should the listeners perceive utterances with a consistent meter as representative of one language and utterances with an inconsistent meter as representative of a different language, then their judgments could be driven by the perception of interacting agents speaking either the same (cooperating) or different (not cooperating) languages. If this were the case then we should have observed similar responses to irregular–irregular and regular–regular stimuli, which is contrary to the data. Thus, although we cannot rule out the possibility that linguistic experience influenced participants’ decisions in Experiment 2, we believe that this is unlikely. Further investigation is required to determine if there is a contribution of linguistic experience to perceptual judgments.

Analysis of confidence ratings across Experiments 1 and 2

To analyze listeners’ awareness of their ability to map rhythmicity in vocalization onto cooperation, we selected only regular–regular and irregular–irregular trials (eliminating potentially confusing cases when one agent produces regular utterances and the other agent irregular). We estimated awareness as the objective sensitivity of confidence ratings to discriminate between correct and incorrect responses: if listeners are able to judge how well they did in assigning mental states, they are aware of their ability and can use it to modify behavior consciously. An observer is not aware of the message in the signal if confidence ratings are not informative regarding the classification, even if the observer can classify the stimuli correctly above chance⁴⁷. Higher awareness is a sign of more efficient metacognitive processing of the signal. We used the Signal Detection Theory framework to quantify the metacognitive performance of participants. Table 4 presents how we defined hits and false

alarms (FA) in this study for perceptual (here, pragmatic inferences) and confidence decisions.

This conceptualization allows mapping objective sensitivity of a participant's confidence to mapping regular rhythm on cooperation. We used hierarchical Bayesian metacognition modeling to estimate metacognitive performance^{48,49}. The approach basically involves estimating sensitivity to the rhythm regularity available to make confidence judgments. This sensitivity measure is called meta-d'. The M-ratio (meta-d'/d') shows how much of the perceived signal is used to make confidence judgments, controlling for the individual level of perceptual performance and individual differences in bias (the tendency to assign overall higher or lower confidence to one's responses⁴⁹). The measures were estimated using the code by Fleming⁵⁰ available on (<https://github.com/metacoglab/Hmeta-d>).

We estimated d', meta-d' and M-ratios for each participant in both experiments (Figure 4). D' measures were compared with zero using one-sample t-tests. The scores were significantly above zero for Experiment 1, $t(24) = 4.424$. $P < 0.0005$, 95%CI for the difference is 0.914–1.413 and for Experiment 2, $t(26) = 2.642$. $P = 0.014$, 95%CI for the difference is 0.83–0.66. Afterward, we compared these measures between the experiments using 2-tailed t-tests (Table 5). The results show that participants indeed reliably map regular rhythm onto cooperation, and they are more sensitive to the coupling between cooperative behavior and regular rhythm at the level of pulse than meter.

Participants exhibit lower meta-sensitivity than what could be expected from their perceptual performance for ideal observers (77% in Experiment 1 and 65% in Experiment 2), which suggests that perceptual and confidence decisions may rely on different information sources. When participants are asked to report their confidence, they are specifically relying on information that is consciously processed, which can reduce access to an unconscious processing stream. Thus, unconscious processing contributes more to perceptual decisions than to confidence decisions^{51,52}. Metacognitive efficiency is also higher at the level of pulse. Metacognitive capacity of the participants in Experiment 2 is lower than in Experiment 1. These results mean that people are aware of rhythmic regularity and its relation to affiliative behavior, and the level of awareness of stronger for pulse isochrony than for meter consistency.

General discussion

Patterns of syllabic rhythm, based on the variability of syllable and vowel durations, depend on situational context and vary in the degree of isochrony not only between, but also within languages³⁵. Lower variance in durational ratios is found for utterances with a more regular pulse (vowel onsets). Our results show that a rhythmic tendency towards isochrony signals cooperative behavior between interacting agents to a third-party observer.

Vowel onsets—pulses—are used for imposing an isochronous temporal grid on the sequence of syllables. This temporal grid enables perception of pulses as isochronous⁵³. This psychological tendency is known as beat induction. Lower durational jitter in the temporal distribution of pulses facilitates beat induction by the cognitive system. In turn, listeners perceive an acoustic signal that facilitates beat induction as a marker of friendliness and cooperation. Rhythm that makes beat induction more challenging is perceived as a marker of hostility. If the signal emitted by one interacting agent facilitates beat induction, but the signal emitted by the other interactor obstructs it, then third-party listeners randomly (at chance level) respond whether the agents are cooperating or antagonistic to each other. This result suggests that the property of the signal per se, i.e., pulse isochrony in single utterances can be used to make pragmatic inferences.

