



The role of native language and fundamental design of the auditory system in detecting rhythm changes

Journal:	<i>Journal of Speech, Language, and Hearing Research</i>
Manuscript ID	JSLHR-S-18-0299.R1
Manuscript Type:	Research Article
Date Submitted by the Author:	15-Nov-2018
Complete List of Authors:	Ordin, Mikhail; Basque Center on Cognition Brain and Language, Spoken Language; Ikerbasque, Polyanskaya, Leona; Universita degli Studi di Padova, Department of Linguistic and Literary Studies Gomez, David; Universidad de O'Higgins, Institute of Educational Sciences Samuel, Arthur; Stony Brook University, Psychology; Basque Center on Cognition Brain and Language, n/a; Ikerbasque, n/a
Keywords:	Prosody, Speech perception, Central auditory processing, Cognition, Psychoacoustics

SCHOLARONE™
Manuscripts

1 The role of native language and the fundamental design of the auditory system in detecting rhythm
2
3 changes
4

5 RUNNING HEAD: L1 and cognitive mechanisms in rhythm perception
6
7
8
9

10 Mikhail Ordin^{1,2}, Leona Polyanskaya¹, David Maximiliano Gómez^{3,4}, Arthur G. Samuel^{1,2,5}
11

12 ¹Basque Centre for Brain, Language, and Cognition, Paseo Mikeletegi 69, Donostia, 2009, Spain
13

14 ²Ikerbasque, Basque Foundation for Science, Maria Diaz de Haro 3, Bilbao, 48013, Spain
15 m.ordin@bcbl.eu

16 ³Institute of Educational Sciences, Universidad de O'Higgins, Avenida Libertador Bernardo O'Higgins 611,
17 Rancagua, Chile

18 ⁴Center for Advanced Research in Education, Universidad de Chile, Periodista Jose Carrasco Tapia 75,
19 Santiago, Chile

20 ⁵Department of Psychology, Stony Brook University, NY, USA
21
22
23
24

25 Corresponding Author:

26
27 Mikhail Ordin

28 Basque Centre for Brain, Language, and Cognition,

29 Paseo Mikeletegi 69,

30 San Sebastian, 20009, Spain

31 Tel: +34 943 309 300 F: +34 943 309 300

32 Email: mikhail.ordin@gmail.com, m.ordin@bcbl.eu
33
34
35
36

37 Conflict of Interest: The authors report no relevant conflicts of interest related to this manuscript.
38

39 Funding: The authors acknowledge support from the Spanish Ministry of Economy and Competitiveness
40 (MINECO) Grant # PSI2017-82563-P (to AGS), from the 'Severo Ochoa' Programme for Centres/Units of
41 Excellence in R&D (SEV-2015-490), and from the Basque Foundation for Science (IKERBASQUE). DMG was
42 supported by Grant PIA/Basal FB0003 from the Chilean Research Council (CONICYT). LP was supported by
43 the Spanish Ministry of Economy and Competitiveness (MINECO) via Juan de la Cierva fellowship.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Abstract

Purpose: We investigated whether rhythm discrimination is mainly driven by the native language of the listener or by the fundamental design of the human auditory system and universal cognitive mechanisms shared by all people irrespective of rhythmic patterns in their native language.

Method:

In multiple experiments, we asked participants to listen to two continuous acoustic sequences and to determine whether their rhythms were the same or different (AX discrimination). Participants were native speakers of four languages with different rhythmic properties (Spanish, French, English, German), to understand whether the predominant rhythmic patterns of a native language affect sensitivity, bias and reaction time (RT) in detecting rhythmic changes in linguistic (Experiment 2) and in non-linguistic (Experiments 1 and 2) acoustic sequences. We examined sensitivity and bias measures, as well as RTs. We also computed Bayes factors in order to assess the effect of native language.

Results: All listeners performed better (i.e., responded faster and manifested higher sensitivity and accuracy) when detecting the presence or absence of a rhythm change when the first stimulus in an AX test pair exhibited regular rhythm (i.e., a syllable-timed rhythmic pattern) than when the first stimulus exhibited irregular rhythm (i.e., stress-timed). This result pattern was observed both on linguistic and non-linguistic stimuli and was not modulated by the native language of the participant.

Conclusion: We conclude that rhythm change detection is a fundamental function of a processing system that relies on general auditory mechanisms and is not modulated by linguistic experience.

Keywords: rhythm perception, rhythm discrimination, rhythm processing, speech rhythm, linguistic experience

1 INTRODUCTION

Rhythm perception in general and discrimination of rhythmic patterns in particular are essential skills for speech and language processing and for language acquisition in infancy (Langus, Mehler & Nespors, 2018). Rhythmic patterns differ between languages (Gervain et al., 2008; Grabe & Low, 2002; Payne et al., 2012; Ramus & Mehler, 1999; White & Mattys, 2007) and non-native rhythm is a salient aspect of L2 (second language) speech (van Maastricht et al, in press; Ordin & Polyanskaya, 2015; White & Mattys, 2007). Non-native (Polyanskaya, Ordin, & Busa, 2017; Tajima, Port, & Dalby, 1997) or pathological (Kent et al., 1989) rhythmic patterns affect speech accentedness and comprehensibility by disrupting inter-speaker entrainment via speech rhythm (Borrie & Liss, 2014; Peelle, Gross & Davis, 2013). These observations suggest that rhythmic patterns in speech might be processed via the phonological filter of the native language. Alternatively, rhythmic perception could rely on a fundamental neurophysiological mechanism (Gitza, 2011; Greenberg & Ainsworth, 2004; Hickok et al., 2015; Howard & Poeppel, 2012) that is shared by all humans irrespective of their native language. In fact, this type of basic neurophysiological mechanism could underlie rhythm discrimination by animals (Tincoff et al., 2005; Toro, Trobalon, & Sebastian-Galle, 2003) and pre-linguistic babies (Nazzi & Ramus, 2003; Ramus, Nespors & Mehler, 1999) as well.

The existing literature is in fact consistent with two plausible and reasonable hypotheses: Either (a) linguistic experience (primarily, one's native language) shapes rhythm processing, or (b) prosody in general (and rhythmic structures in particular) in natural languages is shaped by the general design of the auditory system, cognitive mechanisms, and neural physiology. The objective of this study is to pit these two hypotheses, both logically coherent and plausible according to prior empirical evidence, against one another. The importance of addressing this question lies in the fact that deficits in rhythm perception are related to dyslexia (Molinaro et al., 2016; Muneaux et al., 2004); various speech disorders are also linked to rhythm abnormalities (Liss, White, Mattys, et al., 2009). It would be beneficial to know to what degree the clinical solutions are language-independent, versus linked to a particular native language of the affected person.

1.1 SPEECH RHYTHM

1 The word *rhythm* implies some degree of periodicity and isochrony (as in music, where certain patterns re-
2 occur at regular intervals). Considering these properties led James (1940), Pike (1945), Abercrombie (1967)
3 and Ladefoged (1993) to divide languages into stress-timed (in which intervals between stressed fragments
4 are of roughly equal duration, e.g., German, English, Dutch, Russian), syllable-timed (in which syllables are
5 perceived as equally long, e.g., French, Italian, Spanish) and mora-timed (with equal moras, e.g., Japanese,
6 Finnish, West Greenlandic, and Austronesia languages such as Gilbertese and Hawaiian). This distinction
7 was initially based only on auditory impressionistic analysis, and acoustic measurements in later studies
8 failed to fully support this claim (Dauer, 1983; Pamies Bertran, 1999; Roach, 1982).
9

10
11
12
13
14
15
16
17
18
19
20 However, despite the elusive nature of acoustic cues underlying rhythmic differences in speech,
21 research has shown that rhythmic differences between languages are indeed perceived (e.g., Gervain et
22 al., 2008; Ramus & Mehler, 1999). Recent psycholinguistic data show that discrimination of language-
23 specific rhythmic patterns is based on timing cues that influence the durational variability of speech
24 constituents: segmental length contrasts, compensatory lengthening and shortening, final lengthening or
25 initial strengthening at the edges of prosodic units (White et al., 2012). Also, languages may differ in regard
26 to the presence or absence of long-short vowel contrasts, vowel harmony, flexibility of stress placement,
27 and the relative contribution of duration and pitch to linguistic prominence. In addition, the degree of
28 vowel reduction that occurs in unstressed syllables, stress-induced lengthening of vowels, and phonotactic
29 constraints that allow longer and shorter consonantal clusters may further enhance or inhibit variability of
30 speech intervals. Such phonetic properties lead to cross-linguistic differences in durational variability of
31 vocalic intervals (V), consonantal clusters (C) and syllables (S), as well as differences in the proportion of
32 vocalic and consonantal material in speech. Languages with more properties that enhance durational
33 variability are those that are traditionally labeled as stress-timed on the basis of their auditory impression
34 (Dauer, 1983; Schiering, 2007). This understanding is based on the idea that rhythm is a serial arrangement
35 of time intervals, and timing relations define the organization and temporal structure of the auditory
36 scenes around the listener at different timescales. The differences in timing organization are captured by
37 various rhythm metrics. In Table 1, we summarize the most widely used metrics that have been proposed
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 to capture the patterns of durational variability. Higher %V and lower values of the other measures signal a
2
3 higher degree of syllable-timing.
4

5 The durational cues to perceived speech rhythm are not binary but continuous; the division of
6
7 languages into stress-versus syllable-timed is therefore not dichotomous. Rather, languages can be
8
9 positioned on a continuum between the opposing rhythm extremes. Moreover, estimating cross-linguistic
10
11 rhythmic differences requires calculating the durational variability of speech intervals within individual
12
13 utterances. Utterances even within the same language can also vary in the extent of stress- or syllable-
14
15 timing. Therefore, we will define rhythmic characteristics in terms of the degree of syllable- versus stress-
16
17 timing, referring to the *degree* of regularity (isochrony) in the duration of V, C and S intervals. The higher
18
19 the degree of durational variability of speech intervals, the more stress-timed the utterance is.
20
21
22
23

24 1.2 GENERAL MECHANISMS OF RHYTHM PERCEPTION

25
26 Abundant empirical evidence shows that adults, pre-linguistic babies and even animals are able to
27
28 differentiate between languages that are perceived to be rhythmically contrastive (e.g., Japanese or French
29
30 vs. Dutch or English), while they have difficulty discriminating rhythmically similar languages (e.g., Japanese
31
32 vs. French, or Dutch vs. English) (Nazzi, et al., 1998; Nazzi & Ramus, 2003; Ramus et al., 1999; Ramus &
33
34 Mehler, 1999; Ramus et al., 2000; Tincoff et al., 2005; Toro et al., 2003). This widespread ability suggests
35
36 that discrimination of rhythmic patterns could be governed by general properties of the auditory system
37
38 and cognitive mechanisms shared by all humans, irrespective of their native language.
39
40
41
42
43

44 In fact, several recent studies have shown the existence of a neural basis for the ability to detect
45
46 rhythmic changes (Hickok et al., 2015), and for the ability of acoustic rhythm to modulate the excitability of
47
48 the auditory system (Ghitza et al., 2013; Greenberg & Ainsworth, 2004; Hickok et al., 2015). These abilities
49
50 are presumably based on entrainment of neural oscillations to the acoustic rhythms (Giraus & Poeppel,
51
52 2012; Gitza, 2011; Ghitza et al., 2013; Howard & Poeppel, 2012). Regular rhythms couple with neural
53
54 oscillations and facilitate attention better than irregular rhythms (Barnes & Jones, 2000; Howard &
55
56 Poeppel, 2010; Jones, et al., 2002). Regularity also allows predicting the onset of the following vowels, and
57
58 a rhythm change is detected when this prediction is not met. Rhythms characterized by a higher degree of
59
60

1 stress-timing make it difficult to build expectations of when the next beat should happen, and this
2
3 uncertainty leads to poorer preparation and slower responding (Ellis & Jones, 2010). McAuley and
4
5 Fromboluti (2015) showed that variability in the timing of tones weakens the onset timing effect, leads to
6
7 distortions in perception of tone durations, and undermines attentional entrainment. They proposed that
8
9 the greater variability in timing impairs the entrainment of acoustic rhythms and neural oscillations and
10
11 disrupts perception of interval durations.
12
13

14
15 The existing evidence suggests that syllable-timed, regular rhythm should lead to a better coupling
16
17 between speech and attentional rhythms and thus facilitate attending compared to stress-timed, irregular
18
19 rhythm. If so, it might be easier to detect a change of rhythm when the listener is first tuned into an
20
21 acoustic stream with regular vowel onsets, which are salient auditory events in speech-like signals. This
22
23 faculty, being purely physiological in nature, should not be affected by linguistic experience and the range
24
25 of speech rhythms in the native language of a listener.
26
27
28

29 1.3 LANGUAGE-BASED HYPOTHESIS OF RHYTHMIC DISCRIMINATION

30
31
32 However, there is also some prior evidence in favor of the hypothesis that rhythmic discrimination can be
33
34 modulated by linguistic experience. Speakers of rhythmically contrastive languages employ different
35
36 weightings of speech constituents (morae, syllables, feet, inter-stress intervals) and durational cues when
37
38 segmenting continuous speech into discrete words and phrases (Cutler & Butterfield, 1992; Kim et al.,
39
40 2008; Murty et al., 2007; Polka & Sundara, 2011; Smith et al., 1989)¹. Erra and Gervain (2016) showed that
41
42 the statistical structure of a particular language is shaped, among other things, by rhythmic patterns,
43
44 which, in turn, can fine-tune the auditory codes to the statistics of the native language for efficiency of
45
46 neural coding of speech. Non-native rhythmic patterns in utterances produced by second language (L2)
47
48 speakers contribute to perceived accentedness and impair intelligibility of L2 speech (Polyanskaya et al.,
49
50 2017; Tajima et al., 1997). Rhythm plays a very important role in acquisition of linguistic features of
51
52 particular languages (Langus et al., 2018) and dialects (Clopper & Smiljanic, 2015; Polka & Sundara, 2011).
53
54 These studies confirm the psychological and linguistic reality of rhythm and the fundamental role of
55
56
57
58
59
60

¹ The consistent use of distinct segmentation units in different languages is challenged by a range of studies e.g., in Content et al. (2001).

