

# Body Height Perception from Connected Speech: Raising the Fundamental Frequency Increases Accuracy

## ABSTRACT

Human listeners can perceive the speaker's body size from its voice to a certain degree. The voice pitch or speaking fundamental frequency (Fo) and the vocal formant frequencies are the voice parameters most intensively studied in research on the perception of body size (particularly height). Pisanki et al. (2014) found that artificially lowering Fo of isolated vowels pronounced by male speakers improved the perception of body height by part of listeners. The authors explained this effect as a denser harmonic spectrum provided by low pitch allowed for better resolution of formants, aiding formant-based size assessment.

In the present study, we tested and extended this research using connected speech (sentences, words) pronounced by speakers of both sexes. Unexpectedly, our results showed that raising Fo, not lowering it, increases the performance in two binary discrimination task of body size. This new finding is explained in the temporal domain by the dynamic and time-varying acoustic properties of connected speech. One possible reason is that an increase of Fo might increase the sampling density of the wave acoustic cycles and provide more detailed information, such as higher resolution, on the envelope shape.

**KEYWORDS:** speech perception, speaker body size, height, fundamental frequency, pitch.

1  
2  
3  
4  
5  
6 Human voice is a multidimensional signal that conveys a great amount of information about  
7  
8 the speaker: his/her sex, age, individual identity, emotional state and physical characteristics,  
9  
10 such as her/his body size (height and weight). Among several voice parameters, both vocal  
11  
12 formants and speaking fundamental frequency (pitch) predict the considerable variation in body  
13  
14 size between humans of different ages and sexes (Peterson & Barney, 1952; Titze, 1989).  
15  
16 However, when individual differences in body size are controlled by age and sex variables (e.g.,  
17  
18 adults of the same sex), only formants reliably predict speaker's size (González, 2004, 2006;  
19  
20 Pisanski, Fraccaro, Tigue, O'Connor, & Feinberg, 2014a; Pisanski et al., 2014b). The formants  
21  
22 are the resonant frequencies of the vocal tract and are dependent on the size of the vocal tube  
23  
24 (Fant, 1960). As the vocal tract length (VTL) is larger, the formant frequencies are lower, and  
25  
26 vice versa. Vocal tract is constrained to some extent by anatomical structures related to body size  
27  
28 (e.g., individual's skull and body height). Consequently, a negative correlation exists between  
29  
30 formant frequencies and speaker's body size, although its value is moderate in general terms  
31  
32 (Bruckert, Liénard, Lacroix, Kreutzer, & Leboucher, 2006; González, 2004, 2006; Pisanski et al.,  
33  
34 2014b). Conversely, the correlation between pitch, or speaking fundamental frequency (Fo), and  
35  
36 body size among adults of the same sex is very weak, or virtually null (Bruckert et al., 2006;  
37  
38 González, 2007; Künzel, 1989; Pisanski et al., 2014b; van Dommelen & Moxness, 1995). Pitch  
39  
40 is produced by vocal folds, which are composed of a soft tissue not subject to the same  
41  
42 anatomical constraints than vocal formants.  
43  
44  
45  
46  
47  
48

49 Listeners can roughly infer the speaker's body size (specifically the body height) based  
50  
51 (correctly) on the formant frequencies of his/her voice. Unexpectedly, research has shown that  
52  
53 listeners incorrectly infer the speaker's height among adults of the same sex based on the pitch or  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 fundamental frequency (Fo) of her/his voice. Listeners consistently associate low speaking Fo  
4 with tall people, and, conversely, high Fo with short people (Feinberg, Jones, Little, Burt, &  
5 Perrett, 2005; Pisanski & Rendall, 2011; Pisanski et al., 2014b; van Dommelen & Moxness,  
6 1995). Presumably, this association could be the result of an erroneous overgeneralization of the  
7 physical principle that large or long objects produce low frequency vibrations and small or short  
8 objects give rise to higher vibrations. This extrapolation is wrong because, as mentioned above,  
9 the size of the vocal cords is hardly related to the body size among adults of the same sex.