Isochrony is processed as a signal of cooperation between interacting agents by a third-party observer who does not share the language with the interactors. Consequently, the signal of cooperation is not (only) in the message, it is also in the signal per se, in the way the signal is emitted. Thus, temporal regularity of the articulatory frames, i.e., syllable frames, is the message that is perceived even when referential communication is not possible. The ability to transmit the message in the absence of common referential semantics shared by the individuals can lead to the emergence of a communication system.

A prerequisite for the development of a referential and symbolic communication system is the receiver perceiving the signal as a communicative signal. The receiver should be aware that there is a message in the signal⁵⁴. The question "How did speech emerge?" can be rephrased as "How did a symbolic referential communicative system emerge from a noncommunicative system?" Rhythmic entrainment, vocal or nonvocal, facilitates cooperation and achieving common goals, and is established with more ease when individuals produce events at regular time intervals, allowing other individuals to synchronize their efforts. Higher temporal regularity thus allows easier synchronization and better synergy, signaling readiness of an individual for cooperation. Temporal regularity per se, i.e., isochrony, is a signal that signals a message and makes the first step in turning a noncommunicative behavior into a communicative one. The characteristic of the emission is passed on to the signal itself, and isochrony acquires referential meaning when it is used to transmit affection or discontent (in case the signal is irregular) for the communicative partners. Temporal regularity can signal invitation and agreement to cooperate, providing the context for the development of a referential communication system. Scott-Phillips, Kirby, and Ritchie⁵⁴ provided experimental evidence that the signal is able to signal its own signalhood only in the scenario of cooperation; Knight⁵⁵ has argued that speech could have emerged only as a cooperative signaling system. These papers suggest that manifesting and perceiving the cooperative urge is an important step in speech evolution.

The first step in speech emergence in phylogenesis could have been the ability to transmit meaning via the characteristics of the signal emission. The capacity to emit, perceive and process the emission of the signal is based on the three components of rhythmic cognition (1) evolutionary advances in development of isochrony detection, (2) generation of periodic motor actions, and (3) integration of sensory input and motor output—in other words, the entrainment component of rhythm cognition^{6,56}. Regular rhythm that facilitates rhythmic entrainment provides a possible foundation for the emergence of a vocal referential system—speech—in the human genus. Some phylogenetic evidence that stems from research with primates yields support for this idea. Pulse isochrony also signals cooperation in nonhuman animals, including avian species^{57,58} and some mammalian species (e.g., rodents, namely, ground squirrels, family *Sciuridae*⁵⁹). Male geladas, for example, produce lip-smacking expressions ("wobbles"), accompanied with vocalizations, during affiliative interactions with females⁶⁰. The degree of isochrony actually corresponds to the degree of affiliation and with the female responsiveness to the wobbles, and the frequency of vocalizations corresponds to the syllabic frequency of human speech (4–6 Hz). Ghazanfar and Takahashi⁶¹ reviewed a series of other studies with baboons, macaque monkeys and marmosets and concluded that perceptual processes are tuned to isochrony at the natural frequencies of communication signals in monkeys' vocalizations and in human speech. Interestingly, deviation from pulse isochrony signals aggression in communicative systems also in evolutionary distant (avian) species, e.g., in *Crex crex* (corncrakes). Information about aggressive motivation is transmitted by variable intervocalization intervals, without changing the spectral characteristics of vocalizations per se. Ręk and Osiejuk⁵⁷ suggested that the temporal organization of vocalizations (i.e., pulses) emitted by *Crex crex* "is an example of

the syntax, equivalent to a very simple Morse code system.” Deviations from simple isochrony lead to larger flexibility of signals and increase the complexity of communication signals in general, making it possible to convey a wider range of messages even by those species whose repertoire of vocalizations is small and genetically coded⁵⁸.

As discussed above, pulse isochrony allows better interpersonal motor and vocal coordination via motor coordination with the acoustic signal emitted by a different individual, which strengthens social bonding and promotes prosocial behavior. Importantly, humans possess a capacity which has so far not been attested in other species: the capacity to coordinate behavioral rhythms according to the acoustic rhythms at the metrical level, beyond the beat. In Experiment 2, we focused on whether meter regularity can be used as a timing device for the perception of interpersonal cooperation between interacting agents. The results showed that an acoustic signal that allows consistent grouping of pulses into stable patterns further enhances the perception of the interpersonal cohesion between the interacting agents. Importantly, in Experiment 2 participants tended to respond that the interacting agents were cooperating. We believe the overall bias to perceive the agents as cooperating is due to temporal isochrony at the pulse level, as it was implemented in regular stimuli in Experiment 1. Listeners already perceive this level of pulse isochrony as a cooperation signal, and meter irregularity cannot override this perception.