1 rhythm in language acquisition and in speech processing, suggesting that processing of speech rhythm can
2
3 be language-based.
4

5 Native speakers of more stress-timed languages have an advantage of being familiar with a wider
6
7 range of linguistic rhythms because utterances in such languages vary in how stress-timed they are.
8
9 Children acquiring a stress-timed language start with syllable-timed rhythm (Bunta & Ingram, 2007; Ordin
10
11 & Polyanskaya, 2014; Payne et al., 2012; Polyanskaya & Ordin, 2015). Child-directed speech also exhibits a
12
13 higher degree of isochrony compared to adult-directed speech, and the stylistic differences are more
14
15 extreme in more stressed-timed languages. Thus, children, in the course of L1 acquisition, are exposed to
16
17 more syllable-timed utterances via child-directed speech and other children's speech. At the same time,
18
19 they get exposure to more stress-timed utterances via adult-directed speech, and their own rhythm
20
21 patterns also become increasingly more stress-timed as acquisition progresses (van Maastricht et al., in
22
23 press; Polyanskaya & Ordin 2015; Prieto et al., 2012). Importantly, languages that are generally more
24
25 stress-timed allow utterances with low durational variability of the speech constituents. In contrast,
26
27 languages that are generally more syllable-timed rarely include stress-timed utterances (Ordin &
28
29 Polyanskaya, 2015a). This is in part due to the strict phonotactic constraints which prohibit complex
30
31 consonantal clusters and lead to utterances consisting of predominantly CV syllables with more
32
33 isochronous speech intervals in syllable-timed languages. In stress-timed languages, with lax phonotactic
34
35 constraints, some utterances may be more syllable-timed, consisting of predominantly CV syllables and
36
37 monosyllabic words, and other utterances may contain words with complex consonantal clusters and
38
39 multiple cases of reduced vowels (Prieto et al., 2012). Consequently, native speakers of stress-timed
40
41 languages have experience with a wider range of possible rhythmic patterns, more and less isochronous,
42
43 including those occurring in more syllable-timed languages. A language-based hypothesis would predict
44
45 that native speakers of more stress-timed languages might benefit from experience with a wider spectrum
46
47 of rhythms and perform better when asked to discriminate between contrastive rhythms. Indeed, Lidji,
48
49 Palmer, Peretz & Morningstar (2011) showed that native speakers of English (a stress-timed language)
50
51 entrain their tapping performance to the rhythm in stress-timed utterances better than native French
52
53
54
55
56
57
58
59
60

1 speakers, while native French speakers did not reveal a better entrainment of their motor output to the
2 syllable-timed utterances compared to French.
3

4 1.4 PREDICTIONS

5
6
7 If perception of speech rhythm depends fundamentally on the aforementioned features of the neural and
8 auditory system, then listeners, irrespective of their native language, should perform better detecting the
9 switch from syllable-timed to stress-timed patterns. If instead (or in addition) linguistic experience matters,
10 then native speakers of less syllable-timed languages (e.g., German/English) should outperform the
11 speakers of more syllable-timed languages (e.g., Spanish/French) due to their exposure to a wider range of
12 rhythms in their native languages, irrespective of the direction of the rhythmic change.
13
14
15
16
17
18
19
20
21

22 To choose between these alternatives, we conducted two AX discrimination experiments. In these
23 experiments, a trial included two successive streams of syllables that could either match or mismatch. For
24 the first experiment, we recruited native speakers of French and German (French is more syllable-timed
25 and German is more stress-timed) and measured the effect of (a) their native language and (b) the rhythm
26 in the paired stimuli on decision reaction time and accuracy. For the second experiment, we recruited
27 native speakers of Spanish and English (Spanish is more syllable-timed and English is more stress-timed),
28 and modified the difficulty level by making the rhythmic differences between the stimuli more subtle. We
29 compared performance in rhythmic discrimination on linguistic and non-linguistic stimuli, providing a wider
30 range of languages and stimuli. In both experiments, the first member of an AX pair was either stress-
31 timed or syllable-timed. A significant effect of this factor (from syllable-timed or from stress-timed, to
32 either a rhythmically similar or a rhythmically contrastive stimulus) would provide evidence in favor of the
33 hypothesis that rhythmic discrimination is governed by the design of the auditory system. A significant
34 effect of native language (syllable-timed vs. stress-timed) would provide evidence that rhythm
35 discrimination is affected by linguistic experience.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55

56 2 EXPERIMENT I

57 2.1 METHOD

58 2.1.1 PARTICIPANTS

59
60

1 We recruited monolingual native speakers of two rhythmically different languages – French and German –
2 French being more syllable-timed than German (Ordin & Polyanskaya, 2015a;b). We performed a power
3 analysis to decide on the number of participants necessary to achieve significance, assuming at least a
4 medium effect size (for $p < 0.05$, partial eta square used for the effect size measure). Based on the power
5 analysis, we recruited 25 French (Paris 3: Université Sorbonne Nouvelle, age: 18-30 y.o., median age=24
6 y.o., 18 females) and 25 German (Bielefeld University, age: 18-35 y.o., median age=24 y.o., 14 females)
7 undergraduates. None of the participants reported any speech or hearing problem or proficiency in a
8 foreign language.

19 2.1.2 STIMULI

21 We used an inventory of five consonants ([s,z,f,v,ʃ]) and five vowels ([a,u,i,e,o]). We decided to limit the
22 consonants to fricatives because, unlike plosives, they allow stretching and compressing without losing
23 naturalness, and unlike sonorants, they are difficult to confuse with vocalic intervals even when their
24 durations are short. In pilot testing, synthesized stimuli with only fricative consonants received the highest
25 ratings for naturalness.

26 Concatenations of the five consonants and five vowels produce 25 possible CV syllables. We created
27 two random sequences, each with 3000 syllables, using these 25 syllables (see below for the synthesis
28 procedure). Each syllable occurred 120 times in a sequence, with at least two different syllables between
29 repetitions of the same syllable. One sequence was used to create stimuli with a relatively high degree of
30 vocalic durational variability and lower %V (typical of stress-timed languages), and the other sequence was
31 used to create stimuli with lower vocalic durational variability and higher %V (typical of syllable-timed
32 languages). Each sequence was split into 120 25-syllable passages. In the stimuli with stress-timed rhythm,
33 vowel durations were between 10 ms. and 150 ms. with zero skewness from a normal distribution
34 centered around an 80 ms mean. In the stimuli with syllable-timed rhythm, vowel durations were between
35 80 ms. and 120 ms. with zero skewness from a normal distribution centered around a 100 ms mean.
36 Durations of consonants were calculated by subtracting the duration of each subsequent vowel from 200
37 ms., thus producing streams of 25 CV syllables, 200 ms. each. Different values of median vowel durations,

1 with constant syllable durations, lead to variations in the %V measure, producing rhythmic differences
2
3 between syllable- and stress-timed patterns. The metric scores, capturing the rhythmic distinctions we
4
5 implemented in our syllable sequences, are presented in Table 2. The implemented differences were
6
7 designed to be more extreme than those typical of natural languages (Grabe & Low, 2002; Ramus &
8
9 Mehler, 1999; White & Mattys, 2007) to ensure that the rhythmic differences would be easily perceivable.
10
11 These durations were fed into the MBROLA speech synthesizer (Dutoit et al., 1996) to synthesize 120
12
13 stimuli (25-syllable sequences) of each rhythm, using the IT4 voice.
14
15

16
17 To construct AX discrimination pairs, we presented the stimuli with a 1-second pause between the
18
19 two members of a pair. There were 30 pairs in which the rhythm before and after the pause was syllable-
20
21 timed (syl-syl), 30 pairs in which the rhythm before and after the pause was stress-timed (str-str), 30 pairs
22
23 in which the stimuli with syllable-timed rhythm were followed by those with stress-timed rhythm (syl-str),
24
25 and 30 pairs in which stimuli with stress-timed rhythm were followed by stimuli with syllable-timed rhythm
26
27 (str-syl). Note that all pairs differed in their phonetic content due to the random concatenation procedure
28
29 for sequences. Each stimulus was paired only once, in one of the 120 pairs.
30
31
32

33 2.1.3 PROCEDURE

34
35 The experiment was carried out in sound-treated booths. The participants were told that they would listen
36
37 to utterances in an “extraterrestrial language” followed by a pause and by utterances either in the same or
38
39 a different language. Their task was to identify whether the utterances before and after the pause were in
40
41 the same or a different language. The answer was given by pressing the button “1” or “2” on the keyboard.
42
43 The buttons “1” and “2” to respond “same” or “different” were counterbalanced between listeners. The
44
45 participants were instructed not to wait until the utterance after the pause finished, and to respond
46
47 immediately when they recognized whether the languages before and after the pause were different or
48
49 the same. The order of pairs was randomized. To familiarize the participants with the procedure, the
50
51 experiment was preceded by a training session with 12 additionally prepared pairs, with feedback as to the
52
53 accuracy of the response given on each trial. Accuracy of responses and RTs (onsets of the recorded RTs
54
55 locked to the onsets of the X stimuli in the AX discrimination pairs) were subjected to statistical analysis.
56
57
58
59
60

2.2 RESULTS

2.2.1 DEALING WITH OUTLIERS

As the listeners need at least three syllables to evaluate the durational variability, we excluded 38 responses with RTs less than 600ms. An additional 33 responses were excluded because the RT exceeded the mean RT for the participant +2SE. In total, 1.18% of all responses were excluded.

2.2.2 RESPONSE ACCURACY

As the AX discrimination task can be seen as a change detection task, we decided to adopt a signal detection theory (SDT) approach. SDT accuracy (or “sensitivity”) measures are designed to be independent of any response bias (e.g., if a participant generally tends to respond “X is different from A”). The statistical procedures separate decision factors from sensory ones, which was our primary rationale for this methodological choice. In addition, sensitivity measures are comparable between experiments because they are expressed in the same units.

Sensitivity and bias measures of participants’ responses were computed both for the full set of experimental items and for the subsets in which the first element of the pair was syllable- or stress-timed. We used the measures A' and B_D'' as presented by Donaldson (1992) and implemented in R by Pallier (2002). These measures represent nonparametric alternatives to the classical measures d' and β of signal detection theory, and were preferred because A' can be computed even in cases where d' is infinite or undefined. A' ranges between 1.0 and 0.0: a value of 1.0 means that detection performance is perfect, whereas a value of 0.5 means that performance is at chance. A' values that are significantly above 0.5 indicate sensitivity in detecting rhythmic changes. The bias index B_D'' ranges between -1.0 and 1.0. A value of 0.0 means that a listener shows no bias in reporting whether A and X differ or not, a value of 1.0 means that the listener always responds that X has the same rhythm as A, and a value of -1.0 means that the listener always responds that X has a different rhythm than A. Thus, biases in detection can be assessed by checking whether B_D'' differs significantly from 0.0. B_D'' is undefined when classification performance is fully correct or fully incorrect (i.e., when the hit rate equals 1.0 and the false alarm rate equals 0.0, or vice versa; this shortcoming also applies to β). Undefined values of B_D'' were obtained for three participants

(two French and one German) when analyzing the items in which A was syllable-timed, and for one German participant when analyzing the items in which A was stress-timed. The degrees of freedom of the statistical tests for B_D'' were adjusted accordingly, reflecting the reduced number of subjects in those cases.

The analysis of all items together showed that detection of rhythmic changes was good in both language groups, $A' = .858$ for German listeners, significantly above 0.5, $t(24) = 15.36$, $p < .0005$; and $A' = .832$ for French listeners, also significantly above 0.5, $t(24) = 11.22$, $p < .0005$ (see **Figure 1**; error bars stand for $\pm 2SE$ in all figures). These A' values were statistically indistinguishable from one another, $t(48) = 0.68$, $p = .5$. Response bias was low for both language groups: $B_D'' = 0.143$ for Germans, indistinguishable from 0, $t(24) = 1.29$, $p = .21$, and $B_D'' = 0.155$ for French, again indistinguishable from 0, $t(24) = 1.58$, $p = .13$. The B_D'' values for German and French listeners were not different from each other, $t(48) = -0.077$, $p = .94$.

