

Title: Cross-language differences in how voice quality and f_0 contours map to affect¹

Authors: Irena Yanushevskaya², Christer Gobl, Ailbhe Ní Chasaide

Phonetics and Speech Laboratory, School of Linguistic, Speech and Communication

Sciences, Trinity College Dublin, Dublin, Ireland

Phone: +353 1 8961348

[The manuscript uploaded to the JASA submission system: October 12, 2018]

Running title: Voice quality and cross-language affect perception

¹ Portions of this work were presented in 'Universal and language-specific perception of affect from voice', Proceedings of the XVIIth International Congress of Phonetic Sciences, Hong Kong, August, 2011.

² Author to whom correspondence should be addressed. Electronic mail: yanushei@tcd.ie

ABSTRACT

The relationship between prosody and perceived affect involves multiple variables. This paper explores the interplay of three: voice quality, f_0 contour, and the hearer's language background. Perception tests were conducted with speakers of Irish English, Russian, Spanish and Japanese using three types of synthetic stimuli: (1) stimuli varied in voice quality, (2) stimuli of uniform (modal) voice quality incorporating affect-related f_0 contours, and (3) stimuli combining specific non-modal voice qualities with the affect-related f_0 contours of (2). The participants rated the stimuli for the presence/strength of affective colouring on six bipolar scales, e.g., *happy-sad*. The results suggest that stimuli incorporating non-modal voice qualities, with or without f_0 variation, are generally more effective in affect cueing than stimuli varying only in f_0 . Along with similarities in the affective responses across these languages, many points of divergence were found, both in terms of the range and strength of affective responses overall and in terms of specific stimulus-to-affect associations. The f_0 contour may play a more important role, and tense voice a lesser role in affect signalling in Japanese and Spanish than in Irish English and Russian. The greatest cross-language differences emerged for the affects *intimate*, *formal*, *stressed* and *relaxed*.

1 I. INTRODUCTION

2 This paper explores how the voice communicates affect and considers whether the
3 language background of the listeners impacts on the affective colouring associated with
4 variations in the voice. It looks specifically at how the dimensions of voice quality and
5 f_0 map individually to affect and also at how certain combinations of these dimensions
6 might work. Perception tests were carried out on native speakers of four languages,
7 Irish English, Russian, Spanish and Japanese.

8

9 It is generally appreciated that voice quality plays a fundamental role in spoken
10 communication, carrying an affective strand of information that conveys not only the
11 speaker's emotion, mood and state, but also her/his attitude towards the interlocutor or
12 the current situation. This aspect of speech communication has nonetheless remained
13 elusive, being relatively neglected within the linguistic and speech sciences, and
14 researched mostly within psychology, e.g. Scherer (2003). As man-machine
15 communication is becoming ever more speech enabled, this aspect of spoken
16 communication has moved centre stage – speech dialogue systems are increasingly
17 needed that are sensitive to the affective dimension of the speaker's voice and generate
18 speech with a voice quality that is accordingly modulated in appropriate ways
19 (Burkhardt and Stegmann, 2009). Increasingly, the detection of the speaker's affective
20 state is also important to applications monitoring affective and mental wellbeing
21 (Cummins *et al.*, 2015).

22

23 Despite the increased research focus on this topic in recent decades, summarised in
24 several reviews, e.g., Murray and Arnott (1993), Scherer (2003), Gobl and Ní Chasaide
25 (2003b), Juslin and Laukka (2003), Juslin and Scherer (2005), Laukka (2008), Cowie

1 (2009), Kreiman and Sitdis (2011), Scherer (2013), reliable information on the specific
2 voice source correlates of affect is relatively scarce. This area of research is beset with
3 numerous conceptual and methodological difficulties. To begin with, defining
4 affects/emotions is a complex issue (Cowie and Cornelius, 2003). A simple term like
5 ‘angry’ can be interpreted to mean rather different emotions (‘hot’ and ‘cold’ anger are
6 commonly differentiated in the literature, e.g., Juslin and Laukka (2001), but other
7 states belonging to the same *family* of this emotion can be identified). Obtaining
8 reliable speech samples with specific affects is difficult, particularly if one needs
9 recordings that are amenable to voice source analysis. Likewise, the terms that are
10 commonly used to describe voice quality are often used in vague and inconsistent ways
11 in the literature – see discussion of this issue in Gobl and Ní Chasaide (2003b),
12 although the use of Laver’s (1980) classification framework enables a more rigorous
13 treatment. Probably the greatest roadblock in the field concerns the difficulty in
14 obtaining reliable, accurate production measurements of voice source features,
15 particularly for non-modal voice qualities and in connected speech (Gobl and Ní
16 Chasaide, 2010) – which is precisely the type of data that is needed. These issues are
17 elaborated on in some detail in Section II below.

18

19 The picture is further complicated by the fact that the voice is not uniquely a carrier of
20 affect: modulation of the voice source is an integral part of the linguistic (non-affective)
21 prosody about which rather little is known – but see for example Heldner (2003), Iseli
22 *et al.* (2006), Ní Chasaide *et al.* (2013), as linguistic-phonetic research on prosody is
23 largely focused on f_0 modulation. Furthermore, the segmental context induces
24 substantial microprosodic effects on the source (Ní Chasaide and Gobl, 1993) – and this
25 can also be an important consideration when interpreting voice source data.

1
2 Cross-language studies (see Section II) typically entail the presentation of an utterance
3 produced with affect in one language to listeners of other languages, who are asked to
4 identify the affect they hear. This provides many insights but does not elucidate which
5 aspects of the voice are cueing listener's perceptions. The present experiment follows a
6 different approach, used in earlier studies by the authors to explore voice-to-affect
7 mapping for speakers of Irish English, e.g., Gobl *et al.* (2002), Gobl and Ní Chasaide
8 (2003b), Ryan *et al.* (2003). It entails the presentation of stimuli synthesised with
9 different voice qualities to listeners of different language backgrounds and elicitation of
10 the affective colouring they impart. The stimuli were constructed to provide exemplars
11 of specific voice qualities chosen to provide a representative sample of qualities most
12 frequently mentioned in the literature. The synthesised utterances incorporate the kinds
13 of voice modulations that reflect the linguistic prosody as well as segment-dependent
14 perturbations that have been alluded to above.

15
16 As in those earlier experiments, Laver's (1980) theoretical framework provides a basis
17 for the synthesis of the voice quality stimuli used here: Tense Voice, Whispery Voice,
18 Breathy Voice and Lax-creaky Voice. Laver's framework links auditory vocal qualities
19 to their underlying physiological mechanisms (laryngeal tension settings) as well as to
20 their main acoustic features. Past analytic studies by the authors of voice qualities
21 produced within Laver's system have guided stimulus construction (Gobl, 1989; Gobl
22 and Ní Chasaide, 1992; Ní Chasaide and Gobl, 1995). It is important to note that a
23 specific voice quality is not a set point (like a cardinal vowel) but may vary on a
24 continuum, and strength of the associated affect appears to vary accordingly (Ryan *et*
25 *al.*, 2003). The voice quality used here are not extreme exemplars: the point was to

1 explore how moderate shifts, not necessarily extreme ones, in voice quality might
2 influence perception of affect.

3

4 Differences in the f_0 contour (level, range and dynamics) have also been widely
5 reported production correlates of affect, even though there are questions as to how
6 perceptually important they may be (see Section II.C). To consider the role of f_0 in the
7 present study, two further sets of stimuli were included: a series where the f_0 contours
8 were varied and the voice quality kept constant and a series where these f_0 contours
9 were combined with different voice qualities from the first series. The f_0 contours used
10 here were based on production data, elicited for affective utterances of Dutch
11 (Mozziconacci, 1995). The contours varied in terms of their level, range and dynamics.
12 Unlike the voice quality series, it did include rather extreme contours, with large
13 deviations from the neutral.

14

15 The two main research questions addressed by the present study are:

16 (1) how do voice quality and f_0 dimensions of the voice communicate affect?

17 (2) how does the language background of the listeners influence the mapping of voice
18 to affect?

19

20 The underlying principle of this experiment was to present the listeners with a ‘palette’
21 of vocal stimuli and have them ‘paint the picture’ as to their affective colouring.

22 Consequently, the focus was not on testing specific hypotheses. However, there are
23 certain expectations regarding results, arising from our understanding of the

24 physiological basis for the different voice qualities (following Laver’s framework), and
25 based also on the findings of earlier studies.

1
2 On the one hand, one would expect that tense voice (with high laryngeal tension
3 settings) would be associated with high activation, and breathy (or lax) voice with low
4 activation affects. This expectation is rather like the ‘effort code’ proposed to account
5 for aspects of linguistic intonation by Gussenhoven (2004) which suggests that
6 language prosody ‘rules’ are built on underlying production correlates. This might
7 emerge therefore as an expected common feature, a potential candidate for universality.
8 This expected trend did occur for tense voice in Gobl and Ní Chasaide (2003b), but in
9 that and subsequent studies the lax-creaky setting was found to be rather more effective
10 than breathy voice in conjuring low activation states, and so, both of these qualities are
11 included here. As whispery voice entails more laryngeal tension than breathy voice, it
12 was expected that this quality would be relatively less strongly associated with low
13 activation. Language specific effects might also be expected (see Section II). For
14 example, whispery voice has been mentioned as being associated with fear in English
15 (Boula de Mareüil *et al.*, 2002), but with different qualities in other languages; breathy
16 voice is traditionally associated with intimacy in English (Laver, 1980) but with
17 formality/politeness in Japanese (Ito, 2004; Ishi *et al.*, 2008). Creaky voice tends also to
18 be associated with different affective states in different languages.

19
20 The f_0 contours were based on affective production data, elicited in Dutch for a number
21 of affects (sad, bored, etc.), and these affective labels are retained in this study. Given
22 our ‘broad palette’ approach, we would refrain from strong predictions concerning how
23 speakers of other languages might perceive these contours. As discussed in Section II,
24 past studies have found such affective f_0 contours to be relatively ineffective in cuing
25 the specific affects of the original production data. However, an expectation based on

1 the ‘effort code’ would be that the contours with the more extreme deviations from
2 neutral would be associated with greater affective responses, e.g., the higher the f_0
3 contour the more strongly it would be associated with high activation states.

4

5 As regards the stimuli combining both voice quality and f_0 manipulations, the broad
6 expectation would be that these dimensions would work synergistically, so that
7 combined stimuli would be more effective in signalling affect than those stimuli
8 involving simply voice quality or f_0 modifications on their own. Furthermore, there is a
9 possibility that even though f_0 manipulations might not be effective in signalling the
10 affect (of original production data) when in association with an appropriate voice
11 quality setting such affects (sad, bored, etc.) would more clearly emerge. Although not
12 generally attested in the earlier experiments on English, this could be the case for other
13 languages.

14

15 Although the phonetic literature has typically assumed that specific voice qualities map
16 directly to specific affects (e.g., creaky voice associates with boredom), our past
17 experiments suggest that there is no one-to-one mapping of voice to affect. Rather, a
18 single stimulus was found to be typically associated with several affects, not always
19 related in any obvious way. As a corollary, a given affect might be conjured by more
20 than one stimulus.

21

22 **II. BACKGROUND AND LITERATURE REVIEW**

23 Cross-language production and perception studies are discussed here, along with
24 methodological difficulties that arise, particularly in elucidating the voice source
25 correlates of affect. It is hoped that this broad review can help provide insight into the

1 complexities and difficulties pertaining to research in this field. It is also aims to
2 provide a backdrop that explains the motivations for the present approach.

3

4 **A. Expression of affect: universal or language specific?**

5 As with the facial expression of emotion (Ekman *et al.*, 1987; Ekman, 1993), the
6 spontaneous vocal expression of full-blown ‘basic’ (primarily negative) emotions is
7 generally assumed to be universal (Scherer, 2000; Zinken *et al.*, 2008; Sauter *et al.*,
8 2010b; Scherer *et al.*, 2011; Laukka *et al.*, 2016) reflecting the influence of
9 physiological ‘push effects’ (e.g., changes in rate of respiration and muscular tension)
10 on the mechanism of voice production. However, in real interactions the affective
11 expression of mood, attitude and interpersonal stance is not so extreme, involving much
12 more subtle modulations of the voice than ‘full blown’ emotions and less directly
13 linked to physiological ‘push’ effects.

14

15 Furthermore, the spontaneous expression of affect is constrained by ‘pull effects’ –
16 external factors related to the socially accepted ‘display rules’, the listener’s
17 expectations, the speaker’s self-presentation, etc. (Johnstone and Scherer, 2000; Juslin
18 and Scherer, 2005). Cultural factors codified in language may thus intervene in
19 determining how acceptable it may be to express or suppress the display of affect
20 (Mesquita and Walker, 2003; Oatley *et al.*, 2006). Complex cognitive states such as
21 irony or sarcasm seem to involve intentional mismatching of voice and verbal content:
22 these are particularly likely to be governed by culturally informed display rules (Zinken
23 *et al.*, 2008; Sauter *et al.*, 2010b; Scherer *et al.*, 2011). The pragmatic context of the
24 interaction is also important, regardless of culture. We can and routinely do control our

1 voice quality. As pointed out by Scherer (2013), the listener's attribution of affect and
2 the speaker's underlying emotion do not need to coincide, and this is a feature of
3 normal human interactions.

4

5 **B. Production studies**

6 ***1. Data elicitation***

7 Obtaining reliable and valid affectively coloured speech data is not a straightforward
8 task; e.g., discussion in Juslin *et al.* (2017). Naturally occurring spontaneous expression
9 can be recorded during TV games or reality shows (Douglas-Cowie *et al.*, 2003) or in
10 everyday social interactions (Campbell, 2002). However, the recording conditions and
11 quality and the high variability of the spoken content tend to render these kinds of data
12 unsuited to acoustic analysis that would isolate those voice source measures that can be
13 correlated to affect. Speech samples obtained with mood induction techniques
14 (Westermann *et al.*, 1996) allow control over the verbal and emotional content but
15 produce relatively low in intensity and not clearly differentiated emotional expressions
16 (Scherer, 2003).

17

18 The dominant approach remains the collection of acted expressions of affect (Bänziger
19 and Scherer, 2007; Scherer, 2013; Laukka *et al.*, 2016). Typically, the utterances used
20 are content-neutral, i.e. repetitions of a single word (e.g., a name) or specially
21 constructed nonsense sentences (Scherer *et al.*, 2001). Such simulated data can produce
22 intense prototypical expressions but allows experimental control over the lexical
23 content of the utterances and affects expressed as well as cross-speaker comparison of
24 results. A study on two emotions (*positive/happy* and *negative/sad*) by Scherer (2013)
25 found no difference between acted and induced affective vocalisations. However, the

1 approach has been criticised (Bachorowski and Owren, 2003; Russell *et al.*, 2003;
2 Juslin *et al.*, 2017) mostly on the grounds that portrayals may reflect exaggerated,
3 culture-specific stereotypes rather than genuine real-life expressions.

4 **2. Production correlates of affect**

5 Not surprisingly therefore, empirical analytic studies to date have yielded relatively
6 little information on the voice quality correlates of affect and have tended to focus
7 rather on aspects that are more easily measured, particularly f_0 and intensity, as well as
8 speech rate. (Speech rate emerges as important but is not covered in the present
9 review.) Scherer and colleagues (Pakosz, 1983; Scherer *et al.*, 1984; Ladd *et al.*, 1986;
10 Carlson *et al.*, 1992; Mozziconacci, 1995; Banse and Scherer, 1996; Mozziconacci and
11 Hermes, 1999; Paeschke *et al.*, 1999; Mozziconacci, 2002; Paeschke, 2004; Bänziger
12 and Scherer, 2005; Grandjean *et al.*, 2006) have provided extensive data on how the
13 level, range and dynamics of f_0 and intensity are correlated with affect in speech
14 production. Typically, global statements about these acoustic parameters are reported,
15 e.g., greater level and range of f_0 for high activation affects (Banse and Scherer, 1996).
16 However, when the perceptual relevance of production findings on f_0 differences was
17 tested in a number of the above studies, results were disappointing, leading to a
18 conclusion that the key to the differentiation of emotion lies in the voice quality
19 (Scherer, 1986).

20

21 More recent studies (Patel *et al.*, 2011; Sundberg *et al.*, 2011) have emphasised the
22 need to move towards more direct measurement of the voice source. Scherer (2013) has
23 pointed out that the acoustic measures hitherto most commonly used in emotion
24 research and inherited from phonetics, such as f_0 and intensity, are not in fact optimal
25 indicators of the vocal changes that occur with affect. Unique acoustic patterns for

1 different emotions do exist, but the most appropriate acoustic measures to describe
2 them have not yet been identified, suggesting the need for more complex acoustic
3 measures, based on direct source measurements (Patel *et al.*, 2011; Bänziger *et al.*,
4 2015).

5
6 Direct measurements of the voice source are not easy to obtain, however, given the
7 difficulty in obtaining the glottal flow signal itself. A non-invasive technique for
8 estimating the glottal flow involves inverse filtering of the speech pressure waveform,
9 whereby the effect of the vocal tract transfer function is cancelled (e.g., Gobl and Ní
10 Chasaide, 2010; Alku, 2011). This type of analysis generally requires stringent
11 recording conditions (Gobl, 2003; Alku, 2011), and as automatic methods tend to yield
12 considerable errors, manual editing is ideally required – a technically complex and
13 time-consuming task, which greatly limits the quantity of data that can be analysed
14 (Gobl and Ní Chasaide, 2010). Indirect inferences about the source can also be made
15 based on measurements carried out on the speech signal (Hanson, 1997; Hanson and
16 Chuang, 1999; Keller, 2005; Alku, 2011). However, caution is needed in interpreting
17 such data as they are influenced by the vocal tract filter; see discussion in Gobl and Ní
18 Chasaide (2010).