Analysis of awareness indicators show that participants are more aware that isochrony at the level of pulse signals cooperation and social bonding, although regularity at the level of meter can also be used as a marker of social bonding (but less efficiently). Awareness of one’s abilities is reflected in a feeling of confidence⁴⁹. This feeling enables assigning more or less credit to the information source, and thus calibrating behavior accordingly, which could explain why the development of rhythmic cognition was boosted alongside the development of other cognitive abilities like intentionality(43ff)⁶². Intentionality and mind reading have important implications for the fitness of the organisms in evolution because they allow for 1) active manipulation of the behavior of the others and 2) adaptive behavioral adjustments by selective attention to the relevant properties of the communicative signal. The further interplay between the signaler and the perceiver lead to more complex communicative systems.

Although different components of rhythmic cognition—beat induction, detection of grouping patterns, generation of rhythmic output, and rhythmic entrainment—are inseparable functionally, their biological, neural basis and evolutionary histories are likely to be different³⁸. Thus, separate components of rhythmic cognition—periodic motor pattern generation, pulse extraction, beat entrainment—can be observed in a wide variety of species, yet these vary in the degree to which the separate components are developed^{63,64}. Generation of periodic motor patterns and synchronization of motor behaviors between individuals to signal alliance and cooperation has been observed in different mammalian species^{65,66}. However, it is not yet clear whether the emission of a periodic signal that facilitates entrainment and synchrony is an adaptive signal in nonhuman species⁶⁵, or is a byproduct of affiliative prosocial behavior that is perceived by the nonspecifics, but is not controlled by the signaling individual⁶⁶. In humans, on the other hand, synchrony has been claimed to be an adaptive signal of cooperation^{40,67,68}. All periodic patterns, i.e., isochronous sequences of events, are rhythmic, but the reverse is not always true²¹. It is possible that rhythmic patterns based on periodicity acquire the functionality of an adaptive signal, thus enabling active expression of affiliation and a cooperative urge via vocalization, providing a foundation for referential communication via vocalizations⁵³.

Acknowledgements

The authors acknowledge financial support from the Spanish Ministry of Economy and Competitiveness (MINECO), through the “Severo Ochoa” Programme for Centres/Units of Excellence in R&D (SEV-2015-0490) to the BCBL, from European Commission as Marie Skłodowska-Curie fellow DLV-792331 to LP, from Ministerio de Ciencia E Innovacion by Grant PSI2017-82563-P to AS and Grant RTI2018-098317-B-I00 to MO.

Competing interests

The authors declare no competing interests

Figure captions

Figure 1. The prosodic structure of the stimuli. This sample stimulus (one IP) consists of four PPs, the first and second PPs consist of two PWs, the third and fourth PPs are made up of three PWs. Each PW is a combination of a filler (frequent) syllable and a statistical word.

Figure 2. Number of “partnership” and “hostile” responses for regular–regular and irregular–irregular stimuli pairs in Experiment 1. Error bars show ± 2 SE around the mean. The ratio of the number of “partnership” to the number of “hostile” responses for each pair type represents the odds how more likely the particular type of pairs to receive “partnership” than “hostile” responses. The odds ratio (4.2) shows how more likely regular–regular pairs are to receive “partnership” responses than irregular–irregular pairs.

Figure 3. Number of “partnership” and “hostile” responses for regular–regular and irregular–irregular stimuli pairs in Experiment 2. Error bars show ± 2 SE around the mean.

Figure 4. D' , meta- d' and M-ratios for experiment 1 and 2, for the regular–regular and irregular–irregular trials only. Error bars show ± 2 SE around the mean.

