

Running Head: SMILES AND FROWNS IN SPEECH

Audible smiles and frowns affect speech comprehension

Hugo Quené

Utrecht institute of Linguistics OTS, Utrecht University,

Trans 10, 3512 JK Utrecht, The Netherlands

h.quene@uu.nl

Gün R. Semin and Francesco Foroni

Faculty of Social and Behavioral Sciences, Utrecht University,

Heidelberglaan 1, 3584 CS Utrecht, The Netherlands

g.r.semin@uu.nl

f.foroni@uu.nl

article to appear in *Speech Communication*

version R2

14 March 2012

Corresponding author:

Dr Hugo Quené

Utrecht inst of Linguistics OTS, Utrecht University

Trans 10, 3512 JK Utrecht, The Netherlands

Tel: +31 30 2536070, Fax: +31 30 2536000

Email: h.quene@uu.nl

Abstract

Motor resonance processes are involved both in language comprehension and in affect perception. Therefore we predict that listeners understand spoken affective words slower, if the phonetic form of a word is incongruent with its affective meaning. A language comprehension study involving an interference paradigm confirmed this prediction. This interference suggests that affective phonetic cues contribute to language comprehension. A perceived smile or frown affects the listener, and hearing an incongruent smile or frown impedes our comprehension of spoken words.

Keywords:

Smiles; Speech comprehension; Emotion; Affect perception; Motor resonance;

Audible smiles and frowns affect speech comprehension

1. Introduction

In spoken language, vowels and consonants convey the linguistic meaning intended by the speaker. In addition, and unlike written language, speech also conveys a speaker's emotional state [1,2,3,4,5] mainly by means of its prosody. In addition, the audible properties of the speaker's vocal tract, in particular its second spectral resonance (formant F2) as well as the dispersion between formants, convey whether or not a speaker is smiling while talking [6,7]. Emotionally and affectively nuanced utterances play a central role in speech communication, by conveying importance, relevance, urgency, and attitude, in addition to the spoken semantic content. Listeners can decode audible affective cues such as smiles and frowns [8], even with unfamiliar speakers [9] and in foreign languages [10].

In this study we hypothesize that comprehension of a word's semantic meaning and affect perception based on its phonetic form, are not separate, but interacting components of spoken word processing. The presumed causal mechanism for this interaction is motor resonance [11,12] which is involved in listeners' retrieval of linguistic meaning [13,14], as well as in perception of affect [15,16,17,18]. Thus we investigate whether affectively meaningful phonetic features, related to affective facial expressions such as smiles and frowns, also influence spoken word recognition. We predict that spoken word perception will be faster if the semantic meaning and the affective phonetic cues of a word are congruent, relative to spoken words with incongruity between semantic content and

affective phonetic form. This incongruence would yield a phonetic Stroop effect [19]: if the positive word *pleasant* is spoken with an incongruent affective phonetic form (i.e., frown), its semantic evaluation is predicted to be slower than if this positive word *pleasant* is spoken with a congruent smile.

Previous studies have already shown that emotionally and socially incongruent phonetic forms do indeed have a negative effect on speech processing. For example, semantic evaluation of happy, neutral and angry words was found to be slower if these emotional words were spoken with incongruent emotional prosody [20; cf. 21, 22, 23]. Similarly, naming latencies for happy, neutral and sad words were longer if the emotional words were spoken with incongruent emotional prosody [24]. In an eye-tracking study with visually presented faces expressing various emotions, listeners gazed more frequently to faces with emotions congruent to the prosody of the speech stimuli [25]. Spoken sentence comprehension was also affected (as indicated by an N400 effect) by inconsistency or incongruence between the semantic content and the speaker characteristics (gender, age, and social status) expressed by the speaker's voice [26, 27].