19
20 The most commonly used indirect measures of voice quality in emotional speech
21 research mainly describe spectral slope or spectral balance, i.e. the relative amount of
22 acoustic energy above and below a certain cutoff frequency, e.g., the alpha ratio
23 (Sundberg and Nordenberg, 2006) or the Hammarberg index (Hammarberg *et al.*,
24 1980), and also formant bandwidths (Juslin and Scherer, 2005; Laukka *et al.*, 2005;
25 Goudbeek and Scherer, 2010; Laukka *et al.*, 2011). High activation or negative valence

1 states have been reported to have a less steep spectral slope (likely to be indicative of
2 tensor phonation) while low activation or positive valence states show the opposite
3 trend (Pittam *et al.*, 1990; Laukka *et al.*, 2005; Goudbeek and Scherer, 2010; Guzman
4 *et al.*, 2013). Given the potential influence of the vocal tract filter setting, such
5 measures may only be partially indicative of what is happening at the level of the voice
6 source.

7

8 Given the technical complexity of voice source analysis, the quantity and scope of
9 emotional speech voice source production data available is very limited. Most studies
10 have tended to focus on a single vowel (or sometimes on a single glottal pulse)
11 extracted from utterances produced with different affects, by typically one or very few
12 speakers. Although the studies provide important insights, such measures cannot
13 capture the important dynamic aspects of voice modulation in affect expression.

14

15 Analysis of glottal waveforms, reported in Cummings and Clements (1995), Laukkanen
16 *et al.* (1996), Murphy and Laukkanen (2009) revealed significant differences in the
17 glottal pulse characteristics/shape in speech produced with different speaking styles and
18 different affective colouring. Similar findings were reported in a more extensive study
19 (Patel *et al.*, 2011; Sundberg *et al.*, 2011), where several voice source parameters were
20 measured in the [a] vowel, as produced by 10 speakers for a range of affective states.

21 In a different approach, the dynamic variation across entire utterances was analysed for
22 a single male speaker portraying a range of affective states (Yanushevskaya *et al.*,
23 2007; Yanushevskaya *et al.*, 2009). Results show distinct patterns of source parameter
24 settings for each affect in terms of the within-utterance averages, but there was
25 considerable overlap in parameter trajectories in parts of these utterances.

1

2 One of the issues arising in comparing different studies concerns the number and choice
3 of source parameters measured. The ‘global’ parameter NAQ (Alku and Vilkmann,
4 1996, see also Fant, 1995; 1997 on the global waveshape parameter R_d) intended to
5 capture the overall characteristics of the glottal waveshape has been found to correlate
6 to affect (Airas and Alku, 2004; Laukkanen *et al.*, 2008; Waaramaa *et al.*, 2008),
7 particularly to the activation dimension.

8

9 Other studies like Sundberg *et al.* (2011), Patel *et al.* (2011), Cummings and Clements
10 (1995), Murphy and Laukkanen (2009), Yanushevskaya *et al.* (2007), have measured
11 more numerous, but not always the same, source parameters. This raises the question
12 which pertains to all work on the voice source, be it directed at emotional expression or
13 not, of what the optimal set of parameters might be. When large numbers of parameters
14 are measured, there are considerable redundancies, as there are inevitable
15 interdependencies among them. The issue of how many parameters are needed to
16 capture production differences on the one hand and which can on the other hand lead to
17 perceptually distinct changes to voice quality is one that still requires elaboration. It has
18 been suggested by Eyben *et al.* (2016) that studies on the vocal expression of affect
19 should all include comprehensive coverage of 65 acoustic features (not all of which are
20 voice source measures), to ensure comparability across studies. While this is
21 undoubtedly a desirable goal, when it comes to measurements of the voice source, a
22 scattergun approach may be of little value if the accuracy of the source measurements is
23 not assured or indeed for any given acoustic measure whose perceptual relevance has
24 not been demonstrated. An interesting model in this context is that of Kreiman *et al.*

1 (2014), which aims to integrate voice production and perception, and whose acoustic
2 parameters specifically aim to capture variation in voice quality.

3

4 In many studies there is an assumption that variation in a voice source parameter (or set
5 of parameters) can be mapped directly to affect, and this assumption may not be borne
6 out – for several reasons. NAQ is considered to correlate with audible variation in the
7 tense-lax dimension of voice quality (higher NAQ indicating breather voice) (Airas and
8 Alku, 2007) and as such, has been used as a possible source correlate of affect
9 (Waaramaa *et al.*, 2008; Patel *et al.*, 2011). However, as pointed out in Gobl and Ní
10 Chasaide (2003a) a relatively high NAQ value is not invariably indicative of a
11 relatively breathier voice quality and the correlation may only pertain to phonation in a
12 given f_0 range. And, as mentioned above, in interpreting voice source measures, the
13 non-affective (linguistic-prosodic and segment-related) sources of voice source
14 variation need to be borne in mind, as these are not perceived by the listener as shifts in
15 auditory voice quality. Ultimately, voice source variation is unlikely to express affect
16 unless there is a perceptible shift in the auditory voice quality – hence the approach
17 adopted in this study.

18

19 ***3. Cross-language production studies***

20 Several cross-language production studies suggest that basic emotions (in acted
21 portrayals) may not be encoded universally (at least based on the acoustic parameters
22 measured). Anger is usually described as having high f_0 and high intensity, e.g., in
23 German (Banse and Scherer, 1996) and English (Williams and Stevens, 1972); see also
24 Juslin and Scherer (2005), Kappas *et al.* (1991), Juslin and Laukka (2003) for other
25 languages. This was found also in Russian but not in Estonian, where pitch and

1 intensity were found to be lower in the expressions of anger than in the neutral state
2 (Altrov, 2013). Significant cross-language differences in f_0 mean and range in the
3 expression of anger have been reported in Pell *et al.* (2009a), Pell *et al.* (2009b).
4 Laukka *et al.* (2016) reported significant effects of language in the analysis of affective
5 speech for many types of acoustic cues (f_0 , intensity, spectral balance, and temporal
6 characteristics).

7

8 Auditory judgements of portrayals of the same emotion across-languages suggest
9 differences in the voice qualities used. For example, fear was reported to be expressed
10 with falsetto voice by German speakers (Burkhardt and Sendlmeier, 2000; Schröder,
11 2001), but with whispery voice by an English actor (Boula de Mareüil *et al.*, 2002) and
12 with breathy voice in Italian (Drioli *et al.*, 2003). Similarly, the same voice quality may
13 carry quite different connotations for speakers with different language backgrounds.
14 Thus, creaky voice tends to be associated with boredom in English and German (Laver,
15 1980; Burkhardt and Sendlmeier, 2000). In Japanese, where the use of pressed or
16 creaky voice ('rikimi') follows very complex rules of social interaction, creaky voice
17 has been reported to convey a range of affective states, e.g., surprise, irritation, disgust,
18 anger and admiration (Sadanobu, 2004). In Chinese, creaky voice has been found in the
19 expression of anger (Yuan *et al.*, 2002), while in Italian a creaky voice quality was
20 found in expressions of disgust (Drioli *et al.*, 2003).

21

22 While these observations are possibly indicative of language-specific vocal expression
23 of emotion, one must be cautious in interpreting such apparent differences. There is
24 often uncertainty as to whether a given impressionistic voice quality label refers to the
25 same acoustic-auditory phenomenon in different studies. Can one be sure that the

1 breathy voice and the whispery voice (reported for expressions of fear for the Italian
2 and English speakers respectively) are in fact different qualities? Similarly, is the
3 creaky voice reported as associated with boredom in English and German in fact the
4 same quality as the creaky voice of Japanese ‘rikimi’? Furthermore, even in a single
5 language a particular voice quality can be associated with a number of different affects
6 and a particular affect (such as sadness) can be indicated by more than one voice
7 quality. The lack of bi-uniqueness raises the possibility that such cross-language
8 differences might sometimes be more apparent than real – particularly as the context
9 and methodologies of different studies can differ considerably.

10 **C. Perception studies**

11 Although many production studies have described shifts in the level, range and
12 dynamics of f_0 in affectively coloured speech (Scherer *et al.*, 1984; Ladd *et al.*, 1985;
13 Banse and Scherer, 1996; Grandjean *et al.*, 2006), these features alone do not generally
14 appear very effective in cueing affect (Mozziconacci, 1998; Bänziger and Scherer,
15 2005). Such lack of correspondence of the production and perception findings has led
16 to the overall conclusion, mentioned above, that these are not the optimal acoustic
17 features, and that more complex, voice source measures are needed to capture how
18 listeners decode emotion (Scherer, 2013; Bänziger *et al.*, 2015). Not surprisingly, given
19 the sparsity of production data on affect-related voice source features, rather little is
20 known about their influence/salience in affective cueing.

21

22 As described in the Introduction, studies exploring the perceptual relevance of voice
23 quality in the cueing of affect were carried out by the authors with Irish English

1 participants (Gobl and Ní Chasaide, 2000; Gobl *et al.*, 2002; Gobl and Ní Chasaide,
2 2003b). These are the precursors to the present cross-language investigation.

3

4 **1. *Cross-language perception studies***

5 Unlike production studies, the preponderance of research on perception has been
6 directed at cross-language aspects. The standard paradigm uses acted portrayals of a
7 limited set of discrete emotion categories on content free utterances. Perception tests
8 entail forced-choice judgement among the emotion categories being tested and report
9 the accuracy with which the intended expressed emotions are identified (Bänziger *et*
10 *al.*, 2015). Typically, individuals from many language/ethnic groups judge stimuli
11 produced in one language, individuals from one language/ethnic group judge stimuli
12 from many languages or a balanced design is used (Laukka *et al.*, 2016). Many such
13 cross-language perception studies suggest broad similarities in affect perception from
14 voice by listeners with different language background, e.g., McCluskey *et al.* (1975),
15 McCluskey and Albas (1981), van Bezooijen *et al.* (1983), Scherer (2000), Graham *et*
16 *al.* (2001), Pell and Scorup (2008), Koeda *et al.* (2013), Waaramaa and Leisiö (2013),
17 Waaramaa (2014), see also review of perception studies in Laukka *et al.* (2016).

18

19 In general, it is reported that negative emotions tend to be recognised more accurately
20 cross-linguistically than positive emotions (Sauter *et al.*, 2010a; Sauter *et al.*, 2010b)
21 and basic emotions such as sadness and fear are recognised more accurately compared
22 to more complex cognitive affective states, such as contempt or sarcasm, e.g., van
23 Bezooijen *et al.* (1983), Cheang and Pell (2008), Sauter *et al.* (2010a).

24

1 It has been suggested that the same emotion may be expressed ‘more clearly’ and
2 recognized with higher accuracy in certain languages; for example, in the studies by
3 McCluskey *et al.* (McCluskey *et al.*, 1975; McCluskey and Albas, 1981) the Mexican
4 subjects were found to be significantly more sensitive to expressions of happiness,
5 sadness, love, and anger in speech than their Canadian counterparts. Emotions
6 expressed by Mexicans were generally recognised with higher degree of accuracy by
7 both Mexicans and Canadians than those expressed by Canadians.

8

9 As suggested in Kitayama and Ishii (2002), certain cultures may be more sensitive to
10 vocal expression, while for others, the propositional content matters more than the tone-
11 of-voice. Furthermore, due to in-group advantage, listeners tend to recognise emotions
12 vocally expressed in their native language with more accuracy than those expressed by
13 representatives of other cultures (Elfenbein and Ambady, 2002b; a; Bryant and Barrett,
14 2008).

15

16 Language relatedness may be a factor in the accuracy of cross-language emotion
17 attribution. In a large-scale study (Scherer, 2000; Scherer *et al.*, 2001), listeners from
18 differing language backgrounds judged simulated emotions produced by German actors
19 in artificially constructed sentences. The greater the language dissimilarity from
20 German, the less accurate the identification of emotion was. Similar findings are
21 reported more recently by Waaramaa and Leisiö (2013) and Waaramaa (2014).
22 However, no such advantage of the language/culture proximity was found in Pell *et al.*
23 (2009a) and Altrov and Pajupuu (2015).

24

1 These studies do not clarify which aspects of the voice signal affect within a language
2 group, or whether these differ across languages. Perception tests by Burkhardt *et al.*
3 (2006) address such questions by presenting listeners with stimuli of MBROLA-
4 synthesised, emotionally loaded sentences, where pitch range, duration and jitter have
5 been manipulated. Listeners with different language backgrounds (French, German,
6 Greek and Turkish) showed differences in sensitivity to the jitter and pitch range. A
7 narrow pitch range was associated with a friendly affect by Turkish listeners, but not by
8 the speakers from other countries. Turkish listeners also associated an increased amount
9 of jitter with a threatening affect – an association not made by German and French
10 listeners: for the French listeners the jitter-simulation evoked a frightened affect,
11 whereas the converse was found for the Germans.

12

13 As detailed in the Introduction, the study presented here differs from the above
14 approaches by presenting listeners from differing language backgrounds with stimuli
15 containing auditorily distinct voice qualities in a semantically neutral utterance and
16 eliciting how they are associated with affect by the different groups. In addition to
17 stimuli differing in voice quality, stimuli with different f_0 contours (all with modal
18 voice) were also included, as well as combined stimuli where voice quality and f_0
19 contours were varied. Of interest is to see how similar/different these language groups
20 might be in associating the different dimensions of the voice to affect. Although not
21 testing specific hypotheses, our expectations in terms of likely/possible results are
22 described in Introduction.

23

1 III. MATERIAL AND METHOD

2 A. Synthesised stimuli

3 The stimuli used in this cross-language study are based on a recording of a short
 4 Swedish utterance ‘ja adjö’ [‘ja: a‘jø:] produced by a male speaker with a voice quality
 5 conforming to Laver’s (1980) description of modal voice. To create the modal voice
 6 stimulus, high quality copy synthesis of this natural utterance was carried out using the
 7 KLSYN88 formant synthesiser (Sensimetrics Corporation, Boston, MA, described in
 8 Klatt and Klatt, 1990) employing its version of the LF model for the glottal waveform
 9 generation. By modifying the modal voice stimulus, three series of distinct stimulus
 10 types were generated, referred to as the VQ stimuli, the F0 stimuli, and the VQ+F0
 11 stimuli, as outlined in the following sections.

12

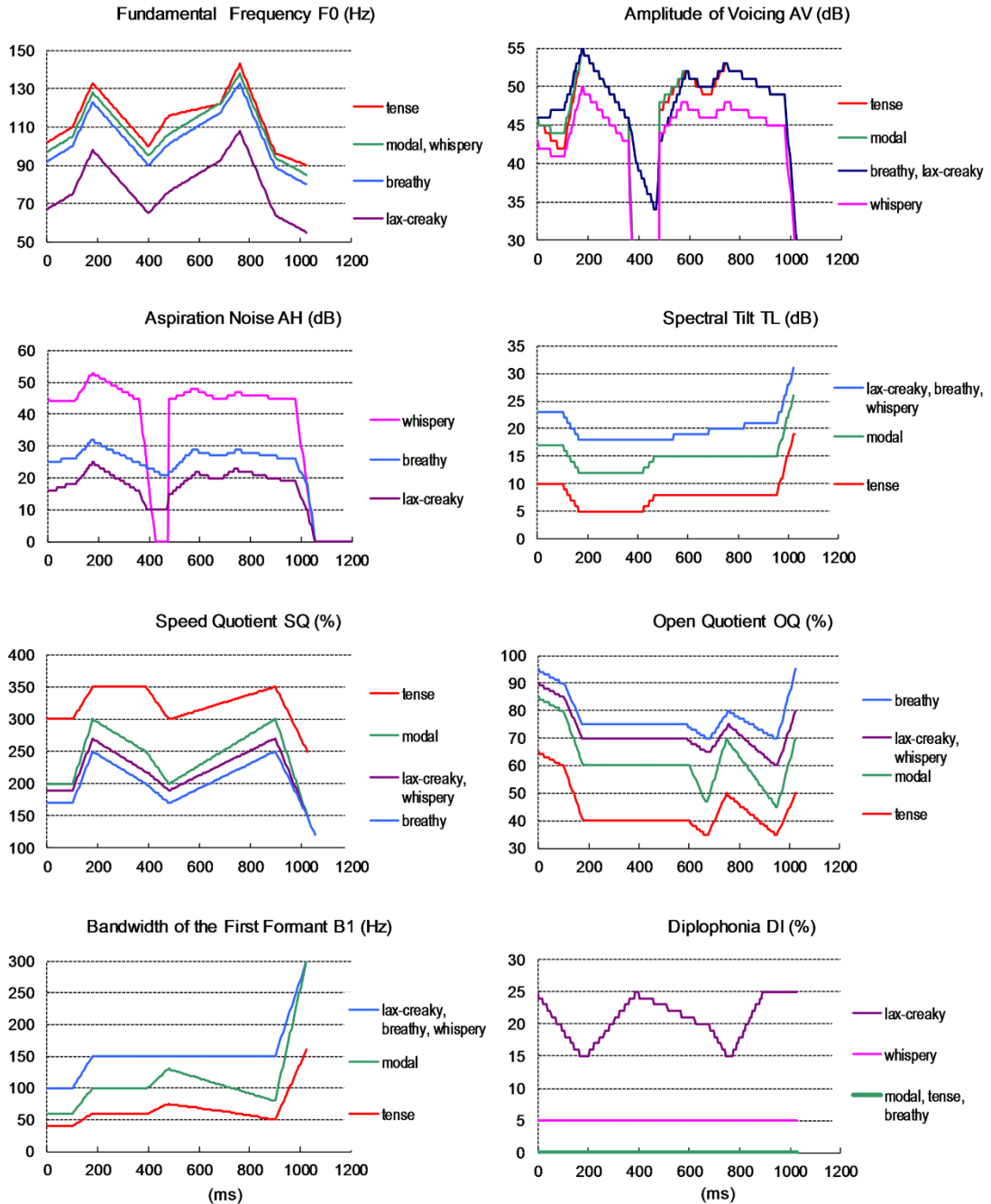
13 To obtain the source and filter settings used in the synthesis for the modal voice
 14 stimulus, the utterance was inverse filtered and the source signal modelled by matching
 15 the LF model (Fant *et al.*, 1985) to the output of the inverse filter. Both processes were
 16 carried out using the manual, interactive technique described in Gobl and Ní Chasaide
 17 (1999) and Gobl and Ní Chasaide (2010), which allows pulse-by-pulse optimisation of
 18 the parameter settings.