References

1. Bargh, J. A., Chen, M. & Burrows, L. Automaticity of social behavior: Direct effects of trait construct and stereotype activation on action. *J. Pers. Soc. Psychol.* **71**, 230–244 (1996).
2. Chartrand, T. L. & Bargh, J. A. The chameleon effect: The perception–behavior link and social interaction. *J. Pers. Soc. Psychol.* **76**, 893–910 (1999).
3. Dijksterhuis, A. & Bargh, J. A. The perception-behavior expressway: Automatic effects of social perception on social behavior. *Adv. Exp. Soc. Psychol.* **33**, 1–40 (2001).
4. Pickering, M. J. & Garrod, S. Toward a mechanistic psychology of dialogue. *Behav. Brain Sci.* **27**, 169–90; discussion 190–226 (2004).
5. Large, E. W. On synchronizing movements to music. *Hum. Mov. Sci.* **19**, 527–566 (2000).
6. Phillips-Silver, J., Aktipis, C. A. & A. Bryant, G. The Ecology of Entrainment: Foundations of Coordinated Rhythmic Movement. *Music Percept.* **28**, 3–14 (2010).
7. Stephens, G. J., Silbert, L. J. & Hasson, U. Speaker-listener neural coupling underlies

- successful communication. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 14425–30 (2010).
8. Reddish, P., Fischer, R. & Bulbulia, J. Let's Dance Together: Synchrony, Shared Intentionality and Cooperation. *PLoS One* **8**, e71182 (2013).
 9. Repp, B. H. & Su, Y.-H. Sensorimotor synchronization: A review of recent research (2006–2012). *Psychon. Bull. Rev.* **20**, 403–452 (2013).
 10. Richardson, M. J., Marsh, K. L., Isenhower, R. W., Goodman, J. R. L. & Schmidt, R. C. Rocking together: Dynamics of intentional and unintentional interpersonal coordination. *Hum. Mov. Sci.* **26**, 867–891 (2007).
 11. von Zimmermann, J. & Richardson, D. C. Verbal Synchrony and Action Dynamics in Large Groups. *Front. Psychol.* **7**, 2034 (2016).
 12. Hove, M. J. & Risen, J. L. It's All in the Timing: Interpersonal Synchrony Increases Affiliation. *Soc. Cogn.* **27**, 949–960 (2009).
 13. Wiltermuth, S. S. & Heath, C. Synchrony and Cooperation. *Psychol. Sci.* **20**, 1–5 (2009).
 14. Brown, S. Evolutionary Models of Music: From Sexual Selection to Group Selection. in 231–281 (Springer, Boston, MA, 2000). doi:10.1007/978-1-4615-1221-9_9
 15. Reddish, P., Bulbulia, J. & Fischer, R. Does synchrony promote generalized prosociality? *Religion. Brain Behav.* **4**, 3–19 (2014).
 16. Auer, P., Couper-Kuhlen, E. & Müller, F. *Language in time : the rhythm and tempo of spoken interaction.* (Oxford University Press, 1999).
 17. Cowley, S. J. Conversational functions of rhythmical patterning: A behavioural perspective. *Lang. Commun.* **14**, 353–376 (1994).
 18. Goswami, U. & Leong, V. Speech rhythm and temporal structure: Converging perspectives? *Lab. Phonol.* **4**, 67–92 (2013).
 19. Nolan, F. & Jeon, H.-S. Speech rhythm: a metaphor? *Philos. Trans. R. Soc. B Biol. Sci.* **369**, 20130396–20130396 (2014).
 20. Ravnani, A. & Norton, P. Measuring rhythmic complexity: A primer to quantify and compare temporal structure in speech, movement, and animal vocalizations. *J. Lang. Evol.* **2**, 4–19 (2017).
 21. McAuley, J. D. Tempo and Rhythm. in 165–199 (2010). doi:10.1007/978-1-4419-6114-3_6
 22. Clopper, C. G. & Smiljanic, R. Regional variation in temporal organization in American English. *J. Phon.* **49**, 1–15 (2015).
 23. Large, E. W. & Snyder, J. S. Pulse and Meter as Neural Resonance. *Ann. N. Y. Acad. Sci.* **1169**, 46–57 (2009).
 24. Dellwo, V., Leemann, A. & Kolly, M.-J. Rhythmic variability between speakers: Articulatory, prosodic, and linguistic factors. *J. Acoust. Soc. Am.* **137**, 1513–1528 (2015).