When we analyzed the responses only to the stimuli in which A was stress-timed (i.e., str-str and str-syl stimuli pairs), detection sensitivity remained high: $A' = .825$ for Germans, significantly above 0.5, $t(24) = 10.08$, $p < .0005$, and $A' = .776$ for French, significantly above 0.5, $t(24) = 7.21$, $p = .0005$. A' values for German and French listeners did not differ from each other, $t(48) = 0.99$, $p = .33$. Response bias was low for both groups: $B_D'' = 0.138$ for Germans, indistinguishable from 0, $t(23) = 1.03$, $p = .32$, and $B_D'' = 0.044$ for French, indistinguishable from 0 as well, $t(24) = 0.39$, $p = .7$. B_D'' values for German and French listeners did not differ from each other, $t(47) = 0.54$, $p = .59$.

The analysis of responses to the stimuli in which A was syllable-timed (i.e., syl-syl and syl-str stimuli pairs) yielded a similar pattern of results. Detection of rhythmic changes was good, $A' = .881$ for Germans, significantly above 0.5, $t(24) = 17.90$, $p < .0005$, and $A' = .877$ for French listeners, significantly above 0.5, $t(24) = 14.54$, $p = .0005$. A' values for French and German participants did not differ from each other, $t(48) = 0.13$, $p = .90$. Response bias was low for Germans, $B_D'' = 0.056$, indistinguishable from 0, $t(23) = 0.42$, $p = .68$. It was somewhat higher for French listeners, $B_D'' = 0.222$, but this value was not significantly different from 0 either, $t(22) = 1.84$, $p = .088$. The difference between B_D'' values for the two language groups was not significant, $t(45) = -0.93$, $p = .36$.

1 Contrasting performance in items in which A was syllable-timed vs. items in which A was stress-
2 timed, we confirmed that both French and German listeners showed higher sensitivity in detecting
3 rhythmic changes when A was syllable-timed than when A was stress-timed, $t(24) = 2.30$, $p = .03$ for
4 Germans and $t(24) = 4.27$, $p = .0003$ for French participants.

5
6
7
8
9
10 These sensitivity results provide preliminary answers to the central questions of the current study.
11 First, language background did not affect performance, indicating that rhythm processing here was
12 dominated by language-independent factors. Second, listeners were more accurate when the vowel onsets
13 in the first stimulus in a pair were relatively regular, with little variability in vowel durations.

14 2.2.3 REACTION TIME

15 The error rates were 22.0% and 20.2% for the French and the German listeners, respectively. An ANOVA on
16 the reaction times, with *L1* of the listener and *Accuracy* (correct vs. incorrect) as factors, revealed a
17 significant effect of *Accuracy*, $F(1, 48) = 36.896$, $p < .0005$, $\eta_p^2 = .435$. There was no effect of *L1*, $F(1,48) =$
18 $.437$, $p = .512$ and no significant interaction, $F(1,48) = 1.004$, $p = .321$. Both German and French listeners
19 responded more slowly when they gave an incorrect answer (see **Figure 2**). Therefore, the following
20 analyses of reaction times were performed only using correct responses, to avoid contaminating the
21 results with longer RTs on incorrect responses.

22 To explore the impact of linguistic experience on the RT for discriminating between rhythm types,
23 we performed a 2-way ANOVA with *L1* and *Pair-type* as factors. The four pair-types consisted of the 2x2
24 crossing of stress-timed and syllable-timed sequences with the two positions in an AX trial. **Figure 3** shows
25 RTs split by *L1* and *Pair-type*. The effect of native language was not significant, $F(1,48) = .112$, $p = .739$;
26 German and French listeners did not differ in their time to decide whether the rhythm of the two speech
27 stimuli was the same or different. The effect of *Pair-type* was significant, $F(3,46) = 5.033$, $p = .004$, $\eta_p^2 =$
28 $.247$: Pairs with a syllable-timed stimulus in the first position were responded to more quickly than those
29 with a stress-timed stimulus in that position. There was no interaction between *L1* and *Pair-type*, $F(3,46) =$
30 $.092$, $p = .964$, $\eta_p^2 = .006$, reflecting a similar advantage for the initial syllable-timed item regardless of the
31 listeners' linguistic backgrounds. Pairwise comparisons (Bonferroni corrected) confirmed that an initial

1 syllable-timed item produced faster responses both when the second item in a pair was syllable-timed, $p =$
2
3 .014; and when the second item was stress-timed, $p = .022$. Listeners responded significantly more quickly
4
5 when the first stimulus in the pair had a regular temporal structure. The facilitatory effect of a regular
6
7 temporal structure was not modulated by native language of the listener. Thus, the conclusions from the
8
9 reaction time data are exactly those that the sensitivity analyses supported.
10

11
12 To make sure that the null effect of linguistic background does not stem from a lack of statistical
13
14 power, we computed Bayes factors using the BayesFactor package in R (Morey & Rouder, 2012). Bayes
15
16 Factor (BF) is the ratio between the likelihood of a set of observed data given two statistical models. If
17
18 these two models are assumed to be equally probable a priori, then the BF corresponds to the ratio of
19
20 probabilities of these models, given the observed data. This means that BFs may be useful not only to
21
22 quantify the degree of support that some data provide for the alternative hypothesis, but also to quantify
23
24 whether they support the null hypothesis² (Morey & Rouder, 2011).
25
26
27
28

29
30 We conducted Bayesian analyses on accuracy rates and response times considering each of the four
31
32 item types (syl-syl, syl-str, str-syl, str-str) separately. We analyzed accuracies and response times (for
33
34 correct answers only) using 4×2 mixed models, with item type as a within-subject factor, native language
35
36 as a between-subject factor, and subjects as a random factor. We computed BFs contrasting each full
37
38 regression model (both factors plus interaction) against the corresponding model with neither the
39
40 interaction nor the main effect of native language. The BFs obtained were 0.0287 for accuracy and 0.0248
41
42 for response times. These numbers indicate that outcomes of both measures strongly support the null
43
44 model, a model lacking any effect of native language. That is to say, the null models are 34.9 times and
45
46 40.3 times (for accuracy and response time, respectively) more likely than the full regression models, given
47
48 the data. Moreover, if the full models and the models with no effect of native language are assumed to be
49
50 equally probable (50%) a priori, these BFs indicate that the probability of the model lacking native language
51
52 grows to 97.2% a posteriori for accuracy, and to 97.6% a posteriori for response times.
53
54
55
56
57

58 2.3 DISCUSSION

59
60

²For a simple statistical test (e.g., a t-test), a BF of 1.0 means that the data provide evidence in favor of neither the alternative nor the null hypothesis, whereas a BF smaller (resp. greater) than 1.0 means that the data support the null (resp. alternative) hypothesis.

1 The results showed that people were quite good in general at discriminating rhythms in all conditions, but
2 detecting the rhythm change in one direction (from more regular to irregular) was faster and more reliable
3 than in the opposite direction. When the A stimulus was syllable-timed, participants' sensitivity was higher,
4 and RT was lower – for both French and German listeners – compared to the conditions when the A
5 stimulus was stress-timed. This asymmetry, similar to one reported by McAuley and Fromboluti (2015) for
6 non-linguistic tone sequences, is shared by listeners with rhythmically different native languages,
7 supporting the view that it is grounded in a general auditory mechanism. Although Germans are familiar
8 with a wider spectrum of rhythmic patterns than French listeners, there was no influence of the native
9 language on sensitivity or RT. Therefore, we tentatively conclude that the ability to discriminate utterances
10 with different linguistic rhythms is not biased by linguistic experience.
11
12
13
14
15
16
17
18
19
20
21
22
23
24

25 Experiment 1 only tested the efficiency of rhythm discrimination by native speakers of one pair of
26 rhythmically different languages: German and French. If we wish to draw broad conclusions, it is necessary
27 to show that the observed pattern is not limited to a particular pair of languages. To this end, a second AX
28 discrimination experiment was conducted using a different pair of languages that differ in terms of being
29 more syllable-timed (Spanish) or more stress-timed (English). Another objective was to include two
30 different sets of stimuli: one that is perceived as being made up of utterances in a natural language, and
31 the other, which is perceived as being less like natural language. This manipulation is designed to
32 investigate whether the performance in the rhythm discrimination task is modulated by whether the
33 stimuli are linguistic or non-linguistic. Finally, we made the stimuli more ecological by modelling the
34 durational values based on the values of a real language and by allowing a wider variety of syllable types,
35 which is typical of real languages.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51 3. EXPERIMENT II

52 3.1 METHOD

53 3.1.1 PARTICIPANTS

54 We recruited native speakers of two rhythmically different languages – Spanish and English – Spanish being
55 more syllable-timed than English (Payne et al., 2012; Prieto et al., 2012; Ramus et al., 1999). We recruited
56
57
58
59
60

29 participants per group, but we had to discard two Spanish participants for not performing the task as instructed. None of the participants reported any speech or hearing problem. Based on self-report, Spanish participants were fluent in Basque (another syllable-timed language), and English participants were fluent in Spanish, with proficiency varying from intermediate to high, which increases their linguistic experience with both syllable- and stress-timed rhythmic patterns.

3.1.2 STIMULI

For this experiment, we created two different types of stimuli, which we will refer to as *linguistic* stimuli and *non-linguistic* stimuli. Linguistic stimuli were intended to encourage listeners to engage the processing mechanisms that are brought to bear while listening to real language. Linguistic stimuli should be perceived as plausible utterances of a natural language, manifesting natural prosody and properties of a real human language. Non-linguistic stimuli are intended to be speech-like sequences of syllables, which, nevertheless, are processed at more of a psychoacoustic level, without engaging additional mechanisms involved in listening to real language. For this, we used random syllabic sequences with a monotonous pitch. These were devoid of prosody, recognizable hierarchical linguistic structures, or segmentable discrete constituents, thus making the stimuli sound less like real language.

To prepare non-linguistic stimuli, we used an inventory of five consonants [s,m,n,l,f] and five vowels [a,u,i,e,o]) to create 25 possible CV syllables. We created 240 sequences, in which the 25 possible syllables were randomized, with each syllable occurring only once per sequence. In each sequence, five random syllables were marked as stressed, with the only restriction being that two consecutive syllables could not be stressed.

Half of the sequences were used to prepare stimuli with high durational variability (more stress-timed), and the other half to prepare stimuli with low durational variability (more syllable-timed). In stimuli with higher durational variability, the vowel durations for the stressed positions varied between 80ms and 240 ms with 40-ms steps (one value per syllable per sequence), and vowel durations in unstressed syllables varied between 40ms and 80 ms with 10-ms steps (four values per syllable per sequence). In stimuli with lower durational variability, the vowel durations for the stressed positions varied between 160ms and 200

ms with 10-ms steps (one value per syllable per sequence), and vowel durations in unstressed syllables varied between 80ms and 100 ms with 5-ms steps (four values per syllable per sequence). The total syllable duration was set to 280 ms for stressed syllables and 180 ms for unstressed syllables. Durations of consonants were calculated by subtracting the duration of each subsequent vowel from 280 ms (stressed syllables) or 180 ms (for unstressed syllables). In this way, we ensured that the stimuli differed in rhythmic characteristics (durational variability), but not in speech rate (number of syllables or phonemes per second). The sequences were synthesized with these phoneme durational parameters in MBROLA, with the ES2 voice, F0 set to 200Hz.

To prepare linguistic stimuli, we took the recordings of a native Welsh speaker reading 38 sentences (this was a Welsh-English bilingual recorded for a different experiment reported in Ordin & Mennen, 2017). The sentences were annotated in Praat, and durations of vowels and consonants were measured. We used a Spanish set of phonemes to re-synthesize the Welsh sentences with the durations of vowels and consonants of a Welsh speaker. In case a Welsh phoneme did not exist in Spanish, it was substituted with the closest Spanish phoneme. Then we imposed an intonational contour of a Welsh sentence on resynthesized versions of the sentences. The values of the rhythm metrics are similar to those of Spanish and contrastive to those of English, which leads to Welsh being positioned closer to the syllable-timed end of the rhythm spectrum (Grabe & Low., 2002). This was confirmed by the scores of vowel duration variability calculated on the individual's recording used in our study.