19 1. VQ stimuli

20 The non-modal voice qualities used in this study were *Whispery Voice*, *Breathy Voice*,
 21 *Lax-creaky Voice* and *Tense Voice*. As mentioned earlier, they were not extreme
 22 exemplars of these qualities. With some minor modifications, these stimuli are as
 23 described in detail in Gobl and Ní Chasaide (2003b) and they were generated from the
 24 modal voice stimulus by manipulating the following parameters of the KLSYN88

1 synthesiser: fundamental frequency (F0), amplitude of voicing (AV), spectral tilt (TL),
2 open quotient (OQ), speed quotient (SQ), aspiration noise (AH), diplophonia (DI), and
3 the bandwidth of the first formant (B1). The dynamic variation of parameters was
4 determined by earlier analytic studies of specific voice qualities (Gobl, 1988; Gobl and
5 Ní Chasaide, 1992; Ní Chasaide and Gobl, 1995) and was adjusted on the basis of
6 auditory analysis of the voice qualities (Gobl and Ní Chasaide, 2003b). Informal
7 evaluation of the VQ stimuli by phoneticians trained in Laver's (1980) framework was
8 used to ensure that the intended voice quality was achieved. The parameter variation of
9 the VQ stimuli is illustrated in Figure 1.

10



1

2 Figure 1 (color online). Parameter variation in the synthesised voice quality stimuli

3 (after Gobl & Ní Chasaide, 2003).

4

5 The synthesised voice quality stimuli in Gobl and Ní Chasaide (2003b) also

6 incorporated minor adjustments of f_0 values, deemed likely to be typical of these7 qualities. For example, f_0 is 5 Hz higher for the *Tense Voice* stimulus and 5 Hz lower

1 for the *Breathy Voice* stimulus compared to *Modal Voice*. In one case the f_0 difference
2 was more substantial: for *Lax-creaky Voice* an f_0 value 30 Hz lower than that of the
3 *Modal* (neutral) stimulus was used, as this quality most typically occurs with low f_0 .
4 While this feature was retained for the present experiment, it resulted in a certain
5 anomaly necessitating consideration for the presentation of results. This is discussed
6 further in Section III.A.3.

7

8 As can be seen in Figure 1, the trajectories for all source parameters show considerable
9 utterance-internal dynamic variation. These reflect the variation due to the linguistic
10 prosody of the utterance (Ní Chasaide and Gobl, 2004b) as well as perturbations which
11 are the consequence of consonantal articulation (Ní Chasaide and Gobl, 1993) and these
12 have been discussed in Section II. These variations were retained in the other VQ
13 stimuli.

14 **2. F_0 stimuli**

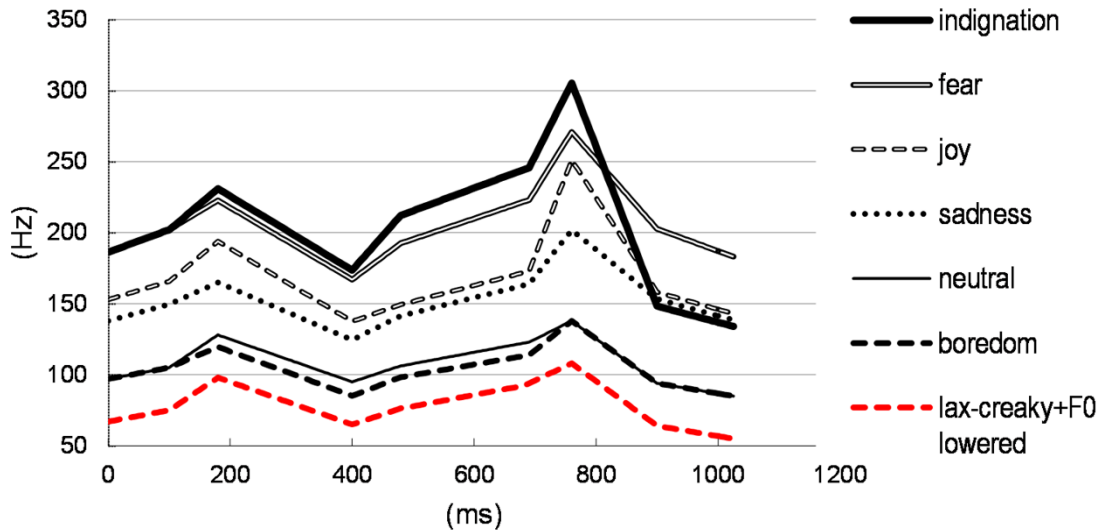
15 The creation of this series of stimuli entailed modifications to the *Modal* stimulus by
16 manipulating the F_0 parameter of the KLSYN88 synthesiser. Five further stimuli were
17 produced with different f_0 contours by adapting the affect-related f_0 contours presented
18 in Mozziconacci (1995). The latter were based on the analysis of a short semantically
19 neutral utterance (containing two accented words) produced by a male Dutch speaker
20 portraying indignation, joy, fear, boredom, sadness as well as a neutral state.

21

22 As the f_0 contour of the modal utterance used in this study was very similar to the
23 neutral contour of Mozziconacci (1995), it served as a neutral reference f_0 contour. The
24 five additional contours were obtained by proportional scaling of the f_0 values of our
25 reference stimulus at nine anchor points, using scale factors derived from

1 Mozziconacci's contours. The contours are shown in Figure 2, and are labelled as
 2 follows: *F0 fear*, *F0 joy*, *F0 sadness*, *F0 boredom*, *F0 indignation*, and *F0 neutral*.

3



4

5 Figure 2 (color online). f_0 contours used in the F_0 stimuli. The lowest contour shows the
 6 f_0 contour used with the lax-creaky voice quality (see text).

7

8 It should be noted that in selecting these contours, and in labelling them according to
 9 the affective states associated with them in Mozziconacci's (1995) study (i.e. *F0 fear*,
 10 *F0 joy*, and so on) no assumption is being made that these f_0 contours should
 11 necessarily be associated with those particular affective states by the speakers of the
 12 languages tested in this study. Rather, they can be regarded simply as a set of f_0
 13 contours that reflect likely affect-related variations, and which encapsulate the kinds of
 14 shifts in f_0 level, range and dynamics that are often mentioned in production-based
 15 studies.

16

17 The f_0 contours have two pitch accents and the overall shape of the contours is
 18 characteristic of declaratives in all four languages selected for this study. These f_0
 19 patterns differ not in terms of the basic contour shape, but in terms of f_0 level, range and

1 dynamics – the kinds of features generally reported in studies on f_0 and emotion
2 (Bänziger and Scherer, 2005; Grandjean *et al.*, 2006). The contours for *boredom* and
3 *neutral* are characterised by a low f_0 level and relatively small f_0 excursions. *Sadness* is
4 represented by a higher f_0 level; *joy*, *fear* and *indignation* are all characterised by
5 increasingly raised f_0 levels, and by an extended dynamic range, particularly in the
6 second peak. *Indignation* has the very highest f_0 level and a rapidly falling f_0 in the final
7 accent (see Figure 2).

8 **3. VQ+F0 stimuli**

9 The stimuli of the third type were generated by combining non-modal voice qualities
10 with the affect-related f_0 contours used in the synthesis of the F0 stimuli. The
11 combination of particular voice qualities with potentially affect-related f_0 contours was
12 guided by the results of the earlier experiments (Gobl *et al.*, 2002; Gobl and Ní
13 Chasaide, 2003b) as well as by the literature, e.g., Laver (1980), Scherer (1986; 2003),
14 Juslin and Laukka (2001; 2003) – the latter is a meta-analysis of 104 studies of vocal
15 expression of emotions, a summary/review in Schröder (2001). Anger, joy and fear tend
16 to be associated with high f_0 mean, variability and range, and generally, sadness and
17 boredom are associated with relatively lower f_0 mean, variability and range. According
18 to Laver (1980), breathy voice has been associated with intimacy and creaky voice with
19 boredom. Anger and joy are characterised by narrow B1 and ‘steep’ glottal waveform
20 and an increase in high frequency spectral energy (suggesting tense phonation), while
21 fear and sadness by broad B1 and ‘rounded’ glottal waveform (suggesting laxer,
22 breathier phonation) (Scherer, 1986; Juslin and Laukka, 2003). Whispery voice was
23 found effective in signalling fear in Gobl and Ní Chasaide (2003b); whispery voice
24 combined with high dynamic f_0 was found in the expression of fear in a production
25 study by Klasmeyer and Sendlmeier (2000).

1

2 Thus, *Whispery Voice* was combined with the *F0 fear* contour (referred to as the
3 *Whispery+F0 fear* stimulus), *Breathy Voice* with *F0 sadness* (*Breathy+F0 sadness*),
4 *Lax-creaky Voice* with *F0 boredom* (*Lax-creaky+F0 boredom*). *Tense Voice* was
5 combined both with *F0 joy* (*Tense+F0 Joy*) and with *F0 indignation* (*Tense+F0*
6 *indignation*).

7

8 As constructed, a potential anomaly arose for the VQ and VQ+F0 stimuli which had a
9 lax-creaky voice quality. The originally constructed *Lax-creaky* stimulus in the VQ
10 series had an f_0 level which deviated considerably from neutral being 30 Hz lower,
11 whereas the combined stimulus *Lax-creaky+F0 boredom* had an f_0 level very close to
12 the neutral. On listening to these two stimuli, it was clear that the lax-creaky quality
13 was fully audible in the stimulus even where f_0 was close to the neutral value. Insofar as
14 the VQ+F0 stimuli are intended to test for the effect of large f_0 deviations from neutral
15 coupled to particular voice qualities (as compared to voice quality alone), it would
16 appear that these two stimuli are the ‘wrong way around’. To remedy this anomalous
17 situation, and to make the presentation of results more straightforward, in the following
18 sections, these two stimuli are transposed, i.e. the *Lax-creaky+F0 boredom* stimulus is
19 treated as belonging to the VQ series of stimuli (as all stimuli in this series entail only
20 minor f_0 deviations from the neutral) and it will be referred to henceforth as *Lax-*
21 *creaky*). The *Lax-creaky* stimulus with the large f_0 lowering is more appropriately
22 treated as belonging to the VQ+F0 stimuli and is renamed as *Lax-creaky+F0 lowered*.
23 This re-classification allows one to assess more directly the contribution of the f_0
24 lowering to the affective signalling of this voice quality, in a way that retains a certain
25 consistency with the overall differentiation of the three stimulus types used here.

1
2 Table I provides a summary of the stimuli used in the perception experiments of this
3 cross-language study. Based on the voice quality and f_0 contour combinations in the
4 VQ+F0 series of stimuli, it is convenient to group the stimuli into five Stimulus
5 Groups, each of which containing one stimulus from each of the three stimulus types.
6 The five groups are named according to the voice quality and f_0 contour that
7 characterise them. As mentioned in the discussion of the F0 stimuli, despite the use of
8 labels such as ‘sadness’, there is no strong expectation that a stimulus thus labelled will
9 yield the highest ratings for the particular emotion, e.g., sadness. In addition to the five
10 groups, the stimulus with modal voice and the neutral f_0 contour provides a baseline to
11 which responses for all other stimuli can be compared.

12

13

14

15

16 Table I. The synthesised stimuli used in the cross-language study. All F0 stimuli have
17 modal voice. All VQ stimuli have neutral F0. Additionally, the stimulus *Modal+F0*
18 *neutral* is included for baseline comparison. For the stimuli with * see text.

Stimulus Group	Stimulus Type (type of manipulation)		
	VQ	F0	VQ+F0
WHISPERY FEAR	<i>Whispery</i>	<i>F0 fear</i>	<i>Whispery+F0 fear</i>
BREATHY SADNESS	<i>Breathy</i>	<i>F0 sadness</i>	<i>Breathy+F0 sadness</i>
LAX-CREAKY BOREDOM	<i>Lax-creaky*</i>	<i>F0 boredom</i>	<i>Lax-creaky+F0 lowered*</i>
TENSE JOY	<i>Tense</i>	<i>F0 joy</i>	<i>Tense+F0 joy</i>
TENSE INDIGNATION	<i>Tense</i>	<i>F0 indignation</i>	<i>Tense+F0 indignation</i>

Baseline Stimulus: *Modal+F0 neutral*

1

2 **B. Listening tests**3 **1. Procedure**

4 The perception experiments were conducted according to the procedure described in
5 Gobl and Ní Chasaide (2003b), as a series of six subtests. In each subtest, the 15 stimuli
6 were presented to the participants 10 times in random order and responses were
7 obtained for a pair of opposite affective attributes (e.g., *happy-sad*). The stimuli were
8 presented to the listeners in a quiet room through a high quality speaker with the
9 volume set at the comfortable listening level. Each subtest took approximately 11
10 minutes. The subjects had a short training session before the tests and were given
11 breaks between the subtests.

12 **2. Rating scales**

13 The participants were asked to judge the stimuli on six bipolar scales defined with
14 contrastive adjectives at each end: *indignant-apologetic*, *interested-bored*, *formal-*
15 *intimate*, *stressed-relaxed*, *happy-sad* and *fearless-scared*. A similar approach where
16 stimuli were rated for affective content on bipolar scales was used, for example, in
17 Uldall (1964) and Ladd *et al.* (1985). This kind of semantic differential scale is
18 commonly used in the study of attitude (Heise, 1970; Osgood *et al.*, 1975; Russell and
19 Carroll, 1999; Streiner and Norman, 2008), and allows one to measure directionality of
20 reaction (e.g., bored vs. interested) as well as intensity (slight to extreme). The scale is
21 usually interpreted as a seven point scale where the neutral attitude (or in our case, ‘no
22 affective colouring’) is assigned the value of zero (Heise, 1970). The data obtained by

1 this kind of ‘summative response scale’ (Gamst *et al.*, 2008) can be analysed
2 statistically using a general linear model, e.g., ANOVA.

3

4 The participants were instructed to judge each stimulus for the presence and strength of
5 affect, and to mark their response on the answer sheet where the opposite affective
6 labels were placed on each side with seven boxes in between. The choice of the centre
7 box implied that no affective colouring was present in the utterance; checking the boxes
8 to the left or right to the centre box indicated the presence and strength of a particular
9 affect, the most extreme ratings being further from the centre box. The scale was later
10 interpreted as ranging from -3 (low activation states) to $+3$ (high activation states),
11 where 0 value corresponded to no affect perceived, ± 1 corresponded to mild affective
12 colouring, ± 2 – to moderate and ± 3 – to strong affective colouring. Numeric values
13 were not presented to the raters during the experiment.

14

15 The affective labels defining the opposite ends of each of the six scales have been
16 chosen to cover a fairly broad range of emotions and milder affective states such as
17 attitudes and interpersonal stances. The affective adjectives used in the scales are
18 among those most frequently found in the lists of affect-related words (Juslin and
19 Laukka, 2003). The choice of affective labels was in part guided by the synthesised
20 voice qualities and by what is known about voice-to-affect mapping in the literature.
21 The scales *stressed-relaxed*, *happy-sad*, *interested-bored*, *formal-intimate* were used in
22 prior work and were adopted here in part to assure comparability of results.

23 **3. Translation of affective labels**

24 The stimuli were presented to speakers of four languages - Irish English, Russian,
25 Spanish and Japanese. For the speakers of languages other than English the affective

1 labels on the answer sheets were given in their respective languages in translation. To
2 render the translations as accurately as possible, we used team translation or the
3 ‘committee approach’ that also involved back translation as recommended in Brislin
4 (1980) and Streiner and Norman (2008). The translation of the labels from English was
5 undertaken by at least two native speakers of the respective languages (university
6 students and staff) who had a good command of English and who were also familiar
7 with the nature of the research and the purpose of the rating scales. As there are usually
8 several translation possibilities, the translators were asked to discuss them and come to
9 consensus regarding the best possible choice. The affective labels used in the listening
10 tests are shown in Table I in supplementary material¹.

11

12 **4. *Participants***

13 The participants were speakers of Irish English, Russian, Spanish and Japanese. As it
14 was desirable to exclude the influence of semantic content of the Swedish utterance
15 used as the basis for the synthesised stimuli on the perception of affect, we made sure
16 prior to the experiments that none of the participants spoke Swedish. It is virtually
17 impossible in the modern world to completely exclude the influence of foreign
18 languages and with this the influence of language-specific ‘foreign’ voice qualities, but
19 such an influence was kept to a minimum in the present experiments as none of the
20 participants had lived in a foreign country for any length of time or used foreign
21 languages in an active way in their professional life (e.g., as an interpreter or a language
22 teacher). The participants were college students recruited in Dublin (n = 20, 13 female,
23 aged 18-35), St. Petersburg (n = 21, 10 female, aged 19-35), Barcelona (n = 20, 10
24 female, aged 18-22) and Tokyo (n = 21, 10 female, aged 18-22).

1 **C. Statistical analyses**

2 To compare the ratings of the three types of stimuli (the VQ stimuli relative to the
3 VQ+F0 and the F0 stimuli) within each stimulus group taking the language factor into
4 account, a separate repeated measures ANOVA was conducted for each subtest using
5 the SPSS software package (IBM Corp., 2013).