The present study aims to expand this converging evidence, in three ways. First, our focus is on smiles and frowns as affective facial gestures, and not on phonetic expressions of basic emotions [4]. Because smile gestures and frown gestures necessarily interfere with speech production, the perception of such affective speech may show stronger motor resonance [11,12,13,14,15,16]. One drawback is that the affective meanings of smiles and frowns may be ambiguous. A smile, for example, might express enjoyment, friendliness, and/or dominance [28].

For similar reasons, secondly, our focus is not on prosody (e.g. [4,24,29]) but on

formant frequencies (and hence formant dispersion) as the main auditory cue. In case of a human speaker, the pattern of formant frequencies may be regarded as the audible effect of a smile or frown gesture produced simultaneously with the speech [6, 8, 30, 31]. Other effects of smiles and frowns, mainly expressed prosodically by means of F0 [6, 8, 30, 31], are ignored in this study, because these prosodic effects cannot be easily related to motor resonance processes. Although this limitation in phonetic cues may result in a conservative study, we note that smiles and frowns are also perceived in whispered speech without F0 [8], and that spectral cues appear to be more important than F0 cues for perception of affect [32].

Thirdly, the effects of affective incongruence are investigated here not by means of acted speech (e.g. [22, 23, 24, 25]) but by means of synthesized speech in which formants were manipulated [30, 31]. This phonetic simulation of smiling and frowning allows stronger experimental control over the affective phonetic cues contributing to spoken word processing. Thus a word's affective meaning and its affective phonetic form were varied orthogonally, yielding congruent and incongruent combinations of affective meaning and form. The listeners' task involved language comprehension of positively and negatively valenced words. Words are predicted to be understood slower if spoken in an incongruent form (e.g., positive words with frown) than in a congruent form (e.g., positive words with smile).

2. Method

2.1. Stimulus words

Experimental stimuli consisted of 60 Dutch words (30 having positive meaning,

e.g. *eerlijk* “honest”, and 30 having negative meaning, e.g. *vijandig* “hostile”). This selection was based on a pre-test in which words were rated for affective value by a sample of 30 Dutch students. Positive words were rated more positively ($M=7.45$, $SD=0.49$) than negative words [$M=2.85$, $SD=0.50$, $t(58)=32.62$, $p<.001$] on a 9-point scale. Positive words had the same length in syllables ($M=2.6$, $SD=1.0$) as negative words [$M=2.6$, $SD=0.8$, $t(58)=0.29$, $p=.774$]. A male native speaker read each word in an affectively neutral manner, using a randomized list of words, and reading each word as a separate utterance (without list intonation). These readings were recorded and then used as targets for speech synthesis.

2.2. Stimulus preparation and selection

Spectral resonances (formants) were computed from the neutral speech recordings, and checked manually before being used for speech synthesis. The corrected formant values were used to control a formant-based speech synthesizer [33, 34]. For neutral phonetic forms, the unshifted frequencies of the formants were used. For smiling forms, the frequency of the lowest spectral resonance (formant F1) was shifted up by 5%, and frequencies of higher formants (F2 to F5) were shifted up by 10% [6]. Conversely, for frowning forms, the F1 was shifted down by 5%, and higher formants were shifted down by 10% [23]. Formants were adjusted throughout the target word. This resulted in phonetically neutral synthetic realizations (positive–neutral, negative–neutral), or congruent realizations (positive–smiling, negative–frowning), or incongruent realizations (negative–smiling, positive–frowning). All other synthesis parameters were identical in

corresponding neutral, congruent and incongruent forms of a word. The pitch contour was copied from the original recording.