We created a second set of sentences in which we multiplied the durations of stressed vowels by 2 and unstressed vowels by 0.6 in order to enhance the durational contrasts between stressed and unstressed vowels. The original re-synthesized sentences represent more syllable-timed rhythms, and the re-synthesized sentences with enhanced durational ratios represent more stress-timed rhythms. The sentences were concatenated into 2-sentence pairs (each pair making one stimulus), 5.5 seconds in duration (± 500 ms), thus matching the duration of non-linguistic and linguistic stimuli. Each sentence was used 5-7 times for the stimuli of each rhythm type. We created 120 linguistic stimuli of each rhythm type. A preliminary pilot experiment was performed to make sure that the linguistic stimuli in our experiment

1 received higher naturalness ratings and were more likely to be perceived as real language compared to the
2 non-linguistic stimuli (see Appendix I for the results of the pilot study).
3
4

5 We constructed 120 pairs of stimuli for each stimulus type (linguistic and non-linguistic) with a 1.5-
6 second pause between the two members of a pair. For each stimulus type, there were 30 pairs in which
7 the rhythm before and after the pause was syllable-timed (syl-syl), 30 pairs in which the rhythm before and
8 after the pause was stress-timed (str-str), 30 pairs in which the stimuli with syllable-timed rhythm were
9 followed by those with stress-timed rhythm (syl-str), and 30 pairs in which stimuli with stress-timed
10 rhythm were followed by stimuli with syllable-timed rhythm (str-syl). Each stimulus was paired only once,
11 in one of the 120 pairs.
12
13
14
15
16
17
18
19
20
21

22 In the first experiment the differences in durational variability between rhythms were pushed to
23 the extremes, whereas in the second experiment, the scores of the rhythm metrics were based on the
24 values typical of natural languages. This made the differences in durational variability more natural and
25 less salient, thus making the discrimination task more challenging.
26
27
28
29
30
31

32 3.1.3 PROCEDURE

33 Experiment 2 consisted of two sessions held with an interval of 10 days (± 3 days) between the sessions. In
34 one of the sessions, participants listened to linguistic stimuli, and in the other session the stimuli were non-
35 linguistic. The order of sessions was counterbalanced across participants. The experiment was carried out
36 in sound-attenuated booths. The participants were told that they would listen to pairs of utterances in an
37 “extraterrestrial language”. Their task was to identify whether the rhythm in the utterances before and
38 after the pause was the same or different. The answer was given by pressing the button “1” or “2” on the
39 keyboard. The buttons “1” and “2” to respond “same” or “different” were counterbalanced between
40 listeners. The order of pairs was randomized. To familiarize the participants with the procedure, the
41 experiment was preceded by a training session with 12 additionally prepared pairs, with feedback as to the
42 accuracy of the response given on each trial during training.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58

59 Unlike the first experiment, when we asked participants not to wait until the end of the second
60 utterance and to make the response as soon as they thought they knew the answer, in this experiment the

1 response could only be given after the participant had finished listening to the second utterance in the
2 trial. This was deemed necessary because the task was subjectively much more challenging, and we
3 wanted to prevent a possible speed-accuracy trade-off by giving the participants the same time on each
4 trial to make their decision. Note that this approach prevented us from using RT as a dependent measure
5 this time.
6
7
8
9
10
11

12 At the end of the second session, participants performed an additional test. Ten randomly chosen
13 stimuli of each type (linguistic and non-linguistic stimuli), half with high variability durations and half with
14 low, were played to the participants in a random order. Upon listening to each stimulus, the participant
15 was asked to mark on an 8-point scale how likely it seemed that the stimulus represents an utterance of a
16 real language. This was done to ensure that individual participants indeed perceived the Welsh-based
17 stimuli as more linguistic than random syllabic concatenations. The analysis of this additional test
18 confirmed that our linguistic stimuli were indeed perceived as more like a real language; the rhythmic type
19 had no effect on the perceived naturalness of the stimuli (see Appendix II for the test results).
20
21
22
23
24
25
26
27
28
29
30
31

32 3.2 RESULTS

33 We kept all the responses and participants; no outliers were removed. As in Experiment 1, we used the
34 non-parametric measures A' and B_D'' to index sensitivity and bias.
35
36
37
38

39 3.2.1 LINGUISTIC STIMULI RESULTS

40 We computed the sensitivity and bias measures of participants' responses for the full set of experimental
41 items and for the subsets in which the first element of the pair was syllable- or stress-timed. **Figure 4a**
42 shows the sensitivity results. The analysis on the full set showed that detection of rhythmic changes was
43 significantly above chance for both language groups, $A' = .642$ for English listeners, $t(28) = 10.62, p < .0001$;
44 and $A' = .590$ for Spanish listeners, $t(26) = 4.79, p < .0001$. A' values were significantly larger for English
45 listeners, $t(54) = 2.30, p = .03$. **Figure 4b** shows that the response bias was low for both languages: $B_D'' =$
46 0.017 for English, indistinguishable from 0, $t(28) = 0.28, p = .78$ and $B_D'' = 0.052$ for Spanish, again
47 indistinguishable from 0, $t(26) = 0.82, p = .42$. Moreover, B_D'' values for English and Spanish listeners were
48 not different from each other, $t(54) = -0.41, p = .68$.
49
50
51
52
53
54
55
56
57
58
59
60

1 The analysis of responses only to the stimuli in which A was syllable-timed revealed good detection
2 of rhythmic changes, $A' = .708$ for English, significantly above 0.5, $t(28) = 10.93$, $p < .0001$, and $A' = .652$ for
3 Spanish listeners, significantly above 0.5 as well, $t(26) = 6.37$, $p < .0001$. A' values for English and Spanish
4 participants marginally differed from each other, $t(54) = 1.86$, $p = .07$. Response bias was low for English,
5 $B_D'' = 0.082$, indistinguishable from 0, $t(28) = 1.06$, $p = .30$, as well as for Spanish listeners, $B_D'' = 0.160$. This
6 value was marginally different from 0, $t(26) = 1.97$, $p = .06$. The difference between B_D'' values for the two
7 languages was not significant, $t(54) = -0.70$, $p = .49$.

8
9
10
11
12
13
14
15
16
17 When we analyzed the responses only to the stimuli in which A was stress-timed, detection
18 sensitivity dropped: $A' = .555$ for English, significantly above 0.5, $t(28) = 2.55$, $p = .02$, and $A' = .516$ for
19 Spanish, statistically indistinguishable from 0.5, $t(26) = 0.73$, $p = .47$. A' values for English and Spanish
20 listeners did not differ from each other, $t(54) = 1.28$, $p = .21$. Response bias was low for both groups: $B_D'' =$
21 -0.050 for English, indistinguishable from 0, $t(28) = -0.82$, $p = .42$, and $B_D'' = -0.061$ for Spanish,
22 indistinguishable from 0 as well, $t(26) = -0.95$, $p = .35$. B_D'' values for English and Spanish listeners did not
23 differ from each other, $t(54) = 0.12$, $p = .91$.

24
25
26
27
28
29
30
31
32
33
34 Contrasting performance on items in which A was syllable-timed vs. items in which A was stress-
35 timed, we confirmed that both English and Spanish listeners showed higher sensitivity in detecting
36 rhythmic changes when A was syllable-timed than when A was stress-timed, $t(28) = 5.04$, $p < .0001$ for
37 English and $t(26) = 4.99$, $p < .0001$ for Spanish participants.

38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
Altogether, the analysis of sensitivity shows that both English and Spanish listeners were good at
discriminating the two rhythms only when the first stimulus in the test pairs was syllable-timed (i.e.,
rhythmically regular). Performance was substantially worse (at the chance level for Spanish participants
and slightly above chance for the English participants, with no significant difference between the groups)
when the first stimulus was stress-timed. This shows that listeners were more accurate when the first
stimulus in the pair is characterized by a regular distribution of vowel onsets, with relatively little variability
in vowel durations.

3.2.2 NON-LINGUISTIC STIMULI RESULTS

1 The analysis of all items together showed that detection of rhythmic changes was above chance in both
2 language groups, $A' = .572$ for English listeners, significantly above 0.5, $t(28) = 4.47$, $p = .0001$; and $A' = .563$
3
4 for Spanish listeners, also significantly above 0.5, $t(26) = 3.24$, $p = .003$ (see **Figure 4a**). The two A' values
5
6 were statistically indistinguishable from one another, $t(54) = 0.36$, $p = .72$. Response bias was low for both
7
8 languages: $B_D'' = 0.048$ for English, indistinguishable from 0, $t(28) = 0.78$, $p = .44$ and $B_D'' = 0.032$ for
9
10 Spanish, again indistinguishable from 0, $t(26) = 0.51$, $p = .62$. Moreover, B_D'' values for English and Spanish
11
12 listeners were not different from each other, $t(54) = 0.18$, $p = .86$.

13
14
15
16
17 The analysis of responses only to the stimuli in which A was syllable-timed yielded a similar pattern
18
19 of results. Detection of rhythmic changes was good, $A' = .631$ for English, significantly above 0.5, $t(28) =$
20
21 6.71 , $p < .0001$, and $A' = .652$ for Spanish listeners, significantly above 0.5, $t(26) = 5.68$, $p < .0001$. A' values
22
23 for English and Spanish participants did not differ from each other, $t(54) = -0.66$, $p = .51$. Response bias was
24
25 low for English, $B_D'' = 0.081$, indistinguishable from 0, $t(28) = 1.13$, $p = .27$, as well as for Spanish listeners,
26
27 $B_D'' = 0.065$, $t(26) = 0.73$, $p = .47$. The difference between B_D'' values for the two languages was not
28
29 significant, $t(54) = 0.14$, $p = .89$.

30
31
32
33
34 When we analyzed the responses only to the stimuli in which A was stress-timed, detection
35
36 sensitivity vanished: $A' = .505$ for English, not significantly above 0.5, $t(28) = 0.18$, $p = .86$, and $A' = .463$ for
37
38 Spanish, not significantly above 0.5, $t(26) = -1.37$, $p = .18$. A' values for English and Spanish listeners did not
39
40 differ from each other, $t(54) = 1.09$, $p = .28$. Response bias was low for both groups: $B_D'' = 0.022$ for
41
42 English, indistinguishable from 0, $t(28) = 0.29$, $p = .78$, and $B_D'' = 0.022$ for Spanish, indistinguishable from 0
43
44 as well, $t(26) = 0.27$, $p = .79$. B_D'' values for English and Spanish listeners did not differ from each other,
45
46
47
48
49 $t(54) = 0.0004$, $p = .99$.

50
51
52 Contrasting performance on items in which A was syllable-timed vs. items in which A was stress-
53
54 timed, we confirmed that both English and Spanish listeners showed higher sensitivity in detecting
55
56 rhythmic changes when A was syllable-timed than when A was stress-timed, $t(28) = 3.75$, $p = .0008$ for
57
58 English and $t(26) = 4.82$, $p < .0001$ for Spanish participants.

1 Altogether, the analysis revealed the same result pattern that we observed with linguistic stimuli:
2 both Spanish and English listeners performed better on the task when the first stimulus in a test pair was
3 characterized by regular rhythm.
4

5 3.2.3 COMPARISON OF SENSITIVITY ON LINGUISTIC AND NON-LINGUISTIC STIMULI

6 We performed repeated-measures ANOVAs on A' scores with the *rhythm in the first stimulus* in the AX pair
7 (stress- vs. syllable-timed) and the *type of the stimuli* (linguistic vs. non-linguistic) as within-subject factors,
8 and with *native language of the listener* (Spanish vs. English) as a between-subject factor. The analysis
9 showed that sensitivity was significantly higher both for linguistic and non-linguistic stimuli when the first
10 stimulus was syllable-timed, $F(1,54) = 63.495$, $p < .0005$, $\eta_p^2 = .54$ (η_p^2 is partial eta squared). Overall,
11 participants performed slightly better on linguistic stimuli than on non-linguistic stimuli, $F(1,54) = 7.361$, p
12 $= .009$, $\eta_p^2 = .12$. The effect of the native language, on the other hand, was not significant, $F(1,54) = 3.05$,
13 $p = .086$, $\eta_p^2 = .05$. Importantly, native language modulated neither the effect of the stimulus type (the
14 language*stimulus type interaction was not significant, $p = .257$), nor the effect of the first stimulus in AX
15 pairs (language*rhythm in the first stimulus of a pair was not significant, $p = .551$). The three-way
16 interaction of *language*preceding stimulus rhythm*stimulus type* was not significant either, $p = .142$.

17 This analysis shows that the results are stable for both stimulus types and language groups. Overall,
18 the effects of the stimulus type (linguistic vs. non-linguistic) and of the rhythmic characteristics of the first
19 stimulus in the test pairs significantly affected sensitivity, whereas the effect of the native language
20 produced a null result. Overall, both language groups yielded similar results, but the native English listeners
21 showed some sensitivity on linguistic stimuli, provided that the first stimulus in the discrimination pair was
22 not stress-timed. To estimate the support for the null hypothesis regarding the effect of the native
23 language on performance in Experiment 2, we analyzed the accuracy rates for this experiment using Bayes
24 factors in a similar manner to the previous experiment. The BFs obtained were 0.0444 for linguistic stimuli
25 and 0.0157 for non-linguistic stimuli, yielding strong evidence for the null models. These are 22.5 times and
26 63.7 times more likely than the full regression models, given the data. If the two models are assumed to be
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 equally probable *a priori*, then the probability of the model lacking native language grows to 95.7% *a*
2
3 *posteriori* for linguistic stimuli, and to 98.5% *a posteriori* for non-linguistic stimuli.