6
7 A complex mixed design was applied. The two within-subjects factors were the
8 Stimulus Group (5 levels: WHISPERY||FEAR, BREATHY||SADNESS, LAX-
9 CREAKY||BOREDOM, TENSE||JOY, TENSE||INDIGNATION) and Stimulus Type (3 levels: VQ,
10 F0, VQ+F0). The between-subjects factor was Language (4 levels: Irish English (E),
11 Russian (R), Spanish (S), Japanese (J)). The dependent measure was mean affective
12 rating (averaged across 10 randomisations for each speaker).

13
14 As *Tense Voice* was used twice, in the generation of two VQ+F0 stimuli in two
15 Stimulus Groups, TENSE||JOY (*Tense+F0 joy*) and TENSE||INDIGNATION (*Tense+F0*
16 *indignation*), it was not possible to apply the $5 \times 3 \times 4$ factorial design to the data. In
17 each subtest, the analysis of variance was therefore done twice as a $4 \times 3 \times 4$ design (4
18 Stimulus Groups \times 3 Stimulus Types \times 4 Languages), first for the stimuli in the
19 Stimulus Groups WHISPERY||FEAR, BREATHY||SADNESS, LAX-CREAKY||BOREDOM and
20 TENSE||JOY (Part I of the ANOVA) and second, for the stimuli in the Stimulus Groups
21 WHISPERY||FEAR, BREATHY||SADNESS, LAX-CREAKY||BOREDOM and TENSE||INDIGNATION
22 (Part II of the ANOVA). Modal voice was not included in the model.

23
24 To assess the ratings of VQ stimuli relative to *Modal* and the ratings of F0 and VQ+F0
25 stimuli relative to *Modal* (with neutral f_0), a separate repeated measures ANOVA test of

1 simple contrasts was conducted for each subtest in which pairwise comparisons of all
2 the stimuli were carried out. The within-subject factor was Stimulus, the between-
3 subject factor was Language. The dependent variable was the mean affective rating for
4 each participant (averaged across 10 randomisations). To correct for multiple
5 comparisons and to minimise the risk of random effects, a Bonferroni correction was
6 applied in the post-hoc tests.

7

8 There are multiple sources of variability in voice stimuli ratings, such as listener factors
9 (sensitivity, bias, error, fatigue, etc.), scale factors (scale resolution and the way the
10 scale is defined), stimulus factors (the quality of voice samples or synthesised stimuli),
11 as well as the interaction of these factors (Kreiman *et al.*, 2007). Consistency in
12 assignment of similar ratings to the same stimulus will suggest stronger voice-to-affect
13 association. Listeners' agreement/consistency in ratings was measured by Intraclass
14 Correlation Coefficients (ICC), the two-way random model (Landis and Koch, 1977;
15 Shrout and Fleiss, 1979; McGraw and Wong, 1996), which were calculated for each
16 subtest for each language group of listeners. As it is of interest to establish whether the
17 judgment of one rater is similar to that of the others, the single measures ICC (r) rather
18 than the average measures ICC (R) will be considered here. The interpretation of ICC is
19 based on Landis and Koch (1977): < 0 poor agreement, 0-0.20 slight, 0.21-0.40 fair,
20 0.41-0.60 moderate, 0.61-0.80 substantial, 0.81-1 almost perfect agreement.

21

22 IV. RESULTS

23 The statistical results are outlined and following a broad overview the detailed findings
24 for the individual stimuli are summarised.

1 **A. Statistical results**

2 **1. Rater's agreement**

3 The participants' agreement, calculated as single measure (r) Intraclass Correlation
4 Coefficients, overall was within the moderate to substantial range (moderate in 12 out
5 of 24 cases, substantial in 6 out of 24 cases). Two groups of listeners, Japanese and
6 Spanish, showed fair agreement in the *formal-intimate* and *fearless-scared* subtests.
7 Spanish listeners also showed fair agreement in the *indignant-apologetic* subtest.
8 Furthermore, the agreement of the Japanese listeners was poor in the *stressed-relaxed*
9 subtest. Going from the greatest to the lowest degree of interrater agreement, the order
10 is $R > E > S > J$ (see supplementary material, Table II²).

11

12 **2. ANOVA results**

13 The results of the omnibus complex design ANOVA are given in supplementary
14 material, Table III³. The results for the within-subject effects - part (a) of the ANOVA -
15 showed significant effects of Stimulus Type and Stimulus Group in all subtests. Two-
16 way Stimulus Group \times Language and Stimulus Type \times Language interactions were
17 largely significant, except in *happy-sad* and *fearless-scared* subtests. A three-way
18 Stimulus Group \times Stimulus Type \times Language interaction was found significant in all
19 but two subtests, *indignant-apologetic* (Part I) and *happy-sad* (Part II). Given the
20 interactions, it is not possible here to assess the relative contribution of individual
21 factors. Thus, for a closer look at cross-language differences, the pairwise comparisons
22 of the individual stimuli are also included.

23

1 In part (b) of the ANOVA, the between-subject effect of Language was found to be
 2 significant in the two subtests: *formal-intimate* and *stressed-relaxed*. It is for these
 3 affects, therefore, that we would expect to see the greatest differences among languages
 4 in terms of voice to affect association.

5
 6 Details of statistically significant cross-language differences of the ratings for the
 7 individual stimuli obtained in the ANOVA test are shown in Figure 3 (explained in
 8 detail below). Pairwise comparisons of within-language differences in ratings for the
 9 different stimuli are not presented here due to space constraints; note that the
 10 differences in affective ratings for *Modal* and *Modal+F0 boredom* and for *Lax-creaky*
 11 *Voice* and *Lax-creaky+F0 lowered* were not found to be statistically significant in any
 12 language for any affective subtest.

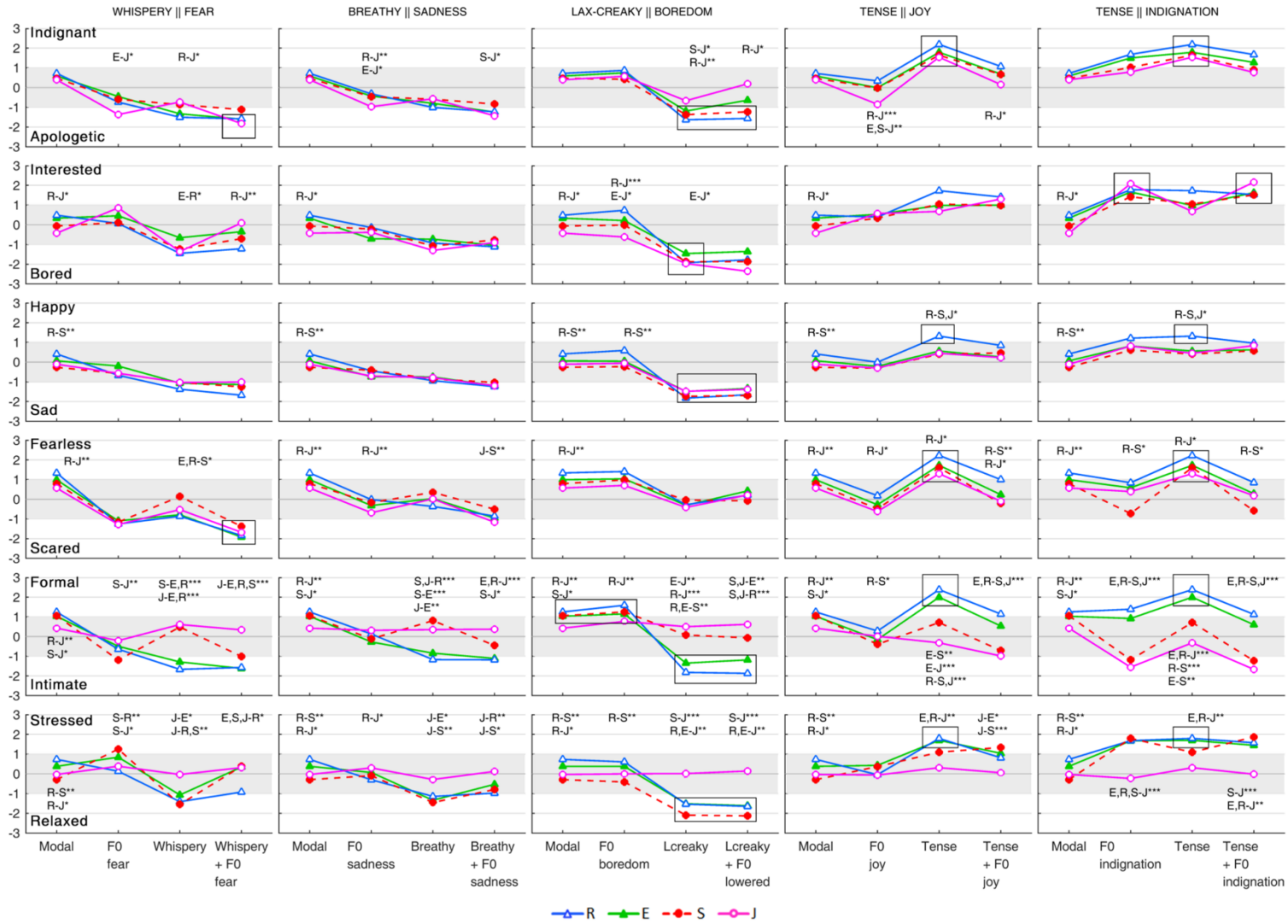
13 **B. Voice-to-affect mapping: broad overview**

14 The five columns in Figure 3 show the ratings obtained organised in terms of the
 15 Stimulus Groups (see also Table I). Each row presents plots with the ratings for one of
 16 the six affective subtests (*indignant-apologetic, etc.*). In each plot, the rating for the
 17 stimulus *Modal Voice (+F0 Neutral)* provides a baseline to which ratings for other
 18 stimuli can be compared. For each affective subtest, the most highly rated stimulus is
 19 marked by a surrounding box. More than one box is used when two stimuli produce
 20 similar ratings or where there are differences among the languages in terms of their
 21 most highly rated stimulus.

22

23 In the individual plots positive values (0 to +3) correspond to high activation/power
 24 states (*indignant, interested, happy, fearless, formal, stressed*), while negative values (0

1 to -3) correspond to low activation/power states (*apologetic, bored, sad, scared,*
2 *intimate, relaxed*). As in previous studies (e.g., Yanushevskaya *et al.*, 2013) discussion
3 focuses on ratings that exceed ± 1 , and this region of weak affect is shaded grey. This
4 ± 1 threshold is admittedly arbitrary and indeed, statistically significant differences can
5 occasionally be found between ratings that are low, but the intention here is to
6 concentrate on the more salient voice-to-affect associations. Significant cross-language
7 differences are indicated by asterisks (* $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$).



- 1 Figure 3 (color online). Affective ratings obtained in the six subtests: each row represents a listening test (e.g., *indignant-apologetic*), each
- 2 column represents a particular stimulus group (e.g., WHISPERY||FEAR). Languages: Irish-English (E); Russian (R), Spanish (S), Japanese (J).
- 3 Statistically significant cross-language differences: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. The most highly rated stimulus for a particular affect
- 4 is marked by a surrounding box.
- 5

1 A broad indication of the extent of cross-language differences/similarities can be
2 gleaned from Figure 3. Two of the affective subtests stand out as having large
3 statistically significant differences. In the *formal-intimate* subtest, languages patterned
4 in terms of two groups, R and E in one, and S and J in the other. In the *stressed-relaxed*
5 subtest, J diverges from the other three languages. In the remaining subtests results
6 appear rather similar overall in terms of voice to affect association, although there are
7 many cases (indicated with asterisks in Figure 3) where the same stimulus yields
8 affective ratings of significantly different strength from listeners of different language
9 groups.

10

11 Of the five Stimulus Groups, TENSE||INDIGNATION accounts for the highest ratings
12 obtained for high activation/power states. Within this group, the VQ stimulus (*Tense*
13 *Voice*) is particularly effective. The F0 stimulus (*F0 indignation*) is strongly associated
14 with *interested* across languages but shows up prominent cross-language differences in
15 its association with other affects – differences which persist in the combined VQ+F0
16 stimulus (*Tense Voice+F0 indignation*).

17

18 The LAX-CREAKY||BOREDOM group is the most effective for signalling low
19 activation/power states. Within this group, the *Lax-creaky Voice* achieves the highest
20 ratings overall, though as mentioned ratings are virtually identical to those obtained for
21 the *Lax-creaky+F0 lowered* stimulus.

22

23 The group WHISPERY||FEAR also appears to be quite effective in signalling specific low
24 activation/power states. The VQ+F0 stimulus *Whispery+F0 fear* emerges across the
25 board as the most effective in conveying the affect *scared*.

1
2 The two stimulus groups, BREATHY||SADNESS and TENSE||JOY, are here relatively
3 ineffective in signalling affect. (In the latter group *Tense Voice* achieves high ratings
4 but note that this appears twice, being also represented in the TENSE||INDIGNATION
5 group.) These two groups are largely omitted from the further discussion of results.

6 **C. Voice-to-affect mapping across languages**

7 Table II presents an overview of all stimuli which achieved ratings above the ± 1
8 threshold. Stimuli associated with high-activation states are grouped in the upper part of
9 the table and low activation states in the lower. Bold font identifies those stimuli which
10 yielded the highest rating for a given affect.

11

12

13

1 Table II. Stimuli that obtained the highest ratings (absolute values shown as subscripts) within each stimulus type for the affects tested. Stimuli
 2 yielding the highest ratings across stimulus types for at least one language are shown in bold type.

Affect	VQ	F0 (all with Modal Voice)	VQ+F0
<i>indignant</i>	Tense <i>R</i> 2.18, <i>E</i> 1.78, <i>S</i> 1.64, <i>J</i> 1.53	F0 indignation <i>R</i> 1.68, <i>E</i> 1.50, <i>S</i> 1.02	Tense+F0 indignation <i>R</i> 1.67, <i>E</i> 1.27
<i>interested</i>	Tense <i>R</i> 1.72, <i>S</i> 1.05	F0 indignation <i>J</i> 2.08, <i>R</i> 1.78, <i>E</i> 1.66, <i>S</i> 1.42	Tense+F0 indignation <i>J</i> 2.15, <i>E</i> 1.62, <i>R</i> 1.53, <i>S</i> 1.51
<i>formal</i> ^a	Tense <i>R</i> 2.38, <i>E</i> 1.99 Modal <i>S</i> 1.05	F0 boredom <i>R</i> 1.58, <i>S</i> 1.25, <i>E</i> 1.14	Tense+F0 joy <i>R</i> 1.43
<i>stressed</i> ^b	Tense <i>R</i> 1.79, <i>E</i> 1.70, <i>S</i> 1.09	F0 indignation <i>S</i> 1.79, <i>R</i> 1.68, <i>E</i> 1.67	Tense+F0 indignation <i>S</i> 1.87, <i>R</i> 1.57, <i>E</i> 1.44
<i>happy</i>	Tense <i>R</i> 1.31	F0 indignation <i>R</i> 1.21	
<i>fearless</i>	Tense <i>R</i> 2.22, <i>E</i> 1.72, <i>S</i> 1.61, <i>J</i> 1.30	F0 boredom <i>R</i> 1.41, <i>E</i> 1.04	Tense+F0 joy <i>R</i> 1.00
<i>apologetic</i>	Lax-creaky <i>R</i> 1.63, <i>S</i> 1.37 Whispery <i>E</i> 1.33	F0 fear <i>J</i> 1.37	Lax-creaky+F0 lowered <i>R</i> 1.57, <i>S</i> 1.23 Whispery+F0 fear <i>J</i> 1.82, <i>E</i> 1.60
<i>bored</i>	Lax-creaky <i>J</i> 1.96, <i>R</i> 1.92, <i>S</i> 1.87, <i>E</i> 1.46		Lax-creaky+F0 lowered <i>J</i> 2.35, <i>S</i> 1.87, <i>R</i> 1.77, <i>E</i> 1.35
<i>intimate</i>	Lax-creaky <i>R</i> 1.81 Whispery <i>E</i> 1.29	F0 fear <i>S</i> 1.19 F0 indignation <i>J</i> 1.55	Lax-creaky+F0 lowered <i>R</i> 1.87 Whispery+F0 fear <i>E</i> 1.62 Tense+F0 indignation <i>J</i> 1.66, <i>S</i> 1.23
<i>relaxed</i> ^c	Lax-creaky <i>S</i> 2.09, <i>R</i> 1.53, <i>E</i> 1.51		Lax-creaky+F0 lowered <i>S</i> 2.13, <i>R</i> 1.64, <i>E</i> 1.61
<i>sad</i>	Lax-creaky <i>R</i> 1.82, <i>S</i> 1.74, <i>E</i> 1.49, <i>J</i> 1.49		Lax-creaky+F0 lowered <i>S</i> 1.69, <i>R</i> 1.65, <i>J</i> 1.38, <i>E</i> 1.34
<i>scared</i>		F0 fear <i>J</i> 1.27, <i>R</i> 1.25, <i>S</i> 1.16, <i>E</i> 1.09	Whispery+F0 fear <i>E</i> 1.91, <i>R</i> 1.83, <i>J</i> 1.67, <i>S</i> 1.37

3

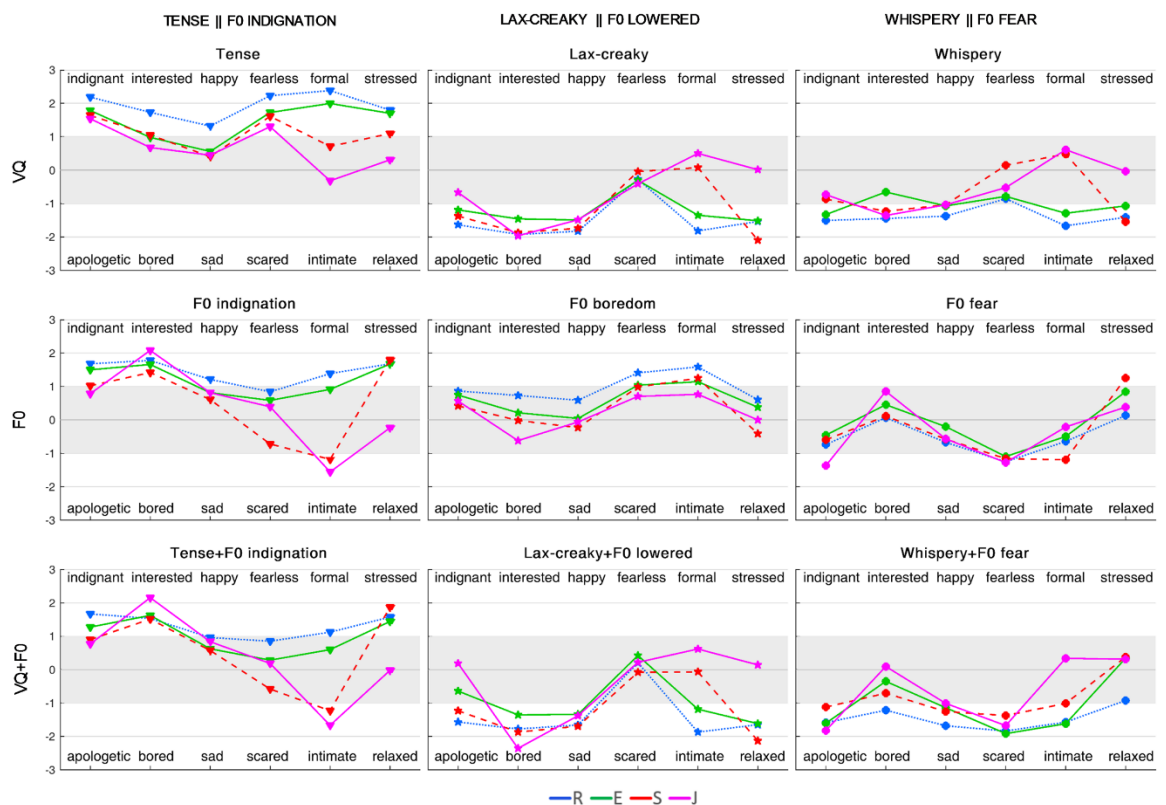
^a None of the stimuli signalled *formal* to Japanese.