10  
11  
12  
13  
14  
15  
16  
17  
18  
19 Given the absence of an objective link between voice pitch and body size at the within-sex  
20 level, pitch has traditionally been hypothesized that its presence in the human voice could make  
21 it difficult to correctly perceive the speaker's body size from the vocal formant frequencies.  
22  
23  
24 Pisanski, et al. (2014a) performed a direct test of how the accuracy of listeners' judgments is  
25 affected when voice pitch is present or absent from the acoustical signal. Counterintuitively, they  
26 found the listeners' accuracy increased in the presence rather than absence of voice pitch in the  
27 experimental stimuli. More specifically, height estimations based on relatively low voice pitch  
28 achieved higher accuracy. These authors hypothesized that the denser harmonic spectrum  
29 provided by low pitch allowed for better resolution of formants, aiding listeners in extracting  
30 reliable formant-based information from the voice. They tested this hypothesis (Exp. 3) and  
31 verified that artificially lowering the fundamental frequency of voice with the help of technology  
32 facilitated assessments of men's body height.

33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47 Pisanski, et al. (2014a, Exp. 3) employed the five Canadian English monophthong vowels  
48 (/a/, /i/, /ε/, /o/, and /u/) recorded from thirty men as experimental stimuli, and subsequently,  
49 their fundamental frequency was manipulated using the PSOLA algorithm of Praat software  
50 (Boersma & Weenink, 2013) maintaining formants constant. The aim of the present work was to  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 extend this research in relation to several relevant aspects: a) using connected speech (sentences,  
4 words) as experimental study, beyond isolated vowels; b) using voices of men and women; and  
5  
6 c) using materials from another linguistic context (Spanish).  
7  
8  
9  
10

## 11 12 **EXPERIMENT 1**

13  
14 The purpose of Experiment 1 was to observe the effect of lowering or raising the  
15 fundamental frequency of a sentence spoken by men and women in assessments of speakers'  
16 body size.  
17  
18  
19  
20

### 21 **METHOD**

#### 22 **Participants.**

23  
24 Participants as listeners were 68 young adults of both sexes (57 females) whose age range  
25 was 18–30 years ( $M = 19.67$ ;  $SD = 2.52$ ). All participants were undergraduate students at the  
26 University Jaume I (Spain), who participated voluntarily with informed consent and were  
27 compensated with course credit.  
28  
29  
30  
31  
32  
33  
34

#### 35 **Materials**

36  
37 The stimuli consisted of a Spanish interrogative sentence ('¿Cuántos años tiene tu primo de  
38 Barcelona?' [How old is your cousin from Barcelona?]) recorded from 40 young adult speakers,  
39  
40 20 men and 20 women, who had been university students several years ago, unknown to the  
41 participants. Voices were selected from a pool of voice recordings (González & Oliver, 2005),  
42 whose speakers had their height measured with a metric tape. The range of heights of the  
43 selected male speakers was 160-189 cm, with an average  $M = 176.9$  cm ( $\pm SD = 6.3$  cm), close  
44 to the general population of Spanish men (176.6 cm; NCDRISC, 2018). The range of heights of  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 the selected female speakers was 154-175 cm, with an average  $M = 163.4$  cm ( $\pm$  SD = 7.5 cm),  
4  
5 very close to the general population of Spanish women (163.4 cm; NCDRISC, 2018).  
6

7  
8 The sentence was recorded in a sound-attenuated booth with a Shure SM58 microphone at  
9  
10 an approximate distance of 12 cm from the mouth, and a Sony-TCD D-8 digital audiotape (DAT)  
11  
12 recorder with a sample frequency of 44.1 kHz. Then, the speech signal was digitally transferred  
13  
14 to a PC computer and converted to 16-bit WAV files. Finally, the files were equated in RMS  
15  
16 (root mean square) amplitude.  
17