4 3.3 DISCUSSION

5
6
7 The result pattern observed in the second experiment confirmed and strengthened the preliminary
8
9 conclusions offered in the first experiment. The ability to discriminate rhythmic patterns in speech is not
10
11 affected by linguistic experience stemming from the native language of the speaker. Rhythm discrimination
12
13 is facilitated if the first stimulus manifests regular rhythm. If the first stimulus has an irregular rhythm, the
14
15 task becomes more challenging and performance drops, even if the native language of the participants
16
17 includes utterances with irregular, stress-timed rhythm. The same pattern was observed for non-linguistic
18
19 and linguistic stimuli. The former may not engage the full set of speech processing mechanisms, while the
20
21 latter presumably do because they are perceived as the utterances of a real language. These results
22
23 suggest that rhythm discrimination in speech is controlled to a large extent by cognitive mechanisms
24
25 shared by all humans, irrespective of their native language, based on the properties and the general design
26
27 of the auditory system.
28
29
30
31
32
33

34 4. GENERAL DISCUSSION

35
36 Although people indeed find it easier to detect slight rhythmic differences in linguistic than in non-linguistic
37
38 stimuli, our results showed that discrimination of rhythms even in linguistic stimuli is not affected by
39
40 linguistic experience. Instead, performance is dependent on the rhythmic characteristics of the first
41
42 stimulus in an AX pair. Our results revealed that utterances characterized by higher regularity of vowel
43
44 onsets and lower variability in vowel durations enhanced attention to the rhythmic details of the acoustic
45
46 signal. As a result, performance on the AX rhythm discrimination task was better when the X stimulus was
47
48 preceded by a syllable-timed A stimulus.
49
50
51
52
53

54 This pattern shows that the native language is not a factor in rhythm discrimination performance,
55
56 which suggests that the stimuli are not filtered via the phonology of the native language, even when these
57
58 stimuli are perceived as speech in a natural language. Given that rhythm discrimination is so important for
59
60 speech processing and language acquisition (Hickok et al., 2015; Mehler et al., 1996; Nazzi & Ramus, 2003),

1 that strategies of speech segmentation and acquisition are so different in rhythmically different languages
2 (e.g., Cutler & Butterfield, 1992; Kim et al., 2008; Langus, Mehler & Nespors, 2018), and that deviations
3
4 from the native norms in rhythmic structure affect speech intelligibility and increase accentedness of
5
6 foreign speech (Polyanskaya et al., 2017; Tajima et al., 1997), it is surprising that linguistic experience plays
7
8 a very minor role in performance on the rhythm discrimination task. As performance is similar on linguistic
9
10 and non-linguistic stimuli, it appears that rhythm discrimination is primarily based on a domain-general
11
12 mechanism defined by the general design of the auditory system. Presumably, rhythm discrimination
13
14 happens at early stages of auditory processing, irrespective of whether the stimuli are linguistic or not,
15
16 before language processing skills are engaged. The auditory system is designed so that acoustic rhythms
17
18 modulate the firing pattern of the auditory nerves (Ghitza et al., 2013; Hickok et al., 2015), which leads to
19
20 coupling of neural cortical oscillations and the acoustic rhythms (Howard & Poeppel, 2012), resulting in a
21
22 facilitatory effect of a regular rhythm in the first stimulus. Rhythms characterized by regularity couple with
23
24 internal oscillators and facilitate attention better than irregular rhythms (Jones, Moynihan, MacKenzie, &
25
26 Puente, 2002).

27
28
29
30
31
32
33
34 This interpretation is in line with the Dynamical Attending Theory (Barnes & Jones, 2000; Howard &
35
36 Poeppel, 2010; Jones, et al., 2002), which predicts that the presentation of a regular rhythm entrains
37
38 attentional oscillations, which in turn creates stronger expectations for upcoming events and thus leads to
39
40 better discrimination performance. Regular rhythm, or the regular occurrence of vowel onsets, supports
41
42 anticipatory mechanisms for temporal prediction of when the following onset is expected to happen;
43
44 deviations from these expectations are therefore detected more reliably and faster (Barnes & Jones, 2000;
45
46 Lakatos et al., 2008; ten Oever, Schroeder, Poeppel, et al., 2014). Rhythm has a direct influence on the
47
48 perception of discrete acoustic events, including salient vowel onsets (Hickok et al., 2015). Selective active
49
50 attending to a stimulus generates a dynamically evolving neural model of the acoustic stream to which the
51
52 listener is attending, in the form of neuronal oscillations (Lakatos et al., 2013). Listeners are constantly
53
54 updating a reference pattern as the stimulus develops in time, based on expectations built on rhythmic
55
56 regularity (Lakatos et al., 2013), and building the reference is more difficult when the reference stimulus is
57
58
59
60

1 more stress-timed, i.e., exhibits irregular vowel onsets. The detection of rhythmic changes is then more
2
3 difficult because the reference is still being updated when the rhythmic change happens. The listener,
4
5 having no clear expectations, can only react to the change that has already happened. Regular rhythm
6
7 allows future-oriented attending, while irregular rhythm requires direct comparison of the rhythmic
8
9 patterns, using past-oriented attention, in order to compare the current rhythm with the reference rhythm
10
11
12 (ten Oever, Schroeder, Poeppel, et al., 2014).
13

14
15 The stimuli for the first experiment were based on extreme values of variability parameters to
16
17 model rhythmic differences, which reduce the challenge of the discrimination task but also potentially
18
19 reduce the ecological validity of the stimuli. The stimuli in the second experiment were based on the less
20
21 extreme values observed in real languages. This made the rhythmic differences subtler, and substantially
22
23 increased the difficulty of the task; hence the overall performance was lower in the second experiment.
24
25

26
27 Note that the non-linguistic stimuli were composed only of CV syllables. Rhythmically contrastive
28
29 languages also exhibit differences in syllabic complexity: Languages featuring a higher degree of stress-
30
31 timing allow more complex syllables due to loose phonotactic constraints (Prieto et al., 2012). In regard to
32
33 phonotactic complexity, all non-linguistic sequences were more syllable-timed than stress-timed. In
34
35 principle, this asymmetry could explain the observed facilitation effect for syllable- over stress-timed
36
37 stimuli. The linguistic stimuli, in contrast, have the properties characteristic of natural languages (e.g.,
38
39 prosody, statistical cues, and extractable constituents), thus engaging the additional mechanisms that are
40
41 involved in listening to real speech.
42
43
44

45
46 Although linguistic stimuli were indeed perceived as though they could have been utterances of
47
48 natural languages, the pattern of results for linguistic stimuli was not different from that for non-linguistic
49
50 ones. This equivalence illustrates the crucial role played by the general design of the auditory system and
51
52 low-level physiological mechanisms that are not affected by individual experience. The mammalian
53
54 auditory system has not evolved specifically for speech processing, and linguistic experience is unlikely to
55
56 shape the general design of the auditory system. On the contrary, it is more likely that the structure of
57
58 speech is shaped by the design of the auditory system, which ensures that the speech stream is
59
60

1 processable by a general auditory mechanism that allows entrainment of the neural oscillations to the
2 environmental rhythms. However, linguistic experience can modulate the *output (neuronal firing rate)* of
3 early auditory mechanisms, when the *output* is passed forward for processing at a higher level, e.g., via the
4 phonological filter of the listener's native language.
5
6
7
8

9
10 In the current study, we observed the outcome of such general perceptual mechanisms in adults
11 processing speech-like stimuli without access to higher-level linguistic information. The auditory
12 mechanism underlying rhythm discrimination, and probably rhythmic cognition in general, is not specific to
13 language processing. It is shared by humans irrespective of their native language. Together with other
14 rhythm-based general perceptual mechanisms (e.g., iambic-trochaic grouping regularities), it manifests
15 itself in the behavior of pre-linguistic babies (Mehler et al., 1996) and non-human species (de la Mora et
16 al., 2013; Tincoff, et al., 2005; Toro & Nespors, 2015). Universal rhythm processing mechanisms have even
17 been mentioned as the precursors to speech emergence and development in phylogenesis (Ghazanfar &
18 Takahashi, 2014 a; b; MacNeilage, 1998; Merker, Madison, & Eckerdal, 2009). The existence of general,
19 rather than language-specific mechanisms for extracting and classifying rhythmic patterns may be a pre-
20 requisite of language acquisition in ontogenesis (Mehler & Nespors, 2004). Psychoacoustic, low-level
21 processing of rhythmic patterns provides pre-linguistic infants recourse to language-independent cues for
22 differentiating and classifying utterances and languages (Nazzi, et al., 1998), for extracting linguistic
23 structure from the ambient language (Mehler, et al., 1996), and for segmentation of the continuous
24 acoustic stream into discrete constituents (Nazzi & Ramus, 2003).
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45

46 The distinction we are making between general auditory mechanisms and language-specific
47 processing is consistent with the assumptions of Werker and Curtin's (2005) PRIMIR model (Processing
48 Rich Information from Multidimensional Interactive Representations). The model posits that the
49 information extracted from speech is processed through filters based on (1) biological biases, (2) the
50 competence level of the listener in the language that is being processed, and (3) the requirements of the
51 specific task. These filters attract the listener's attention to one of three representational dimensions:
52 general perceptual dimensions, a wordform dimension, or a phonemic dimension (extendable to further
53
54
55
56
57
58
59
60

1 linguistic dimensions). Mapping of the signal simultaneously to both a general perceptual dimension and to
2 the wordform dimension is possible via the mechanism of statistical learning. To engage this mechanism,
3 statistical regularities in the input are necessary. The material in Experiment 1 and the non-linguistic
4 stimuli of Experiment 2 were designed to be devoid of these statistical regularities. As a result, the
5 conditions needed to activate the mechanism of statistical learning were not present. Thus, the results
6 only reflect processing at the general perceptual plane through general auditory mechanisms that are
7 based on biological biases, present in all humans irrespective of their native language.
8
9

10
11
12
13
14
15
16
17 If rhythm change detection is indeed controlled by general properties of the human auditory
18 system, and if the ability to discriminate rhythmic patterns of different utterances is a prerequisite for
19 successful language acquisition, then performance in such a task can potentially serve as one of the indices
20 of healthy phonological development. This suggestion is consistent with recent neurophysiological
21 evidence that the disruption of synchronization between acoustic and neural oscillations can lead to
22 deficits in both phonological and reading skills (Molinaro et al., 2016). Dyslexic readers exhibit impaired
23 neural entrainment to speech, with impaired coupling between neural oscillations in the auditory cortex
24 and the left inferior frontal region. This pattern is often accompanied by delays in phonological
25 development (Lallier et al., 2017). Goswami (2011) has proposed an integrated theoretical framework
26 postulating that phonological auditory deficits, stemming from atypical neural oscillations, lead to
27 developmental dyslexia. These theoretical and empirical developments suggest that a behavioral test of
28 rhythm detection on speech-like stimuli, which could be designed as a game for young children, could
29 provide a non-invasive diagnostic for phonological deficits. The procedures in the current study clearly
30 could be modified to develop such a test, though of course considerable further research would be needed
31 to establish clinically-relevant norms.
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

ACKNOWLEDGMENT:

The authors acknowledge support from the Spanish Ministry of Economy and Competitiveness (MINECO) Grant # PSI2017-82563-P (to AGS), from the 'Severo Ochoa' Programme for Centres/Units of Excellence in R&D (SEV-2015-490), and from the Basque Foundation for Science (IKERBASQUE). DMG was supported by Grant PIA/Basal FB0003 from the Chilean Research Council (CONICYT). LP was supported by the Spanish Ministry of Economy and Competitiveness (MINECO) via Juan de la Cierva fellowship. We thank two anonymous reviewers for their constructive suggestions on an earlier version of this paper.

For Peer Review

REFERENCES

- 1
2
3 Abercrombie, D. (1967). *Elements of General Phonetics*. Edinburgh: Edinburgh University Press, pp. 97-99.
4
5 Aslin, R. N., Saffran, J. R. and Newport, E. L. (1998). Computation of conditional probability statistics by 8-
6
7 month-old infants. *Psychological Science* 9,321-324.
8
9
10 Barnes, R., & Jones, M. (2000). Expectancy attention and time. *Cognitive Psychology* 41, 254-311.
11
12 Bunta, F., & Ingram, D. (2007). The acquisition of speech rhythm by bilingual Spanish- and English-speaking
13
14 four-and five-year-old children. *Journal of Speech, Language, and Hearing Research* 50, 999-1014.
15
16
17 Clopper, C. G., & Smiljanic, R. (2015). Regional variation in temporal organization in American English.
18
19 *Journal of Phonetics* 49, 1-15.
20
21
22 Content, A., Meunier, C., Kearns, R.K., & Frauenfelder, U. (2001). Sequence detection in pseudowords in
23
24 French: Where is the syllable effect? *Language and Cognitive Processes* 16, 609-636.
25
26
27 Cutler, A., & Butterfield, S. (1992). Rhythmic cues to speech segmentation: Evidence from juncture
28
29 misperception. *Journal of Memory and Language* 31, 218-236.
30
31
32 Dauer, R. (1983) Stress-timing and syllable-timing reanalyzed. *Journal of Phonetics* 11, 51-62.
33
34
35 De la Mora, D., Nespors, M. & Toro, J.M. (2013). Do Humans and non-human animals share the grouping
36
37 principles of the Iambic - Trochaic Law? *Attention, Perception, & Psychophysics* 75(1), 92-100.
38
39
40 Dellwo, V. (2006). Rhythm and Speech Rate: A Variation Coefficient for deltaC. In P. Karnowski & I. Szigeti
41
42 (Eds.). *Language and language-processing* (pp.231-241). Frankfurt am Main: Peter Lang.
43
44
45 Donaldson, W. (1992). Measuring recognition memory. *Journal of Experimental Psychology: General*,
46
47 121(3), 275-277.
48
49
50 Dutoit, T., Pagel, V., Pierret, N., Bataille, F., & van der Vrecken, O. (1996). The MBROLA project: Towards a
51
52 set of high-quality speech synthesizers free of use for non-commercial purposes. Philadelphia: ICSLP.
53
54
55 Ellis, R.J., & Jones, M. (2010). Rhythmic context modulates foreperiod effects. *Attention, Perception, &*
56
57 *Psychophysics* 72, 2274–2288.
58
59
60 Erra, R., Gervain, J. (2016). The Efficient Coding of Speech: Cross-Linguistic Differences. *PLoS ONE* 11(2),
<https://doi.org/10.1371/journal.pone.0148861>.