^b None of the stimuli signalled *stressed* to Japanese.

^c None of the stimuli signalled *relaxed* to Japanese.

1 Overall, Table II shows that the stimuli most effective in signalling affect are indeed
 2 those which entail voice quality differences, with or without concomitant f_0 shifts (i.e.
 3 the VQ and VQ+F0 series). As expected, *Tense Voice* is particularly associated with
 4 high activation states. *Lax-creaky Voice* is effective in signalling low activation affects,
 5 as is (unsurprisingly) *Lax-creaky+F0 lowered* and the *Whispery+F0 fear* stimuli.
 6 Although *Tense+F0 indignation* also achieves high affective ratings, rather
 7 unexpectedly, it gets associated with both positive and negative affects. The F0 stimuli
 8 were comparatively ineffective on the whole, but the *F0 indignation* stands out as an
 9 exception, being strongly associated with affect.

10

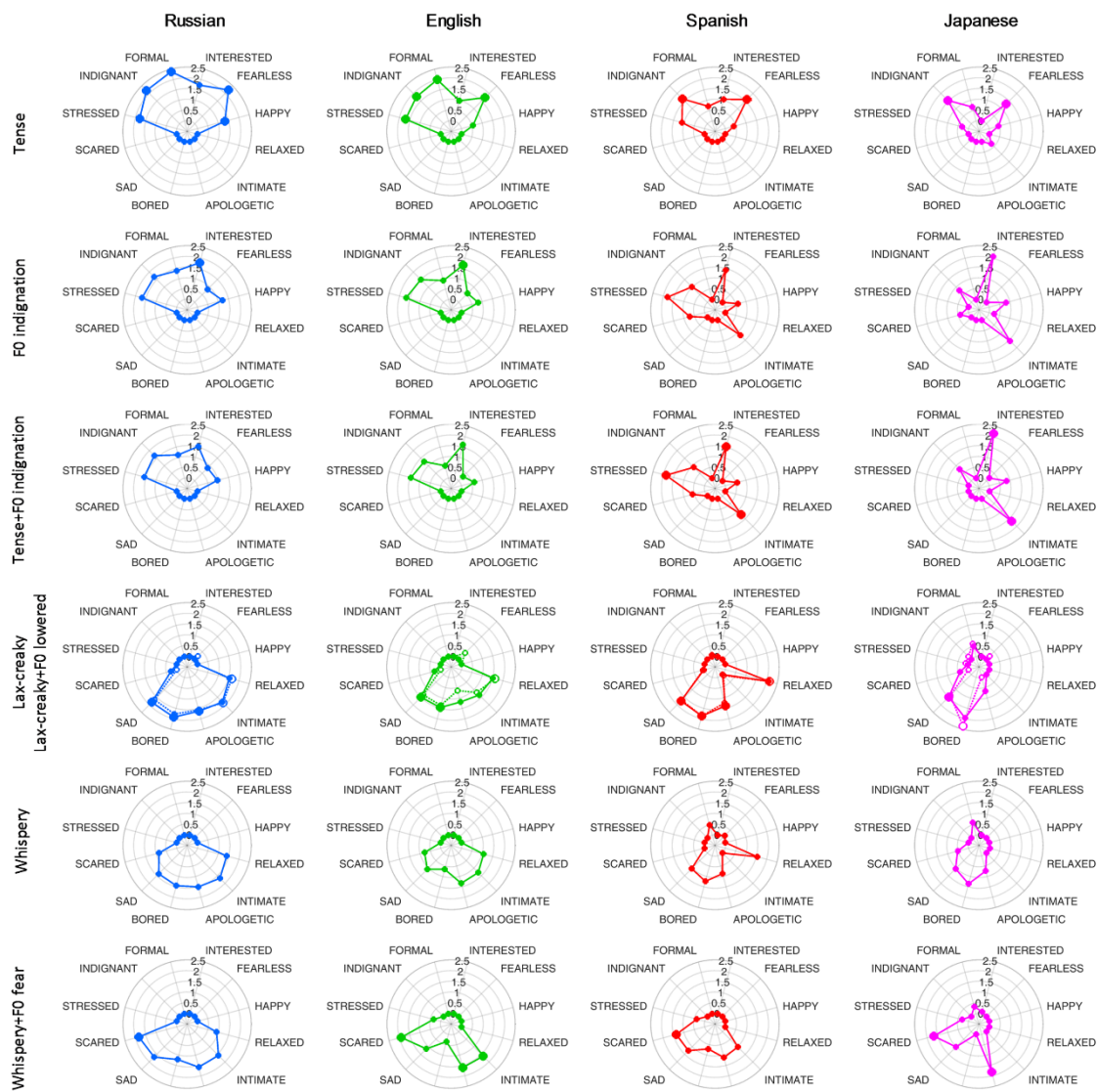


11

12 Figure 4 (color online). Ratings for the three stimulus types, from the three most
 13 effective stimulus groups in each language.

14

1 Figures 4 and 5 allow a closer look at how the individual stimuli map to affect. In
2 Figure 4 results are shown for the subset of stimuli that feature in the three most
3 effective Stimulus Groups of Figure 3, discussed above (TENSE||INDIGNATION, LAX-
4 CREAKY||BOREDOM, WHISPERY||FEAR). Here, for each stimulus, ratings for each
5 affective subtest are shown for the four languages. The spidergram summary plots in
6 Figure 5 show, for the most highly rated stimuli, the network of associated affects in
7 each language. Enlarged data points identify where the highest rating for a particular
8 affect has been obtained (across all stimuli). The affective states have been loosely
9 arranged so that high activation states are in the upper part of the spidergram, and low
10 activation states in the lower. Positive affects are located to the right, and negative to
11 the left.
12



1

2 Figure 5 (color online). Summary plots of voice-to-affect associations for the most
 3 highly rated stimuli. The data for *Lax-creaky* (solid line) and *Lax-creaky+F0 lowered*
 4 (dotted line) stimuli are superimposed. Larger data points show cases when a stimulus
 5 yielded highest ratings compared to other stimuli for a particular affect.

6 1. *The VQ stimuli*

7

8 ***Tense voice:*** *Tense Voice* did emerge as the stimulus most frequently associated with
 9 high activation states. One high activation states *happy* stands out as being only
 10 moderately signalled – and only for a single language, Russian. No other stimulus cues

1 *happy* for any language in this study. The elusive nature of vocal correlates of *happy* is
2 reported in many studies.

3

4 Although results for *Tense Voice* conform to our initial expectations based on the
5 ‘effort code’ and to our earlier findings for English, *Tense Voice* does not appear to
6 carry the same affective load across these languages. In all four languages it elicits the
7 highest ratings for *indignant* and *fearless*. Other than that, there are differences: as can
8 be seen in Figures 4 and 5, it is strongly associated with a wide range of high
9 activation states for Russian, somewhat less strongly in English, while for Japanese
10 and Spanish subjects, the affect ratings are much lower, with many more ‘gaps’ in
11 coverage, e.g., in the cueing of *stressed* and *formal*. (Note that for these two languages
12 *formal* is not strongly signalled by any of the present stimuli, and for J, *stressed* is not
13 rated at all). To sum up, although *Tense Voice* is associated with high activation states
14 in all these languages, the range (number) of affects signalled and the strength of
15 signalling varies considerably – in the order $R > E > S > J$.

16

17 ***Lax-creaky Voice***: (Given the similarity of their ratings, the values for *Lax-creaky*
18 *Voice and Lax-creaky+F0 lowered* have been superimposed in Figure 5.) From the
19 outset, the expectation was to find this quality strongly associated with low activation
20 states. As with *Tense Voice*, the range/number of affects associated varied considerably
21 across the four languages. Again, the most wide-ranging effects are attested for
22 Russian and the least for Japanese.

23

24 This quality yields the highest ratings across the board for both *bored* and *sad*. There
25 are also, however, large cross language differences. *Intimate* is rather well signalled by

1 this quality for Russian and English, but not at all in Japanese or Spanish. The affects
2 *relaxed* and *apologetic* are also associated with this quality for Russian, Spanish, and
3 (more weakly) English, but not at all for Japanese.

4

5 ***Whispery Voice***: As can be seen in Figures 4 and 5, this quality does achieve
6 reasonably high affective ratings for many low-activation states, and its affective
7 profile (Figure 5) is somewhat similar to *Lax-creaky Voice* but more attenuated. The
8 cross-language differences mirror those for *Lax-creaky Voice*, especially for the affects
9 *intimate* and *relaxed*. The initial expectation that whispery voice might be associated
10 with fear is not borne out: while R and E responses for *scared* show weak signalling of
11 this affect (not indicated for the otherwise similar *Lax-creaky Voice*), ratings do not
12 exceed ± 1 threshold.

13

14 ***Breathy Voice***: Overall, this stimulus did not yield strong affective ratings (see Figure
15 3). Although this quality is frequently mentioned as associated with intimacy (e.g.,
16 Laver, 1980) and with sadness (e.g., Juslin and Laukka, 2003) but this did not manifest
17 in a clear way. The suggested association with formality in Japanese (Campbell and
18 Mokhtari, 2003) was not borne out here either.

19

20 **2. *The F0 stimuli***

21 Overall, stimuli involving only f_0 manipulation were relatively ineffective in affect
22 cueing (see Figure 3), and so these stimuli were not associated with the affects used in
23 their labels. This was largely expected, given past studies of f_0 contours in affect
24 perception (see Section II.C). In the Introduction, a possible expectation was discussed,
25 based on the ‘effort code’, i.e. that the strength of activation would mirror the extent to

1 which a particular f_0 contour deviated from the neutral f_0 . Extremely elevated f_0
2 contours might be expected to be associated with high activation states, and lowered f_0
3 contours with low activation states. (The very wide range of f_0 levels and the extreme
4 nature of some of them can be seen Figure 2). The fact that most extreme f_0 contour in
5 the series, *F0 indignation*, did yield strong affective responses could be seen as lending
6 support to this proposition. However, the picture is rather more complex, as is
7 discussed further below.

8
9 ***F0 indignation***: This most extreme f_0 contour stands out as yielding high affective
10 ratings. Across all four languages, this stimulus was the most highly rated for
11 *interested*. Whereas one might have expected this contour to signal *indignant*, such an
12 association only emerged for R and E (though with ratings that were lower than for
13 *Tense Voice*). Overall, the ratings for R and E are rather similar, and show a broad
14 association with high activation states (similar to *Tense Voice* but with much weaker
15 rating strength).

16 In Japanese and Spanish this stimulus was rather unexpectedly associated with
17 *intimate*, an association entirely absent in ratings for Russian and English. Thus, the
18 expectation that extreme f_0 deviation from the neutral would be associated with high
19 activation does not hold across these languages. Though it is true for R and E, the fact
20 that both high and low activation states are signalled in both J and S suggests that the
21 ‘effort code’ explanation is not generalizable. In addition, as can be gleaned in
22 Figure 3, for the other F0 stimuli, there was little evidence of a correlation between the
23 degree of f_0 deviation from the neutral, and degree of activation in responses.

1 **3. The combined VQ+F0 stimuli**

2 The combined stimuli were included to explore whether the combination of voice
3 quality with differing f_0 contours might yield stronger affective responses than when
4 voice quality alone differs (or f_0 alone). If, as many researchers have concluded, f_0
5 contours alone are ineffective in signalling affect because the voice quality dimension
6 is missing, might not the combination of voice quality with f_0 contours provide a
7 synergy to yield the most effective signalling? The present results did reveal synergies
8 but rather less than might be expected. Affective ratings for three combined stimuli did
9 achieve higher ratings than the VQ or F0 series alone: however, the gain was marginal
10 in two of these three cases.

11

12 **Whispery Voice+F0 fear:** This proved to be the most clear-cut case of synergy, where
13 the combination of voice quality and f_0 features was more effective than either
14 dimension on its own. The *F0 fear* contour yielded relatively weak affective ratings
15 overall but slightly above threshold values for *scared* (see Figure 4). Although
16 *Whispery Voice* achieves stronger affect signalling on the whole, it is still much less
17 effective than *Lax-creaky Voice*, for all affects other than *scared*, which is only weakly
18 indicated. However, when *Whispery Voice* is combined with the *F0 fear* contour there
19 is a distinct enhancement of affective signalling (see Figure 5). This combined
20 stimulus yields the highest ratings for *scared* across all four languages; it yields the
21 highest ratings for *apologetic* in E and J and the strongest signalling of *intimate* for E.

22

23 **Tense Voice+F0 indignant:** As *Tense Voice* is associated with high-activation states,
24 ratings were expected to be higher for the combined stimulus. It was also thought
25 likely that the affect *indignant* would emerge strongly.

1

2 These expectations were not borne out, and trends differed according to the language.

3 In Spanish and Japanese, those affects which were associated with *F0 indignant*

4 (*interested, intimate*, and in the case of Spanish, *stressed*) did yield higher ratings

5 when in combination with *Tense Voice*, showing a synergy, even if the extent of the

6 enhancement involved is not great. However, results for Russian and English run

7 contrary to this in that the combined stimulus yields lower affect ratings than either of

8 the VQ or F0 stimuli alone.

9

10 Responses for Spanish and Japanese are overall very like those obtained for the

11 *F0 indignation* stimulus, and very different from those obtained for *Tense Voice*. This

12 leads us to conclude that the extreme f_0 contour is the main determinant of the affects

13 cued by the combined stimulus for these two languages.

14

15 Taking the results for this combined stimulus together with results the *F0 indignation*,

16 we would conclude that the f_0 dimension works differently in the affect signalling of

17 these two language groups, and that it plays a more important role in Japanese and

18 Spanish than in Russian and English.

19

20 These results for the combined stimulus also prompt reflections on the role tense

21 phonation plays in the affect signalling of these two language groups. The fact that a

22 tense voice quality is compatible with, and (even slightly) enhances the impression of

23 *intimate* in Spanish and Japanese underscores a basic difference in how this voice

24 quality functions in the affect-system of the two language groups. While

1 unambiguously associated with high activation states in Russian and English, there is
2 no such necessary linkage in Japanese or Spanish.

3

4 **Lax-creaky+F0 lowered:** although this combined stimulus yielded high affective
5 ratings, these were almost identical to (and never significantly different from) those for
6 *Lax-creaky Voice*. This strongly suggests that, although lax-creaky voice quality is
7 often produced with low f_0 , its affective signalling role appears to be due to the voice
8 quality components other than the low pitch. Although f_0 lowering on its own (with
9 modal phonation) was not included in the stimulus set here, we would tentatively
10 conclude from the results from this combined stimulus that f_0 lowering is not greatly
11 implicated in affect signalling.

12

13 V. DISCUSSION

14 This study set out to (1) explore how the different dimensions of the voice (voice
15 quality and f_0) convey affect and provide insight into whether/how they combine in
16 affect signalling, and (2) to explore the cross-language differences in how affect may
17 be associated with voice.