18  
19 *Manipulation of the Fundamental Frequency (Fo).* We used three sets of stimuli: the natural  
20  
21 speech recordings above mentioned ('Normal'), the natural recordings with their Fo lowered  
22  
23 ('Low'), and the natural recordings with their Fo raised ('High'). Following the same procedure  
24  
25 as Pisanski, et al. (2014a, Exp. 3), 'Low' stimuli were the natural recordings manipulated by  
26  
27 means of PSOLA algorithm of Praat software (Boersma & Weenink, 2016) to subtract 0.5 ERBs  
28  
29 of the baseline Fo, maintaining formants constant. 'High' stimuli were obtained in the same way  
30  
31 adding 0.5 ERBs to the baseline Fo. PSOLA (Pitch Synchronous Overlap and Add) is a digital  
32  
33 signal processing technique initially proposed by Moulines & Charpentier (1990) and used for  
34  
35 manipulating the pitch (or time) of an acoustic speech signal, without altering other parameters.  
36  
37 PSOLA decomposes the speech waveform in small overlapping segments, and moves them  
38  
39 further apart to decrease Fo, or closer together to increase Fo. ERB or equivalent rectangular  
40  
41 bandwidth is a measure used in psychoacoustics, which gives an approximation of the perceived  
42  
43 voice pitch. For example, adding or subtracting 0.5 ERB to 120 Hz is roughly equivalent to  
44  
45 adding or subtracting 20 Hz.  
46  
47  
48  
49

50  
51 *Pairing of speech stimuli.* We created a total of 240 pairs of stimuli, 40 for each Fo  
52  
53 condition (Low, Normal, High) x speaker's sex (males, females). Each pair was formed by two  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 different stimuli belonging to the same Fo condition and speaker's sex, separated by 800 ms of  
4  
5 silence. Normal/male condition was formed by 40 different pairs of male speakers, whose  
6  
7 heights differed in a range of 5 – 19 cm ( $M = 9.3$  cm,  $SD = 3.9$  cm); half of the pairs included the  
8  
9 taller man in the first place, and half in the second place. Normal/female condition was formed  
10  
11 by 40 different pairs of female speakers, whose heights differed in a range of 9 – 17 cm ( $M =$   
12  
13  $13.9$  cm,  $SD = 3.0$  cm); half of the pairs included the taller woman in the first place, and half in  
14  
15 the second place. The pairs were the same for the other two Fo conditions, Low and High.  
16  
17

18  
19 The total of 240 pairs was divided in four sets of 60 pairs (10 per each Fo x speaker sex  
20  
21 condition), ensuring that neither pair of speakers was repeated through the Fo conditions within  
22  
23 the same set.  
24  
25

### 26 **Procedure**

27  
28 Each participant was randomly assigned to one of the four sets of stimuli, resulting in a total  
29  
30 of 60 size assessment trials per each participant. Participants individually performed the  
31  
32 experiment in six short sessions of ten stimuli of the same Fo/speaker sex condition. The order of  
33  
34 the sessions was randomized through the participants. At the beginning of each session,  
35  
36 participants were informed if the speakers were men or women.  
37  
38

39  
40 On each trial, participants were presented through headphones with two speaker's voices of  
41  
42 the same Fo condition (Low, Normal, High) and speaker's sex (males, females). Voices were  
43  
44 played consecutively and separated by 800 ms of silence. After listening to the pair of voices,  
45  
46 participants were asked to indicate which of the two voices belonged to the taller speaker by  
47  
48 selecting the corresponding button on the screen. On each trial, participants could listen to the  
49  
50 pair of voices again at will by clicking on a play button.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## RESULTS AND DISCUSSION

Contrary to expectations, it was raising Fo, rather than lowering Fo, which allowed increasing accuracy in estimating the speaker's body size. The overall percentages of successes of trials in which the taller speaker was correctly identified were the following for each Fo condition. For the 'Normal' condition (in which Fo was not modified), accuracy achieved 59.0%, 95% CI [55.5%, 62.6%], above chance level (50.0%;  $p < .0001$ ). For the Low condition (Fo lowered - 0.5 ERBs), accuracy was 57.0% [53.5%, 60.5%], whereas for High condition (Fo raised + 0.5 ERBs) accuracy increased to 63.6% [60.5%, 66.6%]. Table 1 and Figure 1 shows the percentages separated by the sex of speakers. We can see the same trend for both sexes: accuracy increased for high-Fo sentences and roughly decreased for low-Fo sentences.