25. Polyanskaya, L. & Ordin, M. Acquisition of speech rhythm in first language. *J. Acoust. Soc. Am.* **138**, EL199–EL204 (2015).
26. Prieto, P., Vanrell, M. del M., Astruc, L., Payne, E. & Post, B. Phonotactic and phrasal properties of speech rhythm. Evidence from Catalan, English, and Spanish. *Speech Commun.* **54**, 681–702 (2012).
27. Ramus, F., Nespor, M. & Mehler, J. Correlates of linguistic rhythm in the speech signal. *Cognition* **73**, 265–92 (1999).
28. Greenberg, S. & Ainsworth, W. A. Speech Processing in the Auditory System: An Overview. in *Speech Processing in the Auditory System* 1–62 (Springer-Verlag). doi:10.1007/0-387-21575-1_1
29. Leong, V. & Goswami, U. Acoustic-Emergent Phonology in the Amplitude Envelope of Child-Directed Speech. *PLoS One* **10**, (2015).
30. Ghitza, O., Giraud, A.-L. & Poeppel, D. Neuronal oscillations and speech perception: critical-band temporal envelopes are the essence. *Front. Hum. Neurosci.* **6**, 340 (2012).
31. Ding, N., Melloni, L., Zhang, H., Tian, X. & Poeppel, D. Cortical tracking of hierarchical linguistic structures in connected speech. *Nat. Neurosci.* **19**, 158–164 (2016).
32. White, L., Mattys, S. L. & Wiget, L. Language categorization by adults is based on sensitivity to durational cues, not rhythm class. *J. Mem. Lang.* **66**, 665–679 (2012).
33. Dauer, R. Stress-timing and syllable-timing reanalyzed. *J. Phon.* **11**, 51–62 (1983).
34. Payne, E., Post, B., Astruc, L., Prieto, P. & Vanrell, M. del M. Measuring Child Rhythm. *Lang. Speech* **55**, 203–229 (2012).
35. Hawkins, S. Situational influences on rhythmicity in speech, music, and their interaction. *Philos. Trans. R. Soc. B Biol. Sci.* **369**, 20130398–20130398 (2014).
36. McAuley, J. D. & Fromboluti, E. K. Attentional entrainment and perceived event duration. *Philos. Trans. R. Soc. B Biol. Sci.* **369**, 20130401–20130401 (2014).
37. Ellis, R. J. & Jones, M. R. Rhythmic context modulates foreperiod effects. *Atten. Percept. Psychophys.* **72**, 2274–2288 (2010).
38. Fitch, W. T. The biology and evolution of rhythm: unravelling a paradox. in *Language and Music as Cognitive Systems* 73–95 (Oxford University Press, 2011). doi:10.1093/acprof:oso/9780199553426.003.0009
39. MacNeilage, P. F. The frame/content theory of evolution of speech production. *Behav. Brain Sci.* **21**, 499–511 (1998).
40. Bowling, D. L., Herbst, C. T. & Fitch, W. T. Social Origins of Rhythm? Synchrony and Temporal Regularity in Human Vocalization. *PLoS One* **8**, e80402 (2013).
41. Cross, I. Music, Cognition, Culture, and Evolution. in *The Cognitive Neuroscience of Music* 42–56 (Oxford University Press, 2003). doi:10.1093/acprof:oso/9780198525202.003.0004

42. Dissanayake, E. Antecedents of the temporal arts in early mother-infant interaction. in *The origins of music* (ed. N. L. Wallin, B. Merker, & S. B.) 389–410 (MIT Press, 2000).
43. Merker, B. H., Madison, G. S. & Eckerdal, P. On the role and origin of isochrony in human rhythmic entrainment. *Cortex* **45**, 4–17 (2009).
44. Tomasello, M., Carpenter, M., Call, J., Behne, T. & Moll, H. Understanding and sharing intentions: The origins of cultural cognition. *Behav. Brain Sci.* **28**, 675–691 (2005).
45. London, J. *Hearing in time : psychological aspects of musical meter*. (Oxford University Press, 2004).
46. Nespor, M. & Vogel, I. *Prosodic Phonology*. (DE GRUYTER, 2007). doi:10.1515/9783110977790
47. Persaud, N., McLeod, P. & Cowey, A. Post-decision wagering objectively measures awareness. *Nat. Neurosci.* **10**, 257–261 (2007).
48. Maniscalco, B. & Lau, H. A signal detection theoretic approach for estimating metacognitive sensitivity from confidence ratings. *Conscious. Cogn.* **21**, 422–430 (2012).
49. Fleming, S. M. & Lau, H. C. How to measure metacognition. *Front. Hum. Neurosci.* **8**, 443 (2014).
50. Fleming, S. M. HMeta-d: hierarchical Bayesian estimation of metacognitive efficiency from confidence ratings. *Neurosci. Conscious.* **2017**, (2017).
51. Pasquali, A., Timmermans, B. & Cleeremans, A. Know thyself: Metacognitive networks and measures of consciousness. *Cognition* **117**, 182–190 (2010).
52. Del Cul, A., Dehaene, S., Reyes, P., Bravo, E. & Slachevsky, A. Causal role of prefrontal cortex in the threshold for access to consciousness. *Brain* **132**, 2531–2540 (2009).
53. Ravnani, A. & Madison, G. The Paradox of Isochrony in the Evolution of Human Rhythm. *Front. Psychol.* **8**, 1820 (2017).
54. Scott-Phillips, T. C., Kirby, S. & Ritchie, G. R. S. Signalling signalhood and the emergence of communication. *Cognition* **113**, 226–233 (2009).
55. Knight, C., Studdert-Kennedy, M. & Hurford, J. R. *The Evolutionary emergence of language : social function and the origins of linguistic form*. (Cambridge University Press, 2000).
56. Todd, N. P. M., Lee, C. S. & O’Boyle, D. J. A sensorimotor theory of temporal tracking and beat induction. *Psychol. Res.* **66**, 26–39 (2002).
57. Ręk, P. & Osiejuk, T. S. Sophistication and simplicity: conventional communication in a rudimentary system. *Behav. Ecol.* **21**, 1203–1210 (2010).
58. Ręk, P. & Osiejuk, T. S. Temporal patterns of broadcast calls in the corncrake encode information arbitrarily. *Behav. Ecol.* **24**, 547–552 (2013).
59. Owings, D. H., Gladney, A. B., Hennessey, D. F. & Leger, D. W. Different Functions of