18 A. How the dimensions of the voice signal affect

19 On the whole, results suggest that voice quality dominates the signalling of affect. As
20 in earlier experiments by the authors, stimuli involving VQ – either alone or in
21 combination with a specific f_0 contour – account for most of the highest affect ratings
22 found, while most of F0 series of stimuli were relatively ineffective. However, the high
23 affective ratings for *F0 indignation*, especially in J and S, and the clear synergy of f_0
24 contour and voice quality in the combined stimulus *Whispery+F0 fear* are indicators

1 that f_0 can be very important factor in conveying affect. The fact that the affective
2 profiles of J and S are rather different from R and E for the stimuli *F0 indignation* and
3 *Tense+F0 indignation* (Figure 5) provides an indication that (i) the relative
4 contribution of either dimension is likely to be variable across languages, and that (ii)
5 the ways they combine are consequently likely to differ considerably.

6

7 **An ‘effort code’?** An expectation was expressed in the Introduction that the signalling
8 of activation would mirror degree of underlying laryngeal effort (akin to the ‘effort
9 code’ proposed to account for quasi-universal trends in intonation by Gussenhoven,
10 2004). This suggestion received support from earlier studies for English, where *Tense*
11 *Voice* was clearly associated with high activation states, while *Lax-creaky Voice* was
12 associated with low activation states. It was also suggested that the ‘effort code’ might
13 equally be proposed for the F0 stimuli, where those contours deviating most extremely
14 from the neutral contour would be expected to yield relatively stronger signalling of
15 high activation states.

16

17 At first glance, this trend does emerge in the results for the voice quality materials.
18 *Tense Voice* does signal high activation, while *Lax-creaky Voice* is associated with low
19 activation. However, a close look at cross language differences demands a more
20 nuanced account. *Tense Voice* is clearly associated with high activation states for R
21 and E, but the effect is more limited for J and S. Furthermore, the fact that a tense
22 voice quality, when coupled with *F0 indignation* can in the latter languages signal
23 *intimate* (enhancing the affect ratings vis à vis the f_0 contour alone) means that we
24 cannot expect a linkage to high activation to be universal.

25

1 In a similar vein, the results for the F0 stimuli show that the ‘effort code’ doesn’t hold
2 either in any simple way. Although the extreme f_0 contour *F0 indignation* does achieve
3 high affective ratings, these can simultaneously involve both high (*interested, stressed*)
4 and low activation (*intimate*) states. Also, the other very elevated f_0 contours of the F0
5 series were not more associated with high activation states than were those with lower
6 f_0 contours.

7

8 **Synergies?** The initial proposal was that the combined stimuli would yield affect
9 ratings beyond what the individual VQ or F0 stimuli achieved. This was not
10 resoundingly demonstrated in this experiment, as a clear-cut enhancement was only
11 found for one of the five combined stimuli. Nonetheless, the significant increase in
12 affect signalling in this one case, *Whispery Voice+F0 fear*, provided a clear
13 demonstration that synergies may be needed to signal certain states. As the pairings in
14 this experiment represent a limited set, the present results do not allow of strong
15 conclusions. This is an area that will need to be explored further.

16

17 **One-to-one mappings?** The present results, rather like our earlier experiments
18 indicate that there is no one-to-one mapping of voice to affect: a stimulus such as
19 *Tense Voice* or *Whispery Voice* maps to more than a single affect and may be
20 associated with a cluster of affects. Nonetheless, Figure 5 shows that the mappings are
21 more restricted in some languages than others, with more restricted in J and S than in
22 R and E. For J and S, in the signalling of *intimate* the very extreme f_0 contour appears
23 to be *the* determining factor, but even here, other affects like *interested* are also
24 signalled.

1 **B. Cross-language differences/similarities**

2 Table III provides a simple listing of those stimuli to yield the highest ratings for each
3 affect in each language, and a box surrounds cells where these corresponded across all
4 four languages. The clear-cut cases where the mapping of stimulus-to-affect
5 corresponds across all four languages are as follows:

6 *Tense Voice* → *indignant, fearless*

7 *Lax-creaky Voice* → *bored, sad*

8 *Whispery+F0 fear* → *scared, apologetic*

9 *F0 indignation and Tense voice+F0 indignation* → *interested*

10

11 The affective ‘profiles’ of Figure 5 provide a visual overview of cross-language
12 convergence and divergence. Despite the overlap of values, no single stimulus here
13 yields the same, or even similar, patterns across all four languages.

14

15

16

1

2 Table III. Summary of the stimulus to affect association for the languages tested (**R**=Russian; **E**=Irish-English; **S**=Spanish; **J**=Japanese). Only

3 the stimuli yielding the highest rating are shown.

	Indign.	Interest.	Happy	Fearless	Formal	Stressed	Apologet	Bored	Sad	Scared	Intimate	Relaxed
Tense	RESJ		R	RESJ	RE	RE						
F0 indignation		RE										
Tense+F0 indignation		SJ				S					SJ	
Lax-creaky/ Lax-creaky+F0 lowered							R S	RESJ RESJ			R	RES
Modal/F0 boredom					S							
Whispery+F0 fear							E J			RESJ	E	
Gaps			ESJ		J	J						J

4

5

1 Languages appear to differ in a number of ways:

2 **Entirely different mappings** of voice-to-affect is probably the most striking case. This
3 is exemplified particularly in the signalling of *intimate*. The strong association with the
4 *F0 indignation* stimulus, found for J and S finds no echo in the results for R and E.
5 Similarly, the strong association with *Lax-creaky Voice* for R and E is not found for J
6 and S.

7

8 **The strength** of affective responses can also be quite different across these languages.
9 For example, *Tense Voice* yields much higher affect ratings in R and E than in S or J. In
10 fact, high activation states were on the whole not well signalled for the J and S
11 listeners. It is nonetheless entirely possible that J and S may communicate high
12 activation affects by using voice and f_0 cues, but that the choice of voice qualities and f_0
13 contours of this study did not include the critical ones needed.

14

15 **The range/number** of affects associated with a particular stimulus also differs
16 considerably in these languages. This can again be illustrated by responses for *Tense*
17 *Voice*, which in R and E is associated with a wide range of high activation affects, but
18 with a more restricted set in S and J. This also holds for responses to *Lax-creaky Voice*,
19 which yields wide-ranging low-activation states for R and E, a more limited range for
20 S, and an even more restricted range for J.

21

22 **Gaps in coverage** are a related phenomenon (see lowest row in Table III), which
23 emerged particularly in responses for J, where neither *stressed* nor *relaxed* were
24 signalled at all. This was in striking contrast to the other three languages, where these
25 affects were effectively conveyed. The affect *formal* also failed to emerge for J, and

1 was only weakly signalled for S. Overall, J emerged as having the sparsest signalling
2 of affect from this range of stimuli. It has been suggested (see discussion in Section
3 II.C) that a given emotion may be expressed more clearly and recognised with higher
4 accuracy in some languages than in others. This could be the basis for some of the
5 differences observed here, but other possibilities exist, such as that just mentioned, that
6 the stimulus selection here might not be optimal for a particular language. Furthermore,
7 one must consider that other, linguistic and cultural factors might also be at play, and
8 we return to this issue in the Conclusions.

9

10 **The relative importance of VQ and F0**, and in ways in which these dimensions of the
11 voice combine in affect signalling also appears to be different. As discussed above, the
12 affective impact for J and S of the extreme pitch contour of the *F0 indignation* in the
13 signalling of *intimate*, and the way it combines with *Tense Voice* in these languages to
14 heighten the *intimate* affect lead us to propose f_0 likely plays a different, more
15 important role in J and S than in R and E.

16

17 Broadly speaking, two groupings appear to emerge here: as suggested by many of the
18 above observations and by the affective profiles of Figure 5, R and E show similar
19 trends, and these differ in many respects from the trends observed for J and S. This
20 grouping is of course not absolute: in the signalling of *relaxed* and *stressed* Spanish
21 patterns closely with Russian and English rather than with Japanese.

22

23 Past studies (Scherer, 2000; Scherer, 2001) have proposed that the perception of affect
24 tends to be more similar, the more closely related the languages (see Section II). If this
25 were to pertain here, one should expect to find rather different groupings for these

1 languages. R, E and S are all Indo-European languages, even if they belong to different
2 branches and are geographically fairly distant. Therefore, one might expect these three
3 languages to form the main cluster, differing as a group from J. The grouping emerged
4 in this data suggests that cross-language differences are not easily explained in terms of
5 language relatedness.

6

7 Clearly, this study presents only a very limited contribution to the vast topic of cross-
8 language differences in voice-to-affect signalling. Future extension of this work could
9 include participants from more diverse languages. Of particular interest would be
10 languages where voice quality is exploited for segmental/lexical contrast (e.g., vowels
11 contrasting modal with tense, breathy, creaky voice etc.). Equally, tonal languages
12 would be interesting to explore, where lexical contrasts involve f_0 and, sometimes also
13 voice quality.

14

15 VI. CONCLUSIONS

16 Despite many points of convergence in affect attribution, there are considerable cross-
17 language differences in the affective profile associated with every one of the most
18 affect-carrying stimuli in this study (Figure 5). Results suggest that the two dimensions
19 of the voice, voice quality and f_0 , while they may act synergistically some of the time,
20 are not necessarily coupled in easy-to-predict ways. Thus, for example, the initial
21 expectation that *Tense Voice* would naturally combine with raised pitch to heighten the
22 high-activation affect was not borne out in a clear-cut way: languages differed in (i) the
23 extent to which this combination signalled high activation, and (ii) whether the
24 combined stimulus enhanced affective ratings.

25

1 All in all, these perception results serve as a reminder that voice quality and f_0 are, in
2 production terms, separately controllable. So, even though there are tendencies for
3 voice quality and f_0 to covary, the potential to exploit them individually may be an
4 important feature that allows for the richly nuanced expression of affect in speech. The
5 cross-language differences observed provide indicators that these vocal parameters are
6 being exploited in different ways depending on the language. This is highlighted by the
7 finding here that f_0 appears to play a rather different and more important role in affect
8 signalling in J and S than in R and E and that the way they combine to signal affect may
9 differ considerably.

10

11 The gaps in affective coverage that emerge in these results, where specific affects were
12 not signalled (or only weakly signalled) might be explained by certain factors. First, it
13 should be borne in mind that the selection of voice qualities used in this experiment was
14 not exhaustive, and that these were intentionally designed to be non-extreme exemplars
15 of particular voice qualities. (The f_0 contours included more extreme exemplars.) It is
16 possible that inclusion of extreme voice qualities might result in some of these gaps
17 being filled, and/or in higher affective ratings in certain cases. Likewise, the inclusion
18 of further voice qualities such as harsh voice and falsetto could be important for the
19 signalling of certain affects in one or other language. Furthermore, the combinations of
20 VQ with f_0 were necessarily limited in this study. The fact that the combined stimuli
21 here yielded only a single very clear-cut instance of synergy (*Whispery+FO fear*) may
22 partly be due to this limitation. The number of participants from each language group
23 was also admittedly small and a possibility remains that factors other than
24 language/culture may have caused differences among groups. For all these reasons, the

1 present experiment must be seen as an initial exploration, providing pointers to where
2 future work might be directed for a fuller teasing out of this complex question.

3

4 A major difficulty for both production and perception studies is that of ensuring that the
5 same affect is being targeted by different subjects. This difficulty becomes more
6 pronounced in a cross-language study: achieving conceptual/semantic equivalence of
7 the scale anchor terms in translation is crucial for the interpretation of results, and
8 emotion terms in one language do not always map perfectly onto terms in another
9 language (Russell, 1991; Mesquita *et al.*, 1997; Sabini and Silver, 2005; Ogarkova *et*
10 *al.*, 2009). Although considerable care was taken to mitigate this problem (Section
11 III.B), in any such cross-language study, the possibility must be borne in mind that the
12 affect labels in the different languages do not cover identical semantic fields.

13

14 The absence of a *formal* response in J to any of these stimuli was unexpected, as the use
15 of breathy/whispery voice has been mentioned in the past as a specific marker of
16 formality/politeness in this language (Campbell and Mokhtari, 2003; Ito, 2004).

17 However, the fact that J is also known to codify formality in a system of honorifics –
18 *keigo* ('terms of respect'), e.g., Ofuka *et al.* (2000), Ito (2005), may result in voice
19 quality being of relatively minor importance (or totally unimportant). In other words, if
20 formality is already encoded in the lexicon and grammar, there may be little need to
21 signal it through voice quality. Beyond such kinds of explanations, one must consider a
22 further possibility to do with the broader culture: the rules of affect expression may
23 differ in these languages and it may simply be less acceptable to express particular
24 affects in one culture than another – at least in certain social situations.

25

1 The vocal expression of affect is widely viewed as lying beyond the linguistic system
2 of the language. The present data suggest sometimes quite distinct language-related
3 patterns of voice-to-affect association. Other than in the case of extreme emotions
4 where vocal control gives way to involuntary effects, it seems that we can and do use
5 our voice to encode psychological-social information, integral to the intended message.
6 As argued elsewhere (Ní Chasaide and Gobl, 2004a; b), the temporal modulation of the
7 voice source to express affect is an essential part of the prosody of a language, and is
8 the hallmark not only of affective prosody, but also of linguistic prosody. One would
9 expect to find for affective prosody (as with linguistic prosody) that both quasi-
10 universal tendencies and language specific ‘rules’ (Gussenhoven, 2004, Chapter 5)
11 apply. While children acquire these ‘rules’ in their L1, this is an aspect that becomes
12 problematic in second language learning, being all the more problematic for being
13 poorly understood and difficult to describe in an explicit way. The misinterpretations
14 and misunderstandings that arise from the ‘incorrect’ vocal signalling of affect are all
15 the more important as the listener is not aware of the ‘error’. Whereas grammatical or
16 segmental errors are likely to be identified as part of a foreign accent, the use of
17 language-inappropriate voice prosody is not and tends to be simply interpreted in terms
18 of the L1 code of the listener, impacting on the quality of communication.

19

20 Similar considerations pertain to speech technology. It seems clear that affect
21 expression in speech synthesis needs to be language sensitive as one can readily
22 imagine the misapprehensions that would be occasioned by the use of Japanese
23 ‘intimate’ voice in an English speech synthesis system. There is an increasing demand
24 for applications involving ‘emotionally intelligent dialogue-partners’ capable of
25 providing affectively-appropriate speech output. In such applications and in future

1 technologies such as speech-to-speech translation systems language-sensitivity will be
2 an important consideration.

3

4 Ultimately, a proper understanding of the prosody of languages will depend on being
5 able to integrate their linguistic and affective dimensions within a single framework.

6 For this, being able to capture the different prosodic functions of voice source

7 modulation is key. A holistic understanding of prosody integrating the linguistic and

8 affective prosody will be needed for many practical applications, not only the

9 facilitation of the next generation of speech technologies, but also for language teaching

10 and for an understanding of how voice disorders impact on the ability to communicate

11 both linguistic and affective dimensions of the message.

12

13 As discussed in the Introduction, empirical research in this area is complex and

14 technically challenging. Accurate voice source production data on affect are very

15 difficult to obtain, especially for continuous speech where non-modal voice qualities

16 are used – which is the kind of data that is relevant here. The elusive and difficult

17 nature of the task makes empirical research more complex but does not diminish its

18 importance. In the present study, as in all approaches, there are in-built limitations in

19 scope, but results suggest that this approach nonetheless offers a useful instrument to

20 explore the field and promises new insights into this aspect of spoken communication.

21

22 **ACKNOWLEDGMENTS**

23 This work was partly funded by the FP6 Network of Excellence HUMAINE and was

24 further supported by An Roinn Cultúir, Oidhreachta agus Gaeltachta/Department of

- 1 Culture, Heritage and the Gaeltacht (*ABAIR* project). The authors would like to thank
- 2 the Editor and three anonymous reviewers for their helpful comments and suggestions.

¹ See supplementary material at [URL will be inserted by AIP] for the translations of the affective labels (Table I).

² See supplementary material at [URL will be inserted by AIP] for the data on interrater agreement (Table II).

³ See supplementary material at [URL will be inserted by AIP] for the results of the ANOVA (Table III).

1

2 **REFERENCES**

- 3 Airas, M., and Alku, P. (2004). "Emotions in short vowel segments: effects of the
4 glottal flow as reflected by the normalised amplitude quotient," in *Affective*
5 *Dialogue Systems* (Springer-Verlag Berlin Heidelberg), pp. 13-24.
- 6 Alku, P. (2011). "Glottal inverse filtering analysis of human voice production — A
7 review of estimation and parameterization methods of the glottal excitation and
8 their applications," *Sādhanā: Academy Proceedings in Engineering Sciences*
9 (Indian Academy of Sciences) **36**, 623–650.
- 10 Altrov, R. (2013). "Aspects of cultural communication in recognising emotions,"
11 *Trames* **17(67/62)**, 159-174.
- 12 Altrov, R., and Pajupuu, H. (2015). "The influence of language and culture on the
13 understanding of vocal emotions," *Journal of Estonian and Finno-Ugric*
14 *Linguistics* **6**, 11-48.
- 15 Bachorowski, J.-A., and Owren, M. J. (2003). "Sounds of emotions," *Ann. N. Y. Acad.*
16 *Sci.* **1000**, **Emotions Inside Out: 130 years after Darwin's The Expression of**
17 **Emotions in Man and Animals**, 244-265.
- 18 Banse, R., and Scherer, K. R. (1996). "Acoustic profiles in vocal emotion expression,"
19 *J. Pers. Soc. Psychol.* **70**, 614-636.
- 20 Bänziger, T., Hosoya, G., and Scherer, K. R. (2015). "Path models of vocal emotion
21 communication," *PLoS One* **10**, 1-29.
- 22 Bänziger, T., and Scherer, K. R. (2005). "The role of intonation in emotional
23 expression," *Speech Comm.* **46**, 252-267.
- 24 Bänziger, T., and Scherer, K. R. (2007). "Using actor portrayals to systematically study
25 multimodal emotion expression: the GEMEP corpus," in *Affective Computing*

- 1 *and Intelligent Interaction (ACII 2007)*, edited by A. Paiva, and R. W. Picard
2 (Springer-Verlag, Berlin), pp. 476-487.
- 3 Boula de Mareüil, P., Célérier, P., and Toen, J. (2002). "Generation of emotions by a
4 morphing technique in English, French and Spanish," in *Speech Prosody 2002*
5 (Aix-en-Provence, France).
- 6 Brislin, R. W. (1980). "Translation and content analysis of oral and written material," in
7 *Handbook of Cross-Cultural Psychology*, edited by H. C. Triandis, and J. W.
8 Berry (Allyn and Bacon, Inc, Boston), pp. 389-444.
- 9 Bryant, G. A., and Barrett, H. C. (2008). "Vocal emotion recognition across disparate
10 cultures," *Journal of Cognition and Culture* **8**, 135-148.
- 11 Burkhardt, F., Audibert, N., Malatesta, L., Türk, O., Arslan, L., and Aubergé, V.
12 (2006). "Emotional prosody - does culture make a difference?," in *Speech*
13 *Prosody 2006* (Dresden, Germany).
- 14 Burkhardt, F., and Sendlmeier, W. F. (2000). "Verification of acoustical correlates of
15 emotional speech using formant-synthesis," in *ITRW on Speech and Emotion*,
16 edited by R. Cowie, E. Douglas-Cowie, and M. Schröder (Newcastle, Northern
17 Ireland), pp. 151-156.
- 18 Burkhardt, F., and Stegmann, J. (2009). "Emotional speech synthesis: applications,
19 history and possible future," in *Elektronische Sprachsignalverarbeitung ESSV*
20 2009 (TUDpress, Dresden, Germany).
- 21 Campbell, N. (2002). "Recording techniques for capturing natural everyday speech," in
22 *Language Resources and Evaluation Conference ELREC'02* (Las Palmas,
23 Spain), pp. 2029-2032.