We carried out a 3 (Fo; Low, Normal, High) x 2 (Sex of Speakers) ANOVA. The analysis found a significant Fo effect,  $F(2, 134) = 7.24$ ,  $MS_e = 214.12$ ,  $p = .001$ ,  $\eta^2_p = .098$ , and also a significant effect of the Sex of Speaker factor,  $F(1, 67) = 9.54$ ,  $MS_e = 230.21$ ,  $p = .003$ ,  $\eta^2_p = .125$ , indicating that the female speakers received higher rates (62.2%, 95% CI [58.8%, 65.6%]) than the male speakers (57.5%, 95% CI [54.9%, 60.2%]). The Fo x Sex of Speaker interaction did not reach significance,  $F(2, 134) < 1$ . Posterior within-subject contrasts revealed that accuracy for the High Fo condition was significantly different (higher) from accuracy for the Normal condition,  $F(1, 67) = 14.01$ ,  $MS_e = 301.86$ ,  $p < .0001$ ,  $\eta^2_p = .173$ .

Our data clearly suggested that by raising the fundamental frequency of the voice, listeners increased their accuracy in judging the speakers' body size. To examine this hypothesis, we performed a second experiment applying various degrees of Fo alteration and extending the research with stimuli consisted of words.

-----  
Please, Table 1 and Figure 1 about here  
-----

## EXPERIMENT 2

The aim of Experiment 2 was to study the effect of lowering or raising two degrees ( $\pm 0.5$ ,  $\pm 1$  ERBs) the pitch or speaking fundamental frequency (Fo) of a word in the perception of the speaker's body size by part of listeners. The use of various degrees of Fo multiplies the number of stimulus sets, so the experiment focalized on a sex of the speakers, specifically women, as Pisanski et al. 2014a employed only male speakers.

### METHOD

#### Participants.

Participants as listeners were 64 young adults of both sexes (50 females) whose age range was 18–32 years ( $M = 19.56$ ;  $SD = 2.67$ ). All participants were undergraduate students at the University Jaume I (Spain), who participated voluntarily with informed consent and were compensated with course credit. None of the participants had participated in Experiment 1.

#### Materials

The stimuli consisted of a Spanish word ('importante' [important]) recorded from 20 young female speakers who had been university students several years ago, unknown to the participants. The speakers were the same as the female speakers in experiment 1. The words had been recorded and processed under the same conditions as the sentences from the experiment 1.



1  
2  
3        *Manipulation of the Fundamental Frequency (Fo).* We used five Fo conditions: the natural  
4 speech recordings ('Normal' condition), the natural recordings with their Fo lowered - 0.5 ERB  
5 ('- 0.5 ERB'), the natural recordings with their Fo lowered - 1.0 ERB ('- 1.0 ERB'), the natural  
6 recordings with their Fo raised + 0.5 ERB ('+ 0.5 ERB'), and the natural recordings with their Fo  
7 raised + 1.0 ERB ('+ 1.0 ERB'). The manipulation of Fo was carried out with the same technical  
8 procedure used in experiment 1 (PSOLA algorithm of Praat software, Boersma & Weenink,  
9 2016).

10  
11        *Pairing of speech stimuli.* We created a total of 150 pairs of stimuli, 30 for each Fo  
12 condition (- 1.0 ERB, - 0.5 ERB, Normal, + 0.5 ERB, + 1.0 ERB). Each pair was formed by two  
13 different stimuli belonging to the same Fo condition, separated by 800 ms of silence. Normal  
14 condition was formed by 30 different pairs of female speakers, whose heights differed in a range  
15 of 9 – 19 cm (M = 13.9 cm, SD = 2.9 cm); half of the pairs included the taller woman in the first  
16 place, and half in the second place. The pairs were the same for the other four Fo conditions.

### 32        **Procedure**

33  
34        The total of 150 pairs was divided into three sets of 50 pairs (10 per each Fo condition).  
35 Each participant was randomly assigned to one of the three sets of stimuli, resulting in a total of  
36 50 size assessment trials per each participant. Participants individually performed the experiment  
37 in two short sessions of twenty-five trials in random order.