- 1 Gervain, J. Nespors, M., Mazuka, R., Horie, R., & Mehler J. (2008) Bootstrapping word order in prelexical
2 infants: a Japanese-Italian cross-linguistic study. *Cognitive Psychology* 57(1), 56-74.
- 3
4 Ghazanfar, A., & Takahashi, D. (2014a). Facial expressions and the evolution of the speech rhythm. *Journal*
5
6
7 of Cognitive Neuroscience 26(6), 1196-1207.
- 8
9 Ghazanfar, A., & Takahashi, D. (2014b). The evolution of speech: vision, rhythm, cooperation. *Trends in*
10
11
12 Cognitive Sciences 18(10), 543-553.
- 13
14 Ghitza, O. (2011). Linking speech perception and neurophysiology: speech decoding guided by cascaded
15
16
17 oscillators locked to the input rhythm. *Front. Psychol.* 2:130. <https://doi.org/10.3389/fpsyg.2011.00130>
- 18
19 Ghitza, O., Giraud, A.-L., & Poeppel, D. (2013) Neuronal oscillations and speech perception: critical-band
20
21
22 temporal envelopes are the essence. *Front. Hum. Neurosci.* 6:340.
23
24 <https://doi.org/10.3389/fnhum.2012.00340>
- 25
26 Goswami, U. (2011). A temporal sampling framework for dyslexia. *Trends in cognitive sciences*, 15(1), 3-10.
- 27
28 Grabe, E. & Low, L. (2002). Durational variability in speech and the rhythm class hypothesis. *Laboratory*
29
30
31 Phonology 7, 515-546.
- 32
33 Greenberg, S. & Ainsworth, W. (2004). Speech processing in the auditory system: An Overview. In S.
34
35
36 Greenberg, W. Ainsworth, A. Popper, & R.Fay (Eds.). *Speech Processing in the Auditory System* (pp. 1-62).
37
38
39 New York: Springer Verlag.
- 40
41 Hay, J., & Diehl, R. (2007). Perception of rhythmic grouping: Testing the Iambic/Trochaic law. *Perception &*
42
43
44 *Psychophysics* 69, 113–122.
- 45
46 Hickok, G., Farahbod, H., & Saberi, K. (2015). The Rhythm of Perception: Acoustic Rhythmic Entrainment
47
48
49 Induces Subsequent Perceptual Oscillation. *Psychological Science* 26(7), 1006–1013.
- 50
51 Howard, M., & Poeppel, D. (2010). Discrimination of speech stimuli based on neuronal response phase
52
53
54 patterns depends on acoustic but not comprehension. *Journal of Neurophysiology* 104, 2500-2511.
- 55
56 Howard, M., & Poeppel, D. (2012). The neuromagnetic response to spoken sentences: Co-modulation of
57
58
59 theta band amplitude and phase. *NeuroImage* 60, 2118-2127.
- 60
James, L. (1940). *Speech Signals in Telephony*. London: Sir Isaac Pitman & Sons.

- 1 Jones, M., Moynihan, H. M., MacKenzie, N., & Puente, J. K. (2002). Temporal aspects of stimulus-driven
2 attending in dynamic arrays. *Psychological Science* 13, 313–319.
- 3
4
5 Kent, R.D., Weismer, G., Kent, J.F., & Rosenbek, J.C. (1989). Toward phonetic intelligibility testing in
6
7 dysarthria. *Journal of Speech and Hearing Disorders* 54, 482–499.
- 8
9
10 Kim, J., Davis, C., & Cutler, A. (2008). Perceptual tests of rhythmic similarity: II. Syllable rhythm. *Language*
11
12 *and Speech*, 51(4), 343-359.
- 13
14
15 Ladefoged, P. (1993). *A course in phonetics*. 3rd edition. Fort Worth, TX: Harcourt Brace Jovanovich College
16
17 Publishers.
- 18
19
20 Lakatos, P., Karmos, G., Mehta, A. D., Ulbert, I., & Schroeder, C. E. (2008). Entrainment of neuronal
21
22 oscillations as a mechanism of attentional selection. *Science*, 320, 110–113.
- 23
24
25 Lakatos, P., Musacchia, G., O’Connel, M. N., Falchier, A. Y., Javitt, D. C., & Schroeder, C. E. (2013). The
26
27 spectrotemporal filter mechanism of auditory selective attention. *Neuron* 77, 750–761.
- 28
29
30 Lallier, M., Molinaro, N., Lizarazu, M., Bourguignon, M., & Carreiras, M. (2017) Amodal atypical neural
31
32 oscillatory activity in dyslexia: A cross-linguistic perspective. *Clinical Psychological Science* 5(2), 379-401.
- 33
34
35 Langus, A., Mehler, J., & Nespors, M. (2018). Rhythm in language acquisition. *Neuroscience & Biobehavioral*
36
37 *Reviews* 81, 158-166.
- 38
39
40 Lidji, P., Palmer, C., Peretz, I., & Morningstar, M. (2011). Listeners feel the beat: entrainment to English and
41
42 French speech rhythms. *Psychonomic Bulletin and Review* 18, 1035-1041.
- 43
44
45 Liss, J., White, L., Mattys, S., Lansford, K., Lotto, A., Spitzer, S. & Caviness, J. (2009). Quantifying Speech
46
47 Rhythm Abnormalities in the Dysarthrias. *Journal of Speech, Language and Hearing Research* 52, 1334-
48
49 1352.
- 50
51
52 Low, E.L., Grabe, E., & Nolan, F. (2000). Quantitative characterization of speech rhythm: syllable-timing in
53
54 Singapore English. *Language and Speech* 43(4), 377-401.
- 55
56
57 MacNeilage, P. (1998). The frame/content theory of evolution of speech production. *Behavioral and Brain*
58
59 *Sciences* 21(4), 499-511.
- 60

- 1 McAuley, J.D., & Fromboluti, E.K. (2014). Attentional entrainment and perceived event duration.
2
3 *Philosophical Transactions of the Royal Society B*, 369: 20130401.
4
- 5 Mehler, J. & Nespors, M. (2004) Linguistic rhythm and the development of language. In A. Belletti & L. Rizzi
6
7 (Eds.) *Structures and Beyond: The Cartography of Syntactic Structures*. Oxford: Oxford University Press.
8
9 Pages 213-221.
10
- 11 Mehler, J., Dehaene-Lambertz, G., Dupoux, E. & Nazzi, T. (1996) Coping with Linguistic Diversity: The
12
13 Infant's Viewpoint. In J.Morgan and K. Demuth (Eds.) *Signal to Syntax*. Hillsdale, NJ:LEA, 101-116..
14
- 15 Merker, B., Madison, G., and Eckerdal, P. (2009). On the role and origin of isochrony in human rhythmic
16
17 entrainment. *Cortex* 45, 4–17.
18
- 19 Molinaro, N., Lizarazu, M., Lallier, M., Bourguignon, M., & Carreiras, M. (2016) Out-of-synchrony speech
20
21 entrainment in developmental dyslexia. *Human Brain Mapping*, 37, 2767–2783.
22
- 23 Morey, R. D., & Rouder, J. N. (2012). BayesFactor: An R package for computing Bayes factor for a variety of
24
25 psychological research designs (available in the Comprehensive R Archive Network).
26
- 27 Muneaux, M., Ziegler, J.C., Truc, C., Thomson, J., & Goswami, U. (2004). Deficits in beat perception and
28
29 dyslexia: evidence from French. *Neuroreport* 15, 1255-1259.
30
- 31 Murty, L., Otake, T., & Cutler, A. (2007). Perceptual tests of rhythmic similarity: I. Mora Rhythm. *Language*
32
33 *and Speech*, 50(1), 77-99.
34
- 35 Nazzi, T., & Ramus, F. (2003). Perception and acquisition of linguistic rhythm by infants. *Speech*
36
37 *Communication* 41, 233-243.
38
- 39 Nazzi, T., Bertoncini, J., & Mehler, J. (1998). Language discrimination by newborns: Toward an
40
41 understanding of the role of rhythm. *Journal of Experimental Psychology: Human Perception and*
42
43 *Performance* 24(3), 756-766.
44
- 45 Ordin, M., & Polyanskaya, L. (2014). Development of timing patterns in first and second languages. *System*
46
47 42, 244-257
48
- 49 Ordin, M., & Polyanskaya, L. (2015). Acquisition of English speech rhythm by monolingual children.
50
51 *Proceedings of Interspeech 2015*, 3120-3124.
52
53
54
55
56
57
58
59
60

- 1 Ordin, M., & Polyanskaya, L. (2015a). Acquisition of speech rhythm in a second language by learners with
2 rhythmically different native languages. *The Journal of the Acoustical Society of America* 138(2), 533-544.
3
4 Ordin, M., & Polyanskaya, L. (2015b). Perception of speech rhythm in second language: The case of
5 rhythmically similar L1 and L2. *Frontiers in Psychology* 6: 316, 1-15.
6
7
8
9 Pallier, C. (2002). Computing discriminability and bias with the R software. Available at
10 <http://www.pallier.org/pdfs/aprime.pdf>. Last accessed on November 4, 2018.
11
12
13
14 Pamies Bertran, A., (1999). Prosodic Typology: On the Dichotomy between Stress-Timed and Syllable-
15 Timed Languages. *Language Design*, 2, 103-130.
16
17
18
19 Payne, E., Post, B., Astruc, L., Prieto, P., & del Mar Varnell, M. (2012). Measuring child rhythm. *Language*
20 *and Speech*, 55(2), 203-229.
21
22
23
24 Peelle, J. E., Gross, J., & Davis, M. H. (2013). Phrase-locked responses to speech in human auditory cortex
25 are enhanced during comprehension. *Cerebral Cortex* 23, 1378–1387.
26
27
28
29 Pike, K. L. (1945). *The Intonation of American English*. Ann Arbor. University of Michigan Press.
30
31
32 Polka, L., & Sundara, M. (2011). Word Segmentation in Monolingual Infants Acquiring Canadian English and
33 Canadian French: Native Language, Cross-Dialect, and Cross-Language Comparisons. *Infancy* 17, 198-232.
34
35
36
37 Polyanskaya, L., & Ordin, M. (2015). Acquisition of speech rhythm in first language. *The Journal of the*
38 *Acoustical Society of America* 138(3), 199-204.
39
40
41
42 Polyanskaya, L., Ordin, M., & Busa, M. (2017). Relative Saliency of Speech Rhythm and Speech Rate on
43 Perceived Foreign Accent in a Second Language. *Language and Speech* 60(3), 333-355.
44
45
46
47 Prieto, P., del Mar Vanrell, M., Astruc, L., Payne, E., & Post, B. (2012). Phonotactic and phrasal properties of
48 speech rhythm. Evidence from Catalan, English, and Spanish. *Speech Communication* 54, 681-702.
49
50
51
52 Ramus, F., & Mehler, J. (1999). Language identification with suprasegmental cues: A study based on speech
53 resynthesis. *Journal of the Acoustical Society of America* 105 (1), 512-521.
54
55
56
57 Ramus, F., Hauser, M. D., Miller, C., Morris, D., & Mehler, J. (2000). Language discrimination by human
58 newborns and by cotton-top tamarin monkeys. *Science* 288, 349-351.
59
60

- 1 Ramus, F., Nespors, M., & Mehler, J. (1999). Correlates of linguistic rhythm in the speech signal. *Cognition*
2
3 73(3), 265-292.
- 4
5 Roach P. (1982). On the distinction between 'stress-timed' and 'syllable-timed' languages. In D. Crystal
6
7 (Ed.). *Linguistic Controversies*. (Edward Arnold, London), 73-79.
- 8
9
10 Saffran, J. R., Newport, E. L., and Aslin, R. N. (1996). Word segmentation: The role of distributional cues.
11
12 *Journal of Memory and Language*, 35,606-621.
- 13
14
15 Schiering, R. (2007). The phonological basis of linguistic rhythm. Cross-linguistic data and diachronic
16
17 interpretation. *Sprachtypologie und Universalienforschung* 60, 337-359.
- 18
19
20 Smith, M. R., Cutler, A., Butterfield, S., & Nimmo-Smith, I. (1989). The perception of rhythm and word
21
22 boundaries in noise-masked speech. *Journal of Speech and Hearing Research* 32, 912-920.
- 23
24
25 Tajima, K., Port, R., & Dalby, J. (1997). Effects of temporal correction on intelligibility of foreign-accented
26
27 English. *Journal of Phonetics* 25, 1-24.
- 28
29
30 ten Oever, S., Schroeder, C. E., Poeppel, D., van Atteveldt, N., & Zion-Golumbic, E. (2014). Rhythmicity and
31
32 cross-modal temporal cues facilitate detection. *Neuropsychologia*, 63, 43–50.
- 33
34
35 Tincoff, R., Hauser, M., Tsao, F., Spaepen, G., Ramus, F., & Mehler, J. (2005). The role of speech rhythm in
36
37 language discrimination: further tests with a non-human primate. *Developmental Science* (1), 26-35.
- 38
39
40 Toro, J.M., & Nespors, M. (2015). Experience-dependent emergence of a grouping bias. *Biology Letters* 11.
- 41
42 Toro, J.M., Trobalon, J.B, Sebastian-Galle, N (2003) The use of prosodic cues in language discrimination
43
44 tasks by rats. *Animal Cognition* 6, 131–136.
- 45
46
47 Van Maastricht, L., Krahmer, E., Swerts, M., & Prieto, P. (in press). Learning direction matters. *Studies in*
48
49 *Second Language Acquisition* 1–35.
- 50
51
52 Werker, J., & Curtin, S. (2005) PRIMIR: A developmental framework of infant speech processing. *Language*
53
54 *Learning and Development* 1, 197-234.
- 55
56
57 White, L., & Mattys, S. (2007). Calibrating rhythm: First language and second language studies. *Journal of*
58
59 *Phonetics* 35(4), 501-522.
- 60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

White, L., Mattys, S., & Wiget, L. (2012). Language categorization by adults is based on sensitivity to durational cues, not rhythm class. *Journal of Memory and Language* 66(4), 665-679.