- "Alarm" Calling for Different Time Scales: a Preliminary Report On Ground Squirrels. *Behaviour* **99**, 101–116 (1986).
60. Bergman, T. J. Speech-like vocalized lip-smacking in geladas. *Curr. Biol.* **23**, R268–R269 (2013).
 61. Ghazanfar, A. A. & Takahashi, D. Y. The evolution of speech: vision, rhythm, cooperation. *Trends Cogn. Sci.* **18**, 543–53 (2014).
 62. Dunbar, R. I. M. (Robin I. M. *The human story : a new history of mankind's evolution.* (Faber and Faber, 2004).
 63. Patel, A. D., Iversen, J. R., Bregman, M. R. & Schulz, I. Experimental Evidence for Synchronization to a Musical Beat in a Nonhuman Animal. *Curr. Biol.* **19**, 827–830 (2009).
 64. Schachner, A., Brady, T. F., Pepperberg, I. M. & Hauser, M. D. Spontaneous Motor Entrainment to Music in Multiple Vocal Mimicking Species. *Curr. Biol.* **19**, 831–836 (2009).
 65. Aureli, F., Preston, S. D. & de Waal, F. B. Heart rate responses to social interactions in free-moving rhesus macaques (*Macaca mulatta*): a pilot study. *J. Comp. Psychol.* **113**, 59–65 (1999).
 66. Connor, R. C., Smolker, R. & Bejder, L. Synchrony, social behaviour and alliance affiliation in Indian Ocean bottlenose dolphins, *Tursiops aduncus*. *Anim. Behav.* **72**, 1371–1378 (2006).
 67. Hagen, E. H. & Bryant, G. A. Music and dance as a coalition signaling system. *Hum. Nat.* **14**, 21–51 (2003).
 68. McNeill, W. H. *Keeping together in time : dance and drill in human history.* (Harvard University Press, 1995).

FIGURE LEGENDS

FIGURE 1. Prosodic structure of the stimuli. This sample stimulus (one IP) consists of four PPs, the first and second PPs consist of two PWs, the third and fourth PPs are made up of three PWs. Each PW is a combination of a filler (frequent) syllable and a statistical word.

FIGURE 2. Number of “partnership” and “hostile” responses for *regular-regular* and *irregular-irregular* stimuli pairs in Experiment 1. Error bars show $\pm 2SE$ around the mean. The ratio of the number of “partnership” to the number of “hostile” responses for each pairtype represents the odds how more likely the particular type of pairs to receive “partnership” than “hostile” responses. The odds ratio (4.2) shows how more likely *regular-regular* pairs are to receive “partnership” responses than *irregular-irregular* pairs.

FIGURE 3. Number of “partnership” and “hostile” responses for *regular-regular* and *irregular-irregular* stimuli pairs in Experiment 2. Error bars show $\pm 2SE$ around the mean.

FIGURE 4. D' , meta- d' and M-ratios for experiment 1 and 2, for the *regular-regular* and *irregular-irregular* trials only. Error bars show $\pm 2SE$ around the mean.