- 1 Campbell, N., and Mokhtari, P. (2003). "Voice quality: the 4th prosodic dimension," in
2 *XVth International Congress of Phonetic Sciences* (Barcelona, Spain), pp. 2417-
3 2420.
- 4 Carlson, R., Granström, B., and Nord, L. (1992). "Experiments with emotive speech -
5 acted utterances and synthesized replicas," in *2nd International Conference on*
6 *Spoken Language Processing (ICSLP 92)* (Banff, Alberta, Canada), pp. 671-
7 674.
- 8 Cheang, H. S., and Pell, M. D. (2008). "The sound of sarcasm," *Speech Comm.* **50**,
9 366-381.
- 10 Cowie, R. (2009). "Perceiving emotion: towards a realistic understanding of the task,"
11 *Philosophical Transactions of the Royal Society B* **364**, 3515-3525.
- 12 Cowie, R., and Cornelius, R. R. (2003). "Describing the emotional states that are
13 expressed in speech," *Speech Comm.* **40**, 5-32.
- 14 Cummings, K. E., and Clements, M. A. (1995). "Analysis of the glottal excitation of
15 emotionally styled and stressed speech," *J. Acoust. Soc. Am.* **98**, 88-98.
- 16 Cummins, N., Scherer, S., Krajewski, J., Schnieder, S., Epps, J., and Quatieri, T. F.
17 (2015). "A review of depression and suicide risk assessment using speech
18 analysis," *Speech Comm.* **71**, 10-49.
- 19 Douglas-Cowie, E., Campbell, N., Cowie, R., and Roach, P. (2003). "Emotional
20 speech: towards a new generations of databases," *Speech Comm.* **40**, 33-60.
- 21 Drioli, C., Tisato, G., Cosi, P., and Tesser, F. (2003). "Emotions and voice quality:
22 experiments with sinusoidal modelling," in *VOQUAL'03* (Switzerland), pp. 127-
23 132.
- 24 Ekman, P. (1993). "Facial expression and emotion," *Am. Psychol.* **48**, 384-392.

- 1 Ekman, P., Friesen, W. V., O'Sullivan, M., Chan, A., Diacoyanni-Tarlatzis, I., Heider,
2 K., Krause, R., LeCompte, W. A., Pitcairn, T., Ricci-Bitti, P. E., Scherer, K. R.,
3 Tomita, M., and Tzavaras, A. (1987). "Universals and cultural differences in the
4 judgments of facial expressions of emotion," *J. Pers. Soc. Psychol.* **53**, 712-717.
- 5 Elfenbein, H. A., and Ambady, N. (2002a). "Is there an in-group advantage in emotion
6 recognition?," *Psychol. Bull.* **128**, 243-249.
- 7 Elfenbein, H. A., and Ambady, N. (2002b). "On the universality and cultural specificity
8 of emotion recognition: a meta-analysis," *Psychol. Bull.* **128**, 203-235.
- 9 Eyben, F., Scherer, K. R., Schuller, B. W., Sundberg, J., André, E., Busso, C.,
10 Devillers, L. Y., Epps, J., Laukka, P., Narayanan, S. S., and Truong, K. P.
11 (2016). "The Geneva Minimalistic Acoustic Parameter Set (GeMAPS) for voice
12 research and affective computing," *IEEE Transactions on Affective Computing*
13 **7**, 190-202.
- 14 Fant, G. (1995). "The LF-model revisited: transformations and frequency domain
15 analysis," *STL-QPSR* **2-3**, 119-156.
- 16 Fant, G. (1997). "The voice source in connected speech," *Speech Comm.* **22**, 125-139.
- 17 Fant, G., Liljencrants, J., and Lin, Q. (1985). "A four-parameter model of glottal flow,"
18 *STL-QPSR* **4**, 1-13.
- 19 Gobl, C. (1988). "Voice source dynamics in connected speech," *STL-QPSR* **1**, 123-159.
- 20 Gobl, C. (1989). "A preliminary study of acoustic voice quality correlates," *STL-QPSR*
21 **30**, 9-22.
- 22 Gobl, C., Bennett, E., and Ní Chasaide, A. (2002). "Expressive synthesis: how crucial is
23 voice quality?," in *IEEE Workshop on Speech Synthesis* (Santa Monica,
24 California, USA), pp. 1-4.

- 1 Gobl, C., and Ní Chasaide, A. (1992). "Acoustic characteristics of voice quality,"
2 Speech Comm. **11**, 481-490.
- 3 Gobl, C., and Ní Chasaide, A. (1999). "Techniques for analysing the voice source," in
4 *Coarticulation: Theory, Data and Techniques*, edited by W. J. Hardcastle, and
5 N. Hewlett (Cambridge University Press, Cambridge), pp. 300-321.
- 6 Gobl, C., and Ní Chasaide, A. (2000). "Testing affective correlates of voice quality
7 through analysis and resynthesis," in *ITRW on Speech and Emotion* (Newcastle,
8 Northern Ireland), pp. 178-183.
- 9 Gobl, C., and Ní Chasaide, A. (2003a). "Amplitude-based source parameters for
10 measuring voice quality," in *VOQUAL'03* (Geneva, Switzerland), pp. 151-156.
- 11 Gobl, C., and Ní Chasaide, A. (2003b). "The role of voice quality in communicating
12 emotion, mood and attitude," *Speech Comm.* **40**, 189-212.
- 13 Gobl, C., and Ní Chasaide, A. (2010). "Voice source variation and its communicative
14 functions," in *The Handbook of Phonetic Sciences*, edited by W. J. Hardcastle,
15 J. Laver, and F. E. Gibbon (Blackwell Publishing Ltd, Oxford), pp. 378-423.
- 16 Goudbeek, M., and Scherer, K. (2010). "Beyond arousal: Valence and potency/control
17 cues in the vocal expression of emotion," *J. Acoust. Soc. Am.* **128**, 1322-1336.
- 18 Graham, C. R., Hamblin, A. W., and Feldstein, S. (2001). "Recognition of emotion in
19 English by speakers of Japanese, Spanish and English," *IRAL - International*
20 *Review of Applied Linguistics in Language Teaching* **39**, 19-37.
- 21 Grandjean, D., Bänziger, T., and Scherer, K. R. (2006). "Intonation as an interface
22 between language and affect," *Prog. Brain Res.* **156**, 235-268.
- 23 Gussenhoven, K. (2004). *The Phonology of Tone and Intonation* (Cambridge
24 University Press, Cambridge), p.355.

- 1 Guzman, M., Correa, S., Muñoz, D., and Mayerhoff, R. (2013). "Influence on spectral
2 energy distribution of emotionsl expression," *J. Voice* **27**, 129.e121-129.e110.
- 3 Hammarberg, B., Fritzell, B., Gaufin, J., Sundberg, J., and Wedin, L. (1980).
4 "Perceptual and acoustic correlates of abnormal voice qualities," *Acta*
5 *Otolaryngol.* **90**, 441-451.
- 6 Hanson, H. M. (1997). "Glottal characteristics of female speakers: acoustic correlates,"
7 *J. Acoust. Soc. Am.* **101**, 466-481.
- 8 Hanson, H. M., and Chuang, E. S. (1999). "Glottal characteristics of male speakers:
9 acoustic correlates and comparison with female data," *J. Acoust. Soc. Am.* **106**,
10 1064-1077.
- 11 Heldner, M. (2003). "On the reliability of overall intensity and spectral emphasis as
12 acoustic correlates of focal accents in Swedish," *Journal of Phonetics* **31**, 39-62.
- 13 Iseli, M., Shue, Y.-L., Epstein, M. A., Keating, P., Kreiman, J., and Alwan, A. (2006).
14 "Voice source correlates of prosodic features in American English," in
15 *Interspeech 2006 - ICSLP* (Pittsburgh, Pennsylvania, USA), paper 1933-
16 Thu1931A1933O.1931.
- 17 Ishi, C. T., Ishiguro, H., and Hagita, N. (2008). "The role of breathy/whispery voice
18 qualities in dialogue speech," in *Speech Prosody 2008* (Campinas, Brazil).
- 19 Ito, M. (2004). "Politeness and voice quality: The alternative method to measure
20 aspiration noise," in *Speech Prosody 2004* (Nara, Japan), pp. 213-216.
- 21 Ito, M. (2005). "The contribution of voice quality to the expression of politeness: an
22 experimental study" (unpublished doctoral dissertation), (University of
23 Edinburgh, Edinburgh).

- 1 Johnstone, T., and Scherer, K. R. (2000). "Vocal communication of emotion," in
2 *Handbook of Emotions*, edited by M. Lewis, and J. M. Haviland-Jones (Guilford
3 Press, New York), pp. 220-235.
- 4 Juslin, P., Laukka, P., and Bänziger, T. (2017). "The mirror of our soul? Comparisons
5 of spontaneous and posed vocal expression of emotion," *Journal of Nonverbal*
6 *Behavior*. Advance online publication.
- 7 Juslin, P. N., and Laukka, P. (2001). "Impact of intended emotion intensity on cue
8 utilization and decoding accuracy in vocal expression of emotion," *Emotion* **1**,
9 381-412.
- 10 Juslin, P. N., and Laukka, P. (2003). "Communication of emotions in vocal expression
11 and music performance: different channels, same code?," *Psychol. Bull.* **5**, 770-
12 814.
- 13 Juslin, P. N., and Scherer, K. R. (2005). "Vocal expression of affect," in *The New*
14 *Handbook of Methods in Nonverbal Behavior Research*, edited by J. Harrigan,
15 R. Rosenthal, and K. R. Scherer (Oxford University Press, Oxford), pp. 65-135.
- 16 Kappas, A., Hess, U., and Scherer, K. R. (1991). "Voice and emotion," in
17 *Fundamentals of Nonverbal Behavior*, edited by R. S. Feldman, and B. Rimé
18 (Cambridge University Press, Cambridge), pp. 200-238.
- 19 Keller, E. (2005). "The analysis of voice quality in speech processing," in *Nonlinear*
20 *Speech Modeling and Applications* (Springer Berlin / Heidelberg), pp. 54-73.
- 21 Kitayama, S., and Ishii, K. (2002). "Word and voice: spontaneous attention to
22 emotional utterances in two languages," *Cognition and Emotion* **16**, 29-59.
- 23 Klasmeyer, G., and Sendlmeier, W. F. (2000). "Voice and emotional states," in *Voice*
24 *Quality Measurement*, edited by R. D. Kent, and M. J. Ball (Singular Publishing
25 Group, San Diego), pp. 339-357.

- 1 Klatt, D. H., and Klatt, L. C. (1990). "Analysis, synthesis, and perception of voice
2 quality variations among female and male talkers," *J. Acoust. Soc. Am.* **87**, 820-
3 857.
- 4 Koeda, M., Belin, P., Hama, T., Masuda, T., Matsuura, M., and Okubo, Y. (2013).
5 "Cross-cultural differences in the processing of non-verbal affective
6 vocalizations by Japanese and Canadian listeners," *Front. Psychol.* **4**, 105.
- 7 Kreiman, J., Gerratt, B. R., Garellek, M., Samlan, R., and Zhang, Z. (2014). "Toward a
8 unified theory of voice production and perception," *Loquens* **1**, e009.
- 9 Kreiman, J., Gerratt, B. R., and Ito, M. (2007). "When and why listeners disagree in
10 voice quality assessment tasks," *J. Acoust. Soc. Am.* **122**, 2354-2364.
- 11 Kreiman, J., and Sitdis, D. (2011). "Perception of emotion and personality from voice,"
12 in *Foundations of Voice Studies: an Interdisciplinary Approach to Voice*
13 *Production and Perception* (Wiley-Blackwell, United Kingdom), pp. 302-360.
- 14 Ladd, D. R., Scherer, K. R., and Silverman, K. (1986). "An integrated approach to
15 studying intonation and attitude," in *Intonation in Discourse*, edited by K.
16 Johns-Lewis (Croom Helm, London), pp. 125-138.
- 17 Ladd, D. R., Silverman, K. E. A., Tolkmitt, F., Bergmann, G., and Scherer, K. R.
18 (1985). "Evidence for the independent function of intonation contour type, voice
19 quality, and F_0 range in signaling speaker affect " *J. Acoust. Soc. Am.* **78**, 435-
20 444.
- 21 Landis, J. R., and Koch, G. G. (1977). "The measurement of observer agreement for
22 categorical data," *Biometrics* **33**, 159-174.
- 23 Laukka, P. (2008). "Research on vocal expression of emotion: state of the art and future
24 directions," in *Emotions in the Human Voice*, edited by K. Izdebski (Plural
25 Publishing, San Diego, CA), pp. 153-169.

- 1 Laukka, P., Elfenbein, H. A., Thingujam, N. S., Rockstuhl, T., Iraki, F. K., Chui, W.,
2 and Althoff, J. (2016). "The expression and recognition of emotions in the voice
3 across five nations: A lens model analysis based on acoustic features," *J. Pers.*
4 *Soc. Psychol.* **111**, 686-705.
- 5 Laukka, P., Juslin, P., and Bresin, R. (2005). "A dimensional approach to vocal
6 expression of emotion," *Cognition and Emotion* **19**, 633-653.
- 7 Laukka, P., Neiberg, D., Forsell, M., Karlsson, I., and Elenius, K. (2011). "Expression
8 of affect in spontaneous speech: Acoustic correlates and automatic detection of
9 irritation and resignation," *Computer Speech & Language* **25**, 84-104.
- 10 Laukkanen, A.-M., Alku, P., Airas, M., and Waaramaa, T. (2008). "The role of voice
11 quality in the expression and perception of emotion," in *Emotions in the Human*
12 *Voice*, edited by K. Izdebski (Plural Publishing, San Diego, CA), pp. 171-184.
- 13 Laukkanen, A.-M., Vilkman, E., Alku, P., and Oksanen, H. (1996). "Physical variations
14 related to stress and emotional state: a preliminary study," *Journal of Phonetics*
15 **24**, 313-335.
- 16 Laver, J. (1980). *The Phonetic Description of Voice Quality* (Cambridge University
17 Press, Cambridge), p.186.
- 18 McCluskey, K. W., and Albas, D. C. (1981). "Perception of the emotional content of
19 speech by Canadian and Mexican children, adolescents, and adults," *Int. J.*
20 *Psychol.* **16**, 119-132.
- 21 McCluskey, K. W., Albas, D. C., Niemi, R. R., Cuevas, C., and Ferrer, C. A. (1975).
22 "Cross-cultural differences in the perception of the emotional content of speech:
23 a study of the development of sensitivity in Canadian and Mexican children,"
24 *Dev. Psychol.* **11**, 551-555.

- 1 McGraw, K. O., and Wong, S. P. (1996). "Forming inferences about some intraclass
2 correlation coefficients," *Psychol. Methods* **1**, 30-46.
- 3 Mesquita, B., Frijda, N. H., and Scherer, K. R. (1997). "Culture and emotion," in
4 *Handbook of Cross-Cultural Psychology: Basic Processes and Human*
5 *Development*, edited by J. W. Berry, P. R. Dasen, and T. S. Saraswathi
6 (Longwood Professional Books).
- 7 Mesquita, B., and Walker, R. (2003). "Cultural differences in emotions: a context for
8 interpreting emotional experiences," *Behav. Res. Ther.* **41**, 777-793.
- 9 Mozziconacci, S. (1995). "Pitch variations and emotions in speech," in *XIIIth*
10 *International Congress of Phonetic Sciences* (Stockholm), pp. 178-181.
- 11 Mozziconacci, S. (1998). *Speech Variability and Emotion: Production and Perception.*
12 *Doctoral Thesis* (Technische Universiteit Eindhoven, Eindhoven), p. 210.
- 13 Mozziconacci, S. (2002). "Prosody and emotions," in *Speech Prosody 2002* (Aix-en-
14 Provence, France).
- 15 Mozziconacci, S. J. L., and Hermes, D. (1999). "Role of intonation patterns in
16 conveying emotion in speech," in *The XIVth International Congress of Phonetic*
17 *Sciences* (San Francisco, USA), pp. 2001-2004.
- 18 Murphy, P. J., and Laukkanen, A.-M. (2009). "Electroglottogram analysis of
19 emotionally styled phonation," in *Multimodal Signals: Cognitive and*
20 *Algorithmic Issues* (Springer Berlin/Heidelberg), pp. 264-270.
- 21 Murray, I. R., and Arnott, J. L. (1993). "Toward the simulation of emotion in synthetic
22 speech: a review of the literature on human vocal emotion," *J. Acoust. Soc. Am.*
23 **93**, 1907-1108.
- 24 Ní Chasaide, A., and Gobl, C. (1993). "Contextual variation of the vowel voice source
25 as a function of adjacent consonants," *Lang. Speech* **36**, 303-330.