38  
39        On each trial, participants were presented through headphones with two women's voices of  
40 the same Fo condition. Voices were played consecutively and separated by 800 ms of silence.  
41 After listening to the pair of voices, participants were asked to indicate which of the two voices  
42 belonged to the taller woman by selecting the corresponding button on the screen. On each trial,  
43 participants could listen to the pair of voices again at will by clicking on a play button.

## RESULTS AND DISCUSSION

Table 1 and Figure 2 show the percentages of successes for each Fo condition. We found again that raising the pitch or Fo of an utterance (spoken word) tends to increase the accuracy of body size assessments, whereas lowering Fo tends to show the opposite effect, decreasing the accuracy in judging the body size, even though the effect size was smaller in this experiment than in the first experiment, and the consequences of Fo's manipulation become more evident in the most extreme degrees (- 1 ERB, + 1 ERB).

We performed a 5 (Fo condition; - 1.0 ERB, - 0.5 ERB, Normal, + 0.5 ERB, + 1.0 ERB) ANOVA. The analysis showed a marginal main effect of Fo,  $F(4, 252) = 2.17$ ,  $MS_e = 161.72$ ,  $p = .073$ ,  $\eta^2_p = .033$ , and the polynomial test of within-subject contrasts was significant,  $F(1, 63) = 6.87$ ,  $MS_e = 193.75$ ,  $p = .011$ ,  $\eta^2_p = .098$ . Posterior within-subject contrasts found significant differences between the - 1.0 ERB condition and the mean, and between the - 0.5 ERB vs. - 1.0 ERB conditions.

We performed an ANOVA with only the extreme Fo conditions (- 1.0 ERB, Normal, + 1.0 ERB). The analysis yielded a significant main effect of Fo,  $F(2, 126) = 3.47$ ,  $MS_e = 183.45$ ,  $p = .034$ ,  $\eta^2_p = .052$ . A posterior within-subject contrast found a significant difference between the Normal vs. + 1.0 ERB conditions.

-----  
Please, Figure 2 about here  
-----

## GENERAL DISCUSSION

The initial goal of the present study was to verify and extend the finding obtained by Pisanski et al. 2014a with isolated vowels to connected speech (sentences and words). Pisanski et al. (2014a) found in their third experiment that lowering the voice pitch or fundamental frequency ( $F_0$ ) of a vowel helped listeners improve their accuracy in estimating the body size of the utterers (specifically male speakers). The authors explained this effect by means of the *Harmonic density hypothesis*. More specifically, increasing the density of harmonics in the speech spectrum (a consequence of lowering the fundamental frequency, because harmonics are multiples of  $F_0$ ) would allow for better resolution of formants, aiding formant-based body size perception.

We tested this hypothesis employing connected speech as stimuli instead of isolated vowels. In the first experiment, we increased and decreased the  $F_0$  of a sentence pronounced by a set of male and female speakers. Unexpectedly, listeners perceived better the speaker's body size (height) when  $F_0$  was raised (not lowered), contradicting to the prediction made by the Harmonic density hypothesis. This effect was roughly corroborated in a second experiment based on a word spoken by female speakers, although the effect size was smaller. Apparently, this finding was counterintuitive and contrary to the prediction of the mentioned hypothesis.

One possible explanation could be in the temporal domain of stimuli rather than in the frequency domain. González and Oliver (2004) compared size estimations from connected speech (sentences) versus sustained vowels and found that connected speech yielded more accurate judgments than vowels. Logically, connected speech comprises a much larger voice sample than isolated vowels. However, connected speech is more than the sum of different phonemes and constitutes a dynamic and coarticulated succession of consonants and vowels