2.3. Pre-tests

In order to verify the noticeability of the phonetic manipulations, as well as the resulting intelligibility of the target words, two pre-tests were conducted. In the first pre-test, listeners rated each spoken word on a 9-point scale, to indicate to what extent the word was spoken with a simultaneous frown (scale value 1), in neutral fashion (value 4), or with a simultaneous smile (value 9). The three phonetic forms (unshifted or “neutral”, with formants shifted down or “frowning”, and with formants shifted up or “smiling”) were counterbalanced over three lists. Each list was presented over loudspeakers to a separate group of listeners (of 42, 16, and 16 listeners, respectively), at an inter-stimulus interval of 2 s (as established in pilot tests). Listeners were native Dutch-speaking undergraduate students at Utrecht University, without any knowledge about the purpose of the study. They were asked to rate the presence of simultaneous affective gestures during speech production, on a scale with end points marked with a frowning face symbol (☹, scale value 1) and a smiling face (☺, value 9). These subjective ratings were analyzed by means of mixed-effects regression, with listeners and target words as crossed random effects [35,36]. In the fixed part of the model, the “frowning” and “smiling” forms were compared relative to the neutral (unshifted) phonetic form; semantic valence was also included as a fixed effect.

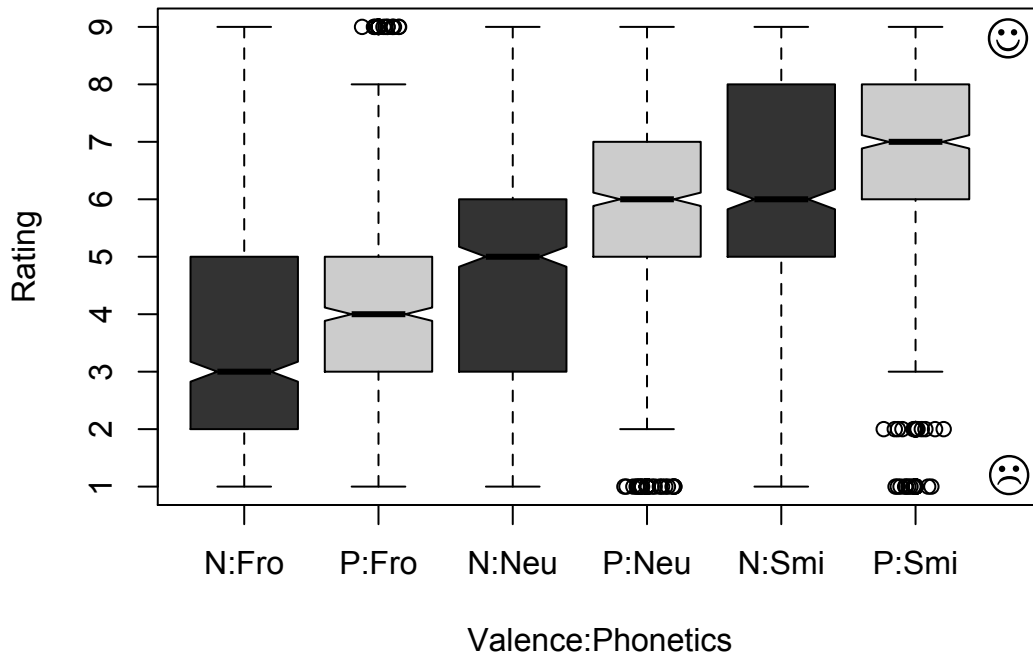


Figure 1. Boxplots of subjective ratings as to what extent words were spoken with a simultaneous smile, in neutral fashion, or with a simultaneous frown, broken down by semantic valence (N:negative, darker boxes, P:positive, lighter boxes) and by phonetic form (“frowning”, neutral, “smiling”). Notches indicate approximate 95% confidence intervals of the box median.

Subjective ratings were significantly lower for the “frowning” forms with formants shifted down (relative to “neutral” forms, $\beta=-1.54$ scale point, $s.e.=0.07$, $p=.0001$), and significantly higher for the “smiling” forms with formants shifted up ($\beta=1.00$ scale point, $s.e.=0.07$, $p=.0001$), as illustrated by the boxplots in Figure 1. This

confirms that the phonetic manipulations of formants in the speech stimuli do indeed successfully convey the desired affective property, viz. of speech being produced with a simultaneous frown gesture or smile gesture. Moreover, subjective ratings were also higher for semantically positive words than for semantically negative words ($\beta=0.78$ scale point, $s.e.=0.17$, $p=.0001$). Semantic valence thus yields a significant effect on subjective ratings of a phonetic affective property. This interesting main effect supports our research hypothesis, as it suggests that subjects' task of affect perception at the phonetic level was not entirely separate from semantic evaluation, thus indicating interacting processes. No interaction effect was observed between phonetic form and semantic valence on subjective ratings: the effects of phonetic manipulations are the same for negative and for positive words.