For Peer Review

FIGURE CAPTIONS

1
2
3 **Figure 1.** Sensitivity (a) and response bias (b) of German and French listeners for different stimuli pairs from
4
5 Experiment I. For example, „str-syl“ refers to performance in AX items, in which A was stress-timed and X was
6
7 syllable-timed. Horizontal line indicates the chance level (50%). Error bars $\pm 2SE$ around the mean.
8

9 **Figure 2.** RT for correct and incorrect answers given by German and French listeners. Error bars $\pm 2SE$ around the
10
11 mean.
12

13 **Figure 3.** RT for correct responses by German and French listeners for different pair-types. Error bars $\pm 2SE$ around
14
15 the mean.
16
17

18 **Figure 4.** Sensitivity (a) and response bias (b) of German and French listeners for different stimuli pairs from
19
20 Experiment II. For example, „str-syl“ refers to performance in AX items, in which A was stress-timed and X was
21
22 syllable-timed. Horizontal line indicates the chance level (50%). Error bars $\pm 2SE$ around the mean.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. Overview of the rhythm measures.

Metric name	Δ (delta)	Varco	%V	nPVI
Description	standard deviation in duration of speech intervals	Coefficient of variability in duration of speech intervals	Proportion of vocalic material in an utterance	Normalized pairwise variability index for speech intervals. $nPVI = \sum_{k=2}^n \frac{ d_k - d_{k-1} }{(d_k + d_{k-1})/2} / (n-1)$
Reference	Ramus, et al., 1999	Dellwo, 2006	Ramus, et al., 1999	Low, Grabe, & Nolan, 2002

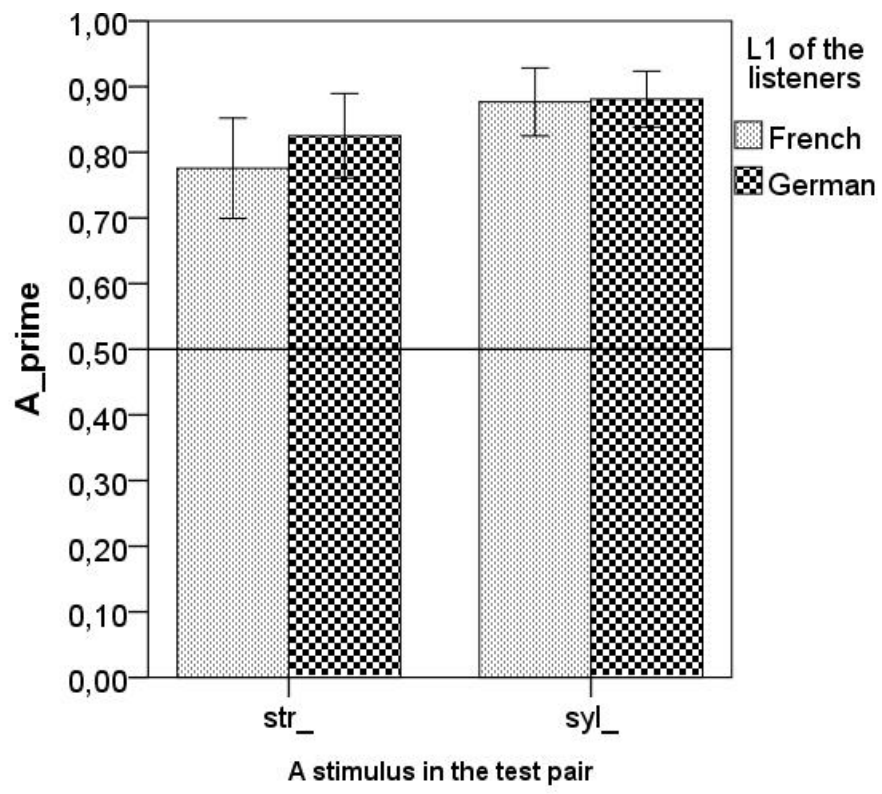
For Peer Review

Table 2. Measures of durational variability of Vs and mean durations of Cs (meanC), Vs (meanV), and the proportion of vocalic material in the stimuli with two contrastive rhythms

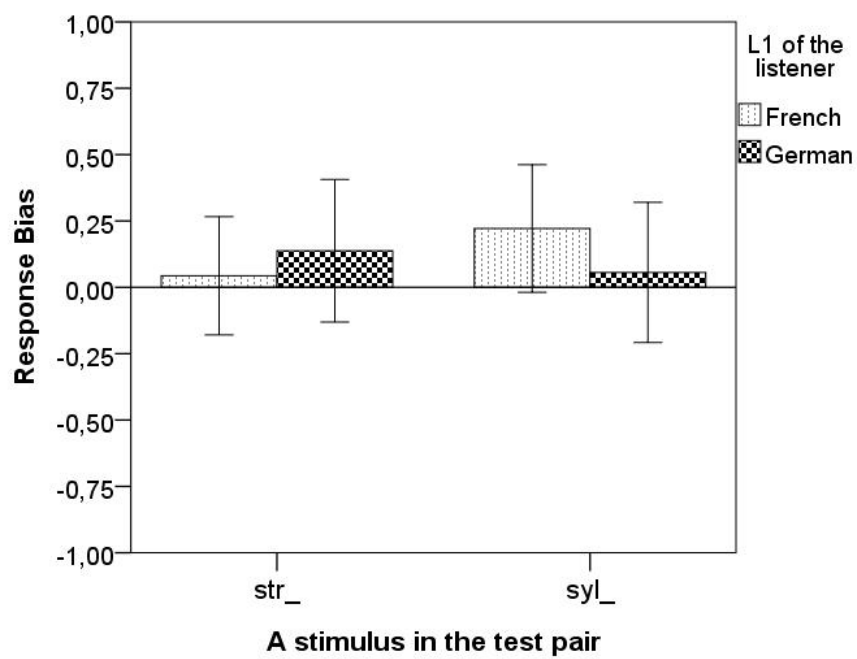
	%V	nPVI	ΔV	meanV	VarcoV	meanC
Stress-timed rhythm	40%	78	34.8	80ms	0.48	120ms
Syllable-timed rhythm	50%	7.7	10.7	100ms	0.107	100ms

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

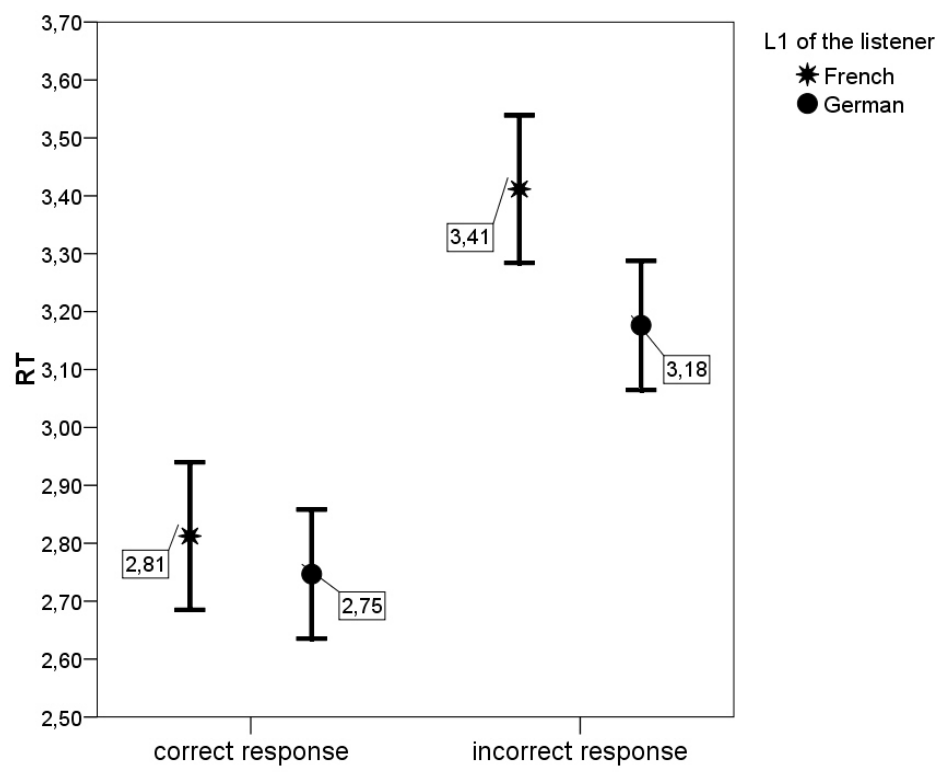


1a



1b

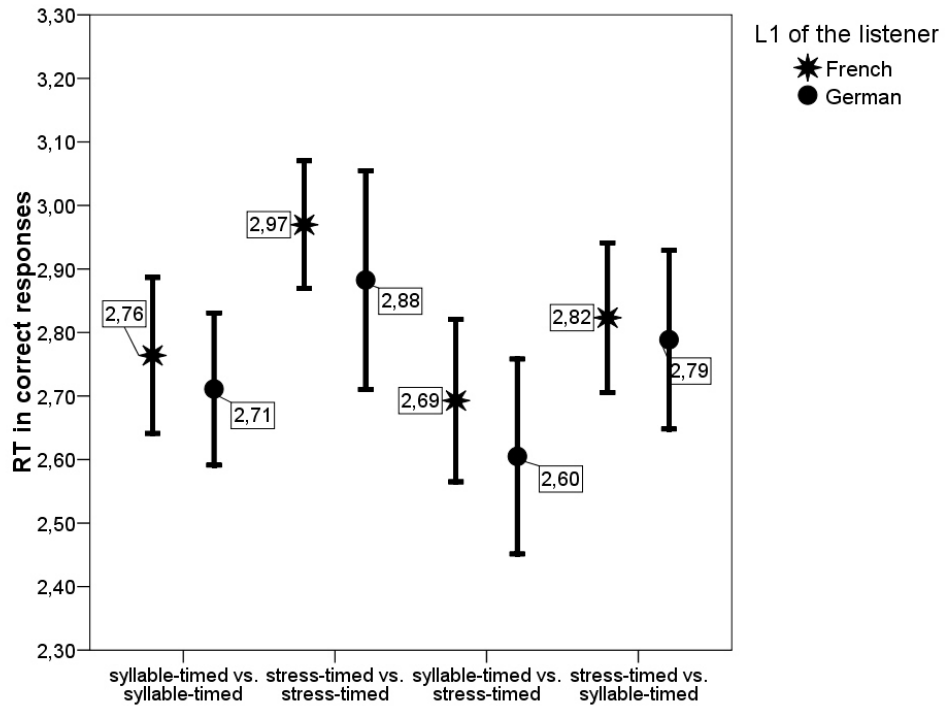
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



2

116x93mm (200 x 200 DPI)