1  
2  
3 forming syllables as units of articulation (particularly in Spanish). There may be other acoustic  
4 predictors of body size beyond the mean frequencies of  $F_0$  and the formants ( $F_1$ - $F_4$ ). González  
5 and Oliver (2004) found that actual height was the strongest predictor of perceptual judgments  
6 from a spoken sentence, in both male and female speakers. This was more to be expected in  
7 female voices, given the low correlation between the acoustic parameters studied ( $F_0$ ,  $F_1$ - $F_4$ )  
8 and height, but it was also evident in male voices, in which the partial correlation between actual  
9 height and judgments of "taller" in a binary discrimination task reached 0.51 after removing the  
10 influence of  $F_0$  and  $F_1$ - $F_4$ . These authors concluded that "partial correlations and multiple  
11 regression analysis showed that listeners efficiently rely on other features present mainly in  
12 connected speech (sentences as opposed to a sustained vowel)" (p. 29). They hypothesized that  
13 such features could be related to dynamic and time-varying acoustic properties of speech that  
14 could serve as sources of inferences about the biomechanical properties – mass, mobility and  
15 inertial characteristics – of the articulator apparatus and its possible relationship with body size.  
16 In this sense, the envelope of the acoustic wave of speech in the temporal domain could play a  
17 role in providing information on such dynamic and time-varying properties. In this case, an  
18 increase in the fundamental frequency would increase the sampling density of the wave acoustic  
19 cycles and provide more detailed information, such as higher resolution, on the envelope shape.

20  
21  
22 Irino, Aoki, Kawahara, and Patterson (2012) studied the interaction between  $F_0$  (glottal  
23 pulse rate, GPR) and mean-formant-frequency (MFF) with synthesized words in the perception  
24 of speaker characteristics and speech recognition. They observed that, within a certain range of  
25 MFF, increasing the  $F_0$  had the effect of improving the performance. The authors interpreted this  
26 result in the temporal domain of speech signal because as the glottal pulse rate ( $F_0$ ) increases, the  
27 number of cycles per unit time also increases. They related this physical fact to the functioning

1  
2  
3 of the auditory system. “If the analysis performed by the auditory system is pitch synchronous  
4 and anchored to the glottal pulses as described in Irino and Patterson (2002), the internal  
5 representation will become much more stable as GPR increases from 0.25 to 4.0, and it seems  
6 reasonable to assume that speech features would be better defined in a more stable  
7 representation” (Irino et al., 2012, p. 1009). Irino and Patterson (2002) considered that human  
8 listeners hear, for example, vowels spoken by men and women as approximately the same vowel  
9 although the length of the vocal tract varies considerably from group to group. Simultaneously,  
10 listeners can identify the speaker group. They suggested that the auditory system can extract and  
11 segregate information about the size of the vocal tract from information about its shape. In  
12 accordance with Irino and Patterson (2002), the auditory system performs a transform (Mellin  
13 transform) creating an auditory image, in which the information on the vocal tract size is  
14 represented apart from the information corresponding to its particular shape. According with  
15 Irino et al. (2012), this auditory image will become more stable and efficient when  $F_0$  increases  
16 within a certain range.

17  
18  
19 Nevertheless, further research will be necessary to confirm the present finding and test the  
20 effect that lowering versus raising the speaking fundamental frequency could have on perception  
21 of speaker’s body size (and vocal tract length) using different speech materials beyond isolated  
22 vowels, and compare voices of male versus female speakers, along with children.  
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

## Footnotes

<sup>1</sup> The effect size interpretations for  $\eta^2_p$  values are: .01 = small, .06 = medium, and .14 = large.

## Ethical Compliance Section

This work was completed with resources provided by the University Jaume I (Spain).

All procedures performed in studies involving human participants were in accordance with the Deontological commission and of the Ethical Committee of Animal Welfare of the University Jaume I (Spain) and with the 1964 Helsinki Declaration and its later amendments or comparable ethical standards.

The authors declare they have no conflict of interest.

Written informed consent was obtained from all individual adult participants included in the study.