- 1 Ní Chasaide, A., and Gobl, C. (1995). "Towards acoustic profiles of phonatory
2 qualities," in *XIIIth International Congress of Phonetic Sciences* (Stockholm),
3 pp. 6-13.
- 4 Ní Chasaide, A., and Gobl, C. (2004a). "Decomposing linguistic and affective
5 components of phonatory quality," in *Interspeech 2004* (Jeju Island, Korea), pp.
6 901-904.
- 7 Ní Chasaide, A., and Gobl, C. (2004b). "Voice quality and f_0 in prosody: towards a
8 holistic account," in *Speech Prosody 2004* (Nara, Japan), pp. 189-196.
- 9 Ní Chasaide, A., Yanushevskaya, I., Kane, J., and Gobl, C. (2013). "The Voice
10 Prominence Hypothesis: the interplay of F0 and voice source features in
11 accentuation," in *Interspeech 2013* (Lyon, France), pp. 3527-3531.
- 12 Oatley, K., Keltner, D., and Jenkins, J. M. (2006). *Understanding Emotions* (Blackwell
13 Publishing Ltd, Oxford), p.536.
- 14 Ofuka, E., McKeown, J. D., Waterman, M. G., and Roach, P. J. (2000). "Prosodic cues
15 for rated politeness in Japanese speech," *Speech Comm.* **32**, 199-217.
- 16 Ogarkova, A., Borgeaud, P., and Scherer, K. R. (2009). "Language and culture in
17 emotion research: a multidisciplinary perspective," *Social Science Information*
18 **48**, 339-357.
- 19 Paeschke, A. (2004). "Global trend of fundamental frequency in emotional speech," in
20 *Speech Prosody 2004* (Nara, Japan).
- 21 Paeschke, A., Kienast, M., and Sendlmeier, W. F. (1999). "F0-contours in emotional
22 speech," in *XIVth International Congress of Phonetic Sciences* (San Francisco,
23 USA).
- 24 Pakosz, M. (1983). "Attitudinal judgments in intonation: some evidence for a theory,"
25 *J. Psycholinguist. Res.* **12**, 311-326.

- 1 Patel, S., Scherer, K. R., Björkner, E., and Sundberg, J. (2011). "Mapping emotions into
2 acoustic space: The role of voice production," *Biol. Psychol.* **87**, 93-98.
- 3 Pell, M. D., Monetta, L., Paulmann, S., and Kotz, S. A. (2009a). "Recognizing
4 emotions in a foreign language," *Journal of Nonverbal Behaviour* **33**, 107-120.
- 5 Pell, M. D., Paulmann, S., Dara, C., Allasseri, A., and Kotz, S. A. (2009b). "Factors in
6 the recognition of vocally expressed emotions: a comparison of four languages,"
7 *Journal of Phonetics* **37**, 417-435.
- 8 Pell, M. D., and Scorup, V. (2008). "Implicit processing of emotional prosody in a
9 foreign versus native language," *Speech Comm.* **50**, 519-530.
- 10 Pittam, J., Gallois, C., and Callan, V. (1990). "The long-term spectrum and perceived
11 emotion," *Speech Comm.* **9**, 177-187.
- 12 Russell, J. A. (1991). "Culture and the categorisation of emotions," *Psychol. Bull.* **110**,
13 426-450.
- 14 Russell, J. A., Bachorowski, J.-A., and Fernández-Dols, J.-M. (2003). "Facial and vocal
15 expression of emotion," *Annu. Rev. Psychol.* **54**, 329-349.
- 16 Ryan, C., Ní Chasaide, A., and Gobl, C. (2003). "Voice quality variation and the
17 perception of affect: continuous or categorical?," in *XVth International*
18 *Congress of Phonetic Sciences* (Barcelona, Spain), pp. 2409-2412.
- 19 Sabini, J., and Silver, M. (2005). "Why emotion names and experiences do not neatly
20 pair," *Psychol. Inq.* **16**, 1-10.
- 21 Sadanobu, T. (2004). "A natural history of Japanese pressed voice," *Journal of the*
22 *Phonetic Society of Japan* **8**, 29-44.
- 23 Sauter, D. A., Eisner, F., Calder, A. J., and Scott, S. K. (2010a). "Perceptual cues in
24 nonverbal vocal expressions of emotion," *Quarterly Journal of Experimental*
25 *Psychology* **63**, 2251-2272.

- 1 Sauter, D. A., Eisner, F., Ekman, P., and Scott, S. K. (2010b). "Cross-cultural
2 recognition of basic emotions through nonverbal emotional vocalizations," Proc.
3 Natl. Acad. Sci. U. S. A. **107**, 2408-2412.
- 4 Scherer, K. R. (1986). "Vocal affect expression: a review and a model for future
5 research," Psychol. Bull. **99**, 143-165.
- 6 Scherer, K. R. (2000). "A cross-cultural investigation of emotion inferences from voice
7 and speech: implications for speech technology," in *6th International
8 Conference on Spoken Language Processing* (Beijing, China), pp. 379-382.
- 9 Scherer, K. R. (2003). "Vocal communication of emotion: a review of research
10 paradigms," Speech Comm. **40**, 227-256.
- 11 Scherer, K. R. (2013). "Vocal markers of emotion: Comparing induction and acting
12 elicitation," Computer Speech & Language **27**, 40-58.
- 13 Scherer, K. R., Banse, R., and Wallbott, H. G. (2001). "Emotion inferences from vocal
14 expression correlate across languages and cultures," J. Cross Cult. Psychol. **32**,
15 76-92.
- 16 Scherer, K. R., Clark-Polner, E., and Mortillaro, M. (2011). "In the eye of the beholder?
17 Universality and cultural specificity in the expression and perception of
18 emotion," Int. J. Psychol. **46**, 401-435.
- 19 Scherer, K. R., Ladd, D. R., and Silverman, K. E. A. (1984). "Vocal cues to speaker
20 affect: testing two models," J. Acoust. Soc. Am. **76**, 1346-1356.
- 21 Schröder, M. (2001). "Emotional speech synthesis," in *Eurospeech 2001* (Aalborg,
22 Denmark), pp. 561-564.
- 23 Shrouf, P. E., and Fleiss, J. L. (1979). "Intraclass correlations: uses in assessing rater
24 reliability," Psychol. Bull. **86**, 420-428.

- 1 Streiner, D. L., and Norman, G. R. (2008). *Health Measurement Scales* (Oxford
2 University Press, Oxford), p. 431.
- 3 Sundberg, J., and Nordenberg, M. (2006). "Effects of vocal loudness variation on
4 spectrum balance as reflected by the alpha measure of long-term-average
5 spectra of speech," *J. Acoust. Soc. Am.* **120**, 453-457.
- 6 Sundberg, J., Patel, S., Björkner, E., and Scherer, K. R. (2011). "Interdependencies
7 among voice source parameters in emotional speech," *IEEE Transactions on*
8 *Affective Computing* **2**, 162-174.
- 9 van Bezooijen, R., Otto, S. A., and Heenan, T. A. (1983). "Recognition of vocal
10 expressions of emotion: a three-nation study to identify universal
11 characteristics," *J. Cross Cult. Psychol.* **14**, 387-406.
- 12 Waaramaa, T. (2014). "Perception of emotional nonsense sentences in China, Egypt,
13 Estonia, Finland, Russia, Sweden, and the USA," *Logopedics Phoniatics*
14 *Vocology*, 1-7.
- 15 Waaramaa, T., Laukkanen, A. M., Alku, P., and Väyrynen, E. (2008). "Monopitched
16 expression of emotions in different vowels," *Folia Phoniatica et Logopaedica*
17 **60**, 249-255.
- 18 Waaramaa, T., and Leisiö, T. (2013). "Perception of emotionally loaded vocal
19 expressions and its connection to responses to music. A cross-cultural
20 investigation: Estonia, Finland, Sweden, Russia, and the USA," *Front. Psychol.*
21 **4**, 344.
- 22 Westermann, R., Spies, K., Stahl, G., and Hesse, F. W. (1996). "Relative effectiveness
23 and validity of mood induction procedures: a meta-analysis," *Eur. J. Soc.*
24 *Psychol.* **26**, 557-580.

- 1 Williams, C. E., and Stevens, K. N. (1972). "Emotions and speech: some acoustical
2 correlates," *J. Acoust. Soc. Am.* **52**, 1238-1250.
- 3 Yanushevskaya, I., Gobl, C., and Ní Chasaide, A. (2009). "Voice parameter dynamics
4 in portrayed emotions," in *6th International Workshop on Models and Analysis
5 of Vocal Emissions for Biometrical Applications (MAVEBA 2009)* (Florence,
6 Italy), pp. 21-24.
- 7 Yanushevskaya, I., Gobl, C., and Ní Chasaide, A. (2013). "Voice quality in affect
8 cueing: does loudness matter?," *Front. Psychol.* **4:335**, 1-14.
- 9 Yanushevskaya, I., Tooher, M., Gobl, C., and Ní Chasaide, A. (2007). "Time- and
10 amplitude-based voice source correlates of emotional portrayals," in *Affective
11 Computing and Intelligent Interaction: Proceedings of the ACII 2007*, edited by
12 A. Paiva, R. Prada, and R. W. Picard (Springer-Verlag, Lisbon, Portugal), pp.
13 159-170.
- 14 Yuan, J., Shen, L., and Chen, F. (2002). "The acoustic realisation of anger, fear, joy and
15 sadness in Chinese," in *7th International Conference on Spoken Language
16 Processing* (Denver, Colorado, USA), pp. 2025-2028.
- 17 Zinken, J., Knoll, M., and Panksepp, J. (2008). "Universality and diversity in the
18 vocalisation of emotions," in *Emotions in the Human Voice*, edited by K.
19 Izdebski (Plural Publishing, San Diego, CA), pp. 185-202.
- 20
21
22
23
24
25

- 1 Table I. The synthesised stimuli used in the cross-language study. All F0 stimuli have
 2 modal voice. All VQ stimuli have neutral F0. Additionally, the stimulus *Modal+F0*
 3 *neutral* is included for baseline comparison. For the stimuli with * see text.

Stimulus Group	Stimulus Type (type of manipulation)		
	VQ	F0	VQ+F0
WHISPERY FEAR	<i>Whispery</i>	<i>F0 fear</i>	<i>Whispery+F0 fear</i>
BREATHY SADNESS	<i>Breathy</i>	<i>F0 sadness</i>	<i>Breathy+F0 sadness</i>
LAX-CREAKY BOREDOM	<i>Lax-creaky*</i>	<i>F0 boredom</i>	<i>Lax-creaky+F0 lowered*</i>
TENSE JOY	<i>Tense</i>	<i>F0 joy</i>	<i>Tense+F0 joy</i>
TENSE INDIGNATION	<i>Tense</i>	<i>F0 indignation</i>	<i>Tense+F0 indignation</i>
Baseline Stimulus: <i>Modal+F0 neutral</i>			

1

2 Table II. Stimuli that obtained the highest ratings (absolute values shown as subscripts) within each stimulus type for the affects tested. Stimuli
 3 yielding the highest ratings across stimulus types for at least one language are shown in bold type.

Affect	VQ	F0 (all with Modal Voice)	VQ+F0
<i>indignant</i>	Tense <i>R</i> 2.18, <i>E</i> 1.78, <i>S</i> 1.64, <i>J</i> 1.53	F0 indignation <i>R</i> 1.68, <i>E</i> 1.50, <i>S</i> 1.02	Tense+F0 indignation <i>R</i> 1.67, <i>E</i> 1.27
<i>interested</i>	Tense <i>R</i> 1.72, <i>S</i> 1.05	F0 indignation <i>J</i> 2.08, <i>R</i> 1.78, <i>E</i> 1.66, <i>S</i> 1.42	Tense+F0 indignation <i>J</i> 2.15, <i>E</i> 1.62, <i>R</i> 1.53, <i>S</i> 1.51
<i>formal</i> ^a	Tense <i>R</i> 2.38, <i>E</i> 1.99 Modal <i>S</i> 1.05	F0 boredom <i>R</i> 1.58, <i>S</i> 1.25, <i>E</i> 1.14	Tense+F0 joy <i>R</i> 1.43
<i>stressed</i> ^b	Tense <i>R</i> 1.79, <i>E</i> 1.70, <i>S</i> 1.09	F0 indignation <i>S</i> 1.79, <i>R</i> 1.68, <i>E</i> 1.67	Tense+F0 indignation <i>S</i> 1.87, <i>R</i> 1.57, <i>E</i> 1.44
<i>happy</i>	Tense <i>R</i> 1.31	F0 indignation <i>R</i> 1.21	
<i>fearless</i>	Tense <i>R</i> 2.22, <i>E</i> 1.72, <i>S</i> 1.61, <i>J</i> 1.30	F0 boredom <i>R</i> 1.41, <i>E</i> 1.04	Tense+F0 joy <i>R</i> 1.00
<i>apologetic</i>	Lax-creaky <i>R</i> 1.63, <i>S</i> 1.37 Whispery <i>E</i> 1.33	F0 fear <i>J</i> 1.37	Lax-creaky+F0 lowered <i>R</i> 1.57, <i>S</i> 1.23 Whispery+F0 fear <i>J</i> 1.82, <i>E</i> 1.60
<i>bored</i>	Lax-creaky <i>J</i> 1.96, <i>R</i> 1.92, <i>S</i> 1.87, <i>E</i> 1.46		Lax-creaky+F0 lowered <i>J</i> 2.35, <i>S</i> 1.87, <i>R</i> 1.77, <i>E</i> 1.35
<i>intimate</i>	Lax-creaky <i>R</i> 1.81 Whispery <i>E</i> 1.29	F0 fear <i>S</i> 1.19 F0 indignation <i>J</i> 1.55	Lax-creaky+F0 lowered <i>R</i> 1.87 Whispery+F0 fear <i>E</i> 1.62 Tense+F0 indignation <i>J</i> 1.66, <i>S</i> 1.23
<i>relaxed</i> ^c	Lax-creaky <i>S</i> 2.09, <i>R</i> 1.53, <i>E</i> 1.51		Lax-creaky+F0 lowered <i>S</i> 2.13, <i>R</i> 1.64, <i>E</i> 1.61
<i>sad</i>	Lax-creaky <i>R</i> 1.82, <i>S</i> 1.74, <i>E</i> 1.49, <i>J</i> 1.49		Lax-creaky+F0 lowered <i>S</i> 1.69, <i>R</i> 1.65, <i>J</i> 1.38, <i>E</i> 1.34
<i>scared</i>		F0 fear <i>J</i> 1.27, <i>R</i> 1.25, <i>S</i> 1.16, <i>E</i> 1.09	Whispery+F0 fear <i>E</i> 1.91, <i>R</i> 1.83, <i>J</i> 1.67, <i>S</i> 1.37

^a None of the stimuli signalled *formal* to Japanese.

^b None of the stimuli signalled *stressed* to Japanese.

^c None of the stimuli signalled *relaxed* to Japanese.

1

2 Table III. Summary of the stimulus to affect association for the languages tested (**R**=Russian; **E**=Irish-English; **S**=Spanish; **J**=Japanese). Only

3 the stimuli yielding the highest rating are shown.

	Indign.	Interest.	Happy	Fearless	Formal	Stressed	Apologet	Bored	Sad	Scared	Intimate	Relaxed
Tense	RESJ		R	RESJ	RE	RE						
F0 indignation		RE										
Tense+F0 indignation		SJ				S					SJ	
Lax-creaky/ Lax-creaky+F0 lowered							R S	RESJ	RESJ		R	RES
Modal/F0 boredom					S							
Whispery+F0 fear							E J			RESJ	E	
Gaps			ESJ		J	J						J

4

1 **The list of figure captions**

2

3 Figure 1 (color online). Parameter variation in the synthesised voice quality stimuli
4 (after Gobl & Ní Chasaide, 2003).

5

6 Figure 2 (color online). f_0 contours used in the F0 stimuli. The lowest contour shows the
7 f_0 contour used with the lax-creaky voice quality (see text).

8

9 Figure 3 (color online). Affective ratings obtained in the six subtests: each row
10 represents a listening test (e.g., *indignant-apolgetic*), each column represents a
11 particular stimulus group (e.g., WHISPERY||FEAR). Languages: Irish-English (E); Russian
12 (R), Spanish (S), Japanese (J). Statistically significant cross-language differences: * $p <$
13 0.05 , ** $p < 0.01$, *** $p < 0.001$. The most highly rated stimulus for a particular affect
14 is marked by a surrounding box.

15

16 Figure 4 (color online). Ratings for the three stimulus types, from the three most
17 effective stimulus groups in each language.

18

19 Figure 5 (color online). Summary plots of voice-to-affect associations for the most
20 highly rated stimuli. The data for *Lax-creaky* (solid line) and *Lax-creaky+F0 lowered*
21 (dotted line) stimuli are superimposed. Larger data points show cases when a stimulus
22 yielded highest ratings compared to other stimuli for a particular affect.

23