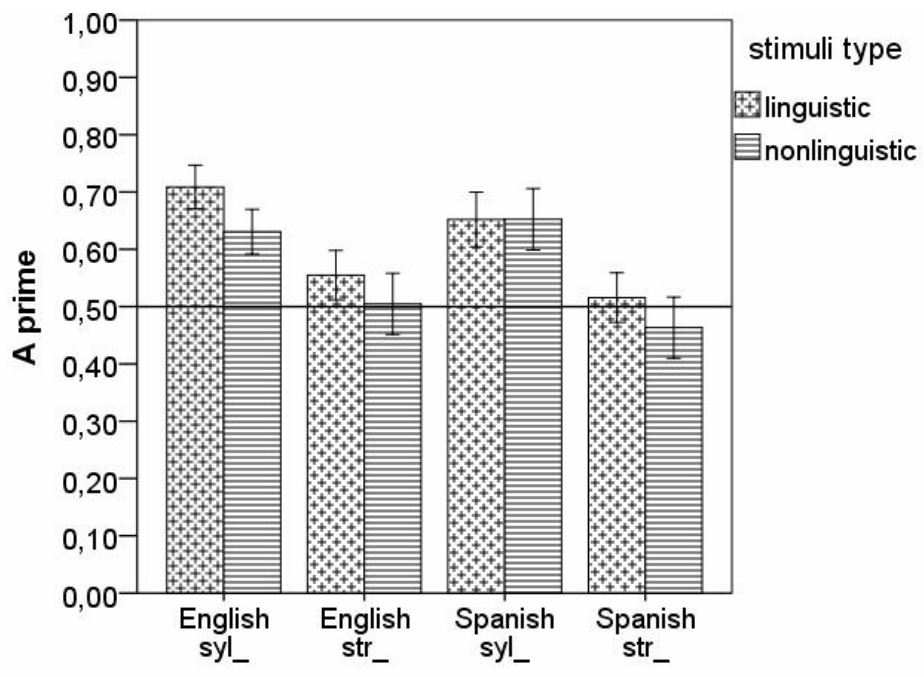
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



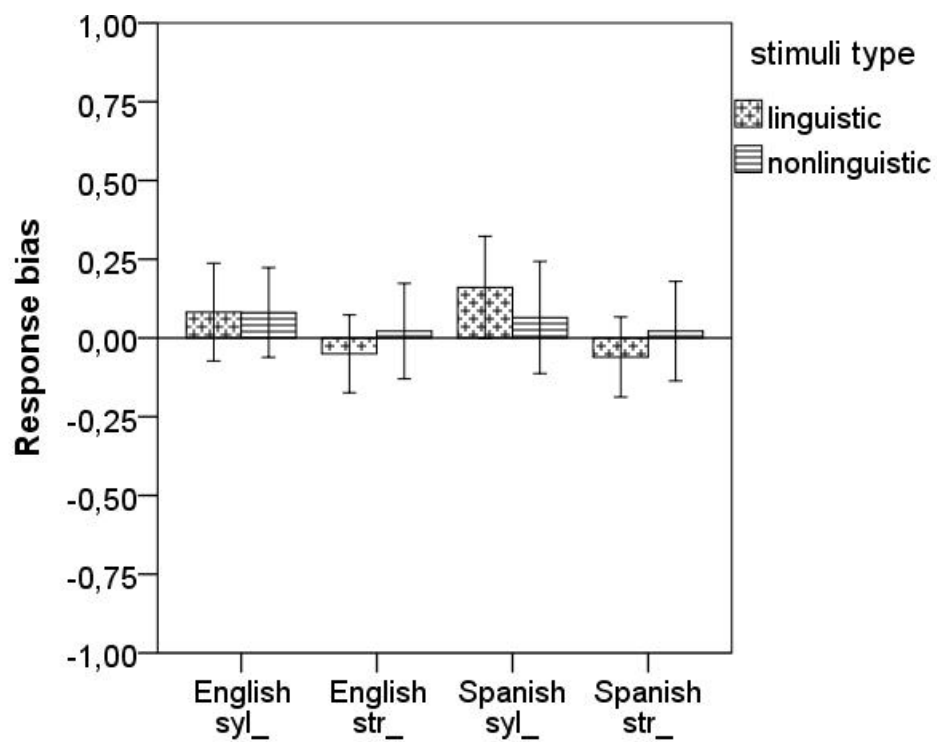
3

123x92mm (200 x 200 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



4a



4b

APPENDIX I:

To ensure that Welsh-based resynthesized sentences were perceived as linguistic stimuli for Experiment 2, we recruited 10 participants for a pilot assessment. We concatenated Welsh-based (i.e., linguistic) stimuli into two 2-minute sequences. Non-linguistic stimuli were also concatenated into two 2-minute sequences. Finally, we created an artificial language (following the approach often used in an artificial language learning paradigm, Aslin, Saffran, & Newport, 1998; Saffran, Newport, & Aslin, 1996). We used the same syllables with simple CV syllabic structure, which also comprised non-linguistic stimuli, and constructed 10 bi-syllabic nonsense words (samu, nelo, noma, namo, fenu, fale, lufe, mesu, sofu, sela). The vowel /i/ was only used in 'filler' syllables (fi, si, mi, ni, li) that were interspersed with the nonsense words and modelled frequent structural words (articles and prepositions). Lexical stress on word-initial syllables was modelled by lengthening the vowel by a factor of 1.5. In a stream FIMESUMISELALISAMUMIFALESINELO..., 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 segmenting continuous syllabic streams into constituents. An intonational contour was imposed on the syllabic stream, with boundary tones aligned with some word-final syllables. This contour allowed clustering the statistical nonsense words into larger constituents, e.g., sentences. In this way, a prosodic hierarchy with smaller discrete constituents embedded into larger units (Nespor & Vogel, 2007) was modelled. This implements into our artificial language a) a hierarchical structure, b) prosody (distribution of stressed and unstressed syllables, intonation); and c) statistical cues to mark discrete units and typological properties (head-prominence) of a language. In this way, we tried to make the stimuli sound similar to a real language. Such stimuli were also concatenated into two 2-minute sequences. The resulting 6 sequences of two minutes each were concatenated into a continuous stream, and we asked participants to listen to this stream. Prior to listening, participants were told that they were going to listen to some speech-like passages, some of which were a real natural language, while other passages were not.

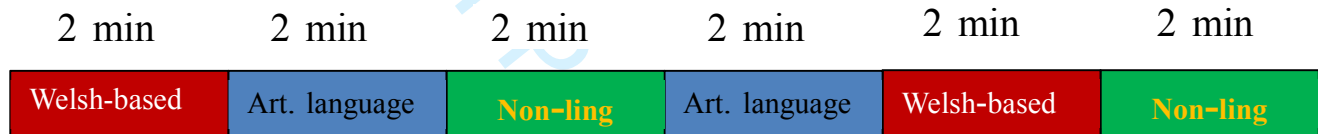


Figure AI-1. Schematic illustration of the familiarization stream used in the pilot experiment for evaluating the stimuli.

After familiarization, participants were played 120 stimuli, 6.6-seconds each, 40 stimuli of each type. Upon listening to each stimulus, they had to indicate how likely it is that the stimulus represents a real language. The responses were registered on an 8-point scale, from 1 – “I am sure it is a language”, to 8 – “I am sure it is not a language”. Figure AI-2 shows the results, error bars indicate 2 SE around the mean.

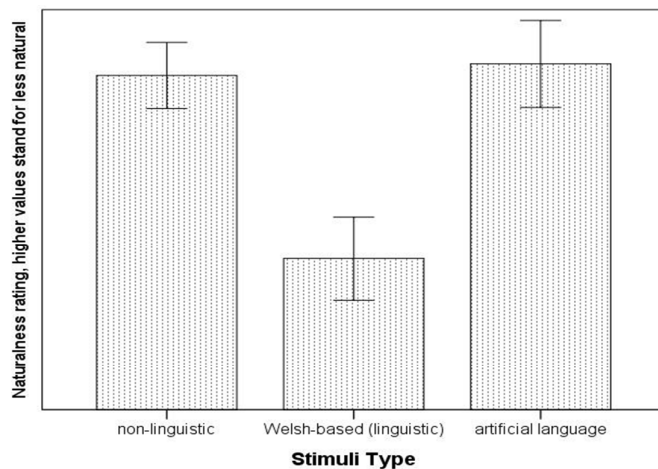


Figure AI-2. Differences in naturalness rating for different stimuli types.

The analysis revealed that stimulus type has a significant effect on the ratings, $\lambda=.908$, $F(2,398)=20.262$, $p<.0005$, $\eta^2=.092$. Pairwise comparisons (with Bonferroni correction applied) showed that Welsh-based stimuli were perceived to be significantly more likely to come from a real language than non-linguistic stimuli, $p<.0005$ for both contrasts. Ratings assigned to non-linguistic stimuli and to the artificial language did not differ, $p=.6$. Welsh-based stimuli were rated as significantly more natural and similar to real language utterances than the non-linguistic and artificial language stimuli. This suggests that, besides the presence of prosody and distributional cues, a larger variation in phonotactic and segmental complexity is needed for the sequence of syllables to be perceived as a real language. Based on these results, we chose Welsh-based resynthesized sentences as linguistic stimuli for Experiment 2.

APPENDIX II:

In order to ensure that linguistic stimuli indeed sound to participants as more representative of a real language, listeners were asked to perform a short test after the second session. For the test, we randomly selected 10 linguistic and 10 non-linguistic stimuli. Half of the stimuli exhibited regular rhythm (syllable-timed) and the other half had irregular rhythm (stress-timed). Participants listened to the stimuli in randomized order and responded to the question "How much does this sound like a real language". The naturalness ratings were given on an 8-point scale, from 1 – I am sure it is a language, to 8 – I am sure it is not a language. Figure AII-1 shows the results, error bars indicate 2 SE around the mean.

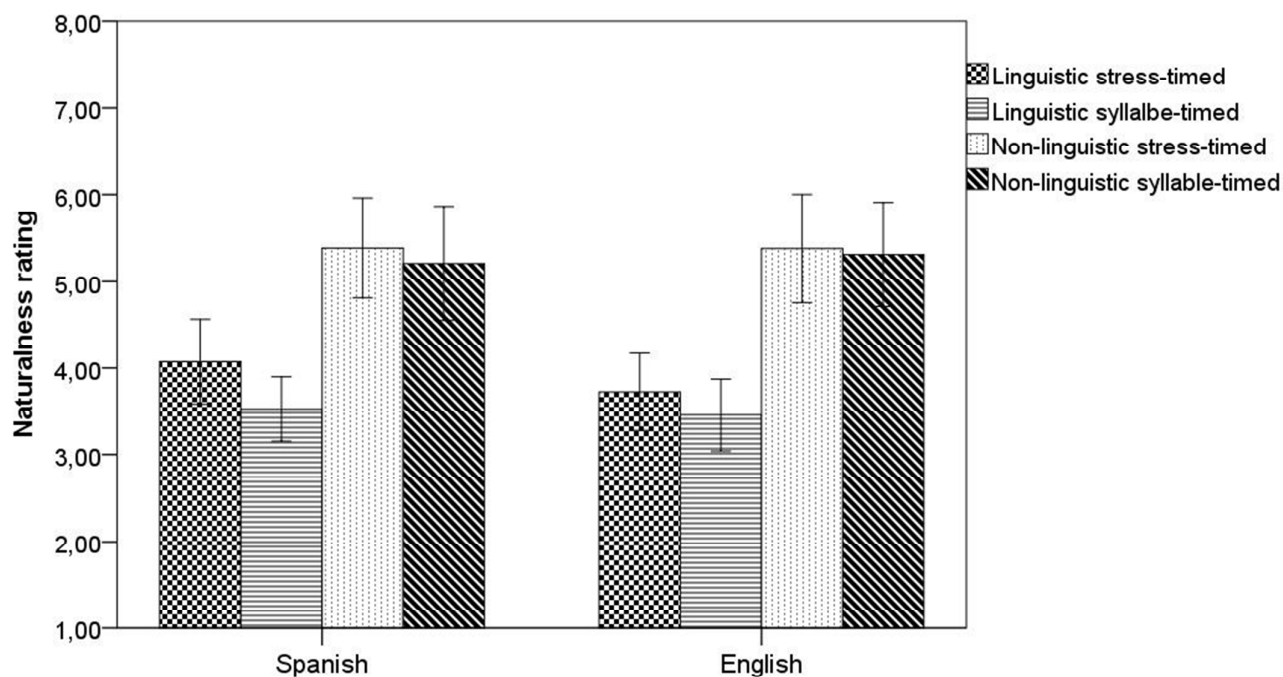


Figure AII-1. Naturalness ratings assigned by participants after the experiment to linguistic and non-linguistic stimuli with different rhythms.

An ANOVA with *L1* of the participant as a between-subject factor and *stimulus type* as a within-subject factor revealed a significant effect of stimulus type, $\lambda=.45$, $F(3,52)=21.19$, $p<.0005$, $\eta^2=.55$. There was no effect of *L1*, $F(1,54)=.089$, $p=.767$, $\eta^2=.002$ and no interaction between *L1* and stimulus type, $\lambda=.967$, $F(3,52)=.43$, $p=.73$, $\eta^2=.024$. Pairwise comparisons (with Bonferroni correction) showed that the ratings assigned to linguistic stimuli were significantly lower (thus rated as more similar to real language) than the ratings assigned to non-linguistic stimuli.

The results also showed that manipulations of duration of the linguistic stimuli did not affect naturalness ratings. Within each stimulus type (linguistic or non-linguistic), the naturalness ratings assigned to the stimuli with irregular and regular rhythmic patterns (i.e., stress- and syllable-timed correspondingly) did not statistically differ. These results confirm that linguistic stimuli were indeed perceived as more representative of a real language and were more likely to engage a fuller set of speech processing mechanisms.