## References

- Boersma, P., & Weenink, D. (2013). Praat: Doing phonetics by computer (Version 5.2.15). [Software]. Retrieved from <http://www.praat.org>.
- Boersma, P., & Weenink, D. (2016). Praat: Doing phonetics by computer (Version 6.0.19). [Software]. Retrieved from <http://www.praat.org>.
- Bruckert, L., Liénard, J. S., Lacroix, A., Kreutzer, M., & Leboucher, G. (2006). Women use voice parameters to assess men's characteristics. *Proceedings of the Royal Society B: Biological Sciences*, 273(1582), 83-89.
- Fant, G. (1960). *Acoustic theory of speech production*. The Hague: Mouton.
- Feinberg, D. R., Jones, B. C., Little, A. C., Burt, D. M., & Perrett, D. I. (2005). Manipulations of fundamental and formant frequencies influence the attractiveness of human male voices. *Animal behaviour*, 69(3), 561-568.
- González, J. (2004). Formant frequencies and body size of speaker: A weak relationship in adult humans. *Journal of Phonetics*, 32, 277-287.
- González, J. (2006). Research in acoustics of human speech sounds: Correlates and perception of speaker body size. *Recent Research Developments in Applied Physics*, 9, 1-15.
- González, J. (2007). Correlations between speakers' body size and acoustic parameters of voice. *Perceptual and motor skills*, 105(1), 215-220.
- González, J., & Oliver, J. C. (2004). Percepción a través de la voz de las características físicas del hablante: identificación de la estatura a partir de una frase o una vocal [Perception through the voice of speaker physical characteristics: identification of height from a sentence or a vowel]. *Revista de psicología general y aplicada*, 57(1), 21-34.

1  
2  
3 Gonzalez, J., & Oliver, J. C. (2005). Gender and speaker identification as a function of the  
4 number of channels in spectrally reduced speech. *The Journal of the Acoustical Society of*  
5  
6 *America*, 118(1), 461-470.  
7  
8

9  
10 Irino, T., Aoki, Y., Kawahara, H., & Patterson, R. D. (2012). Comparison of performance  
11 with voiced and whispered speech in word recognition and mean-formant-frequency  
12 discrimination. *Speech Communication*, 54(9), 998-1013.  
13  
14

15  
16 Irino, T., & Patterson, R. D. (2002). Segregating information about the size and shape of the  
17 vocal tract using a time-domain auditory model: The stabilised wavelet-Mellin transform. *Speech*  
18 *Communication*, 36(3-4), 181-203.  
19  
20  
21

22  
23 Künzel, H. J. (1989). How well does average fundamental frequency correlate with speaker  
24 height and weight? *Phonetica*, 46(1-3), 117-125.  
25  
26

27  
28 NCDRISC (2018). NCD Risk Factor Collaboration (NCD-RisC).  
29  
30 <https://www.ncdrisc.org/height-mean-map.html>.  
31  
32

33  
34 Moulines, E., & Charpentier, F. (1990). Pitch-synchronous waveform processing techniques  
35 for text-to-speech synthesis using diphones. *Speech communication*, 9(5-6), 453-467.  
36  
37

38  
39 Peterson, G. E., & Barney, H. L. (1952). Control methods used in a study of the vowels. *The*  
40 *Journal of the acoustical society of America*, 24(2), 175-184.  
41  
42

43  
44 Pisanski, K., Fraccaro, P. J., Tigue, C. C., O'Connor, J. J., & Feinberg, D. R. (2014a). Return  
45 to Oz: Voice pitch facilitates assessments of men's body size. *Journal of Experimental*  
46 *Psychology: Human Perception and Performance*, 40(4), 1316.  
47  
48

49  
50 Pisanski, K., Fraccaro, P. J., Tigue, C. C., O'Connor, J. J., Röder, S., Andrews, P. W., ... &  
51 Feinberg, D. R. (2014b). Vocal indicators of body size in men and women: a meta-analysis.  
52  
53 *Animal Behaviour*, 95, 89-99.  
54  
55



1  
2  
3           Pisanski, K., & Rendall, D. (2011). The prioritization of voice fundamental frequency or  
4 formants in listeners' assessments of speaker size, masculinity, and attractiveness. *The Journal of*  
5 *the Acoustical Society of America*, 129(4), 2201-2212.  
6  
7

8  
9  
10           Titze, I. R. (1989). Physiologic and acoustic differences between male and female voices.  
11 *The Journal of the Acoustical Society of America*, 85(4), 1699-1707.  
12  
13

14           van Dommelen, W. A., & Moxness, B. H. (1995). Acoustic parameters in speaker height and  
15 weight identification: sex-specific behaviour. *Language and speech*, 38(3), 267-287.  
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

## FIGURE CAPTIONS

Figure 1. Experiment 1 (sentences). Percentages of successes in a binary discrimination task on the speakers' body height, separated by experimental conditions according to the manipulation of  $F_0$  and the speaker sex.

Figure 2. Experiment 2 (words). Percentages of successes in a binary discrimination task on the female speakers' height, separated by experimental conditions according to the manipulation of  $F_0$ . ERB: Equivalent rectangular bandwidth.

## TABLES

Table 1. Data from Experiments 1 (sentences) and 2 (words). Percentages of successes in a binary discrimination task on the speakers' body height, separated by experimental conditions according to the manipulation of  $F_0$ . ERB: Equivalent rectangular bandwidth.

Table 1. Data from Experiments 1 (sentences) and 2 (words). Percentages of successes in a binary discrimination task on the speakers' body height, separated by experimental conditions according to the manipulation of Fo. ERB: Equivalent rectangular bandwidth.

	Fo manipulation				
	-1.0 ERB	-0.5 ERB	Normal	+0.5 ERB	+1.0 ERB
Sentences (1 <sup>st</sup> Exp.)					
male speakers		55.6 (17.8)	55.9 (14.8)	61.2 (14.9)	
female speakers		58.4 (17.6)	62.2 (19.7)	66.0 (16.0)	
Words (2 <sup>nd</sup> Exp.)					
female speakers	61.0 (16.0)	63.1 (15.5)	65.0 (18.2)	65.1 (15.0)	67.3 (12.3)

For Peer Review

Figure 1. Experiment 1 (sentences). Percentages of successes in a binary discrimination task on the speakers' body height, separated by experimental conditions according to the manipulation of Fo and the speaker sex.

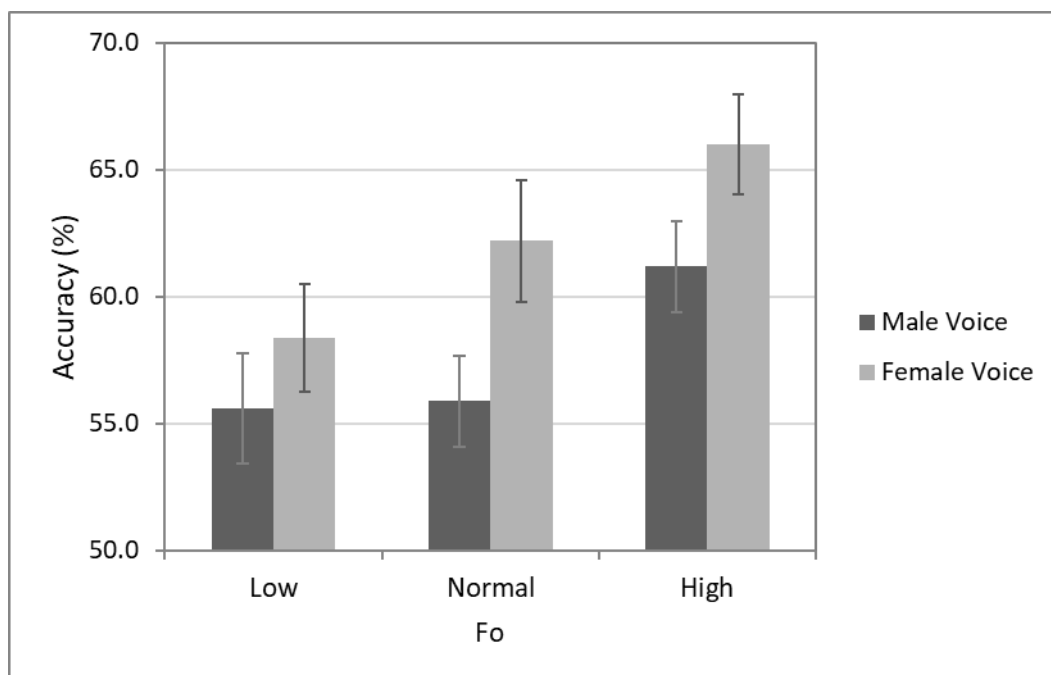


Figure 2. Experiment 2 (words). Percentages of successes in a binary discrimination task on the female speakers' height, separated by experimental conditions according to the manipulation of  $F_0$ . ERB: Equivalent rectangular bandwidth.