In the second pre-test, the phonetically neutral forms of the 60 target words were presented to 15 listeners (from the same sample of participants as below) in randomized order, to assess intelligibility of the synthesized target words. Participants listened to each resynthesized word (with unshifted formants) individually, and typed the word they had heard. Responses were scored for accuracy, with correction of occasional spelling errors (e.g. "interesant" for *interessant*). Participants' typed responses showed poor intelligibility (accuracy < .8) for 7 out of 60 target words (5 positive, 2 negative). These 7 poorly intelligible target words were kept in the main experiment but were excluded from further data analysis.

2.4. Participants and procedure

In the main experiment, 48 native Dutch-speaking students (39 females, 9 males)

with no hearing, language or speech deficits listened to the synthesized words. Participants' ages were between 18 and 27 years (median 21.5 years). The three phonetic forms of a word were balanced over three separate experimental lists. Each list was presented to 16 participants. Listeners' task was to classify the meaning of the spoken word as positive (exemplar *peace*) or negative (exemplar *war*) as quickly and as accurately as possible after the offset of the spoken word. Stimuli were presented with a 5-ms fade-in and fade-out to prevent click sounds. Before the actual task, participants were presented with 5 practice trials (including all 3 phonetic forms, both congruent and incongruent). The actual test started with 10 warm-up items indiscernible from the subsequent stimuli; the list of stimuli was re-randomized for each participant. Listeners responded by pressing one of two response buttons, always using the index finger of their dominant hand. Positive and negative response buttons were balanced over the 16 participants within each experimental list. No instructions were given about the relevant semantic or phonetic properties of the stimuli. The total time of an experimental session was about 12 minutes.

3. Results

The main dependent measure was response time (RT) measured from the onset of the spoken word. Responses with outlier RTs (4%) and incorrect responses (4%) were excluded from the data analysis. The remaining RTs were analyzed by means of mixed-effects regression, with listeners and target words as crossed random effects [24,25]. The resulting optimal model, summarized in Table 1, confirmed the predicted interference pattern. Responses were significantly slower for words with affectively incongruent

phonetic forms (positive–frowning or negative–smiling, see Figure 2) than for phonetically neutral words ($\beta=39.4$ ms, $s.e.=13.8$, $p=.004$). Responses for words with congruent phonetic forms (positive–smiling, negative–frowning) were equally fast as for neutral words ($\beta=1.5$ ms, $s.e.=13.8$, n.s.). RTs were faster for positive than for negative words ($\beta=120$ ms, $s.e.=33.5$, $p<.001$). Interactions were not significant, as confirmed by a likelihood ratio test [$\chi^2(2)<1$, n.s.].

Table 1

Estimated parameters of mixed-effects model of response times. Estimates of fixed parameters are given in ms, with standard error and significance level (in boldface if $p < .05$). Estimates of random parameters are given in standard deviations, with 95% confidence interval of the estimate. N=2343.

<i>fixed coefficients</i>	estimate	<i>s.e.</i>	<i>p</i>
Intercept	1264.2	30.6	.0001
Positive Valence Word	-119.7	33.5	.0002
Incongruent Phonetic Form	39.4	13.8	.0036
Congruent Phonetic Form	1.5	13.8	.9170
Pos. Valence × Incongruent	-4.2	20.1	.8336
Pos. Valence × Congruent	10.4	20.1	.6028
<i>random coefficients</i>	estimate	95% C.I.	
listeners	140.2	96.7, 136.6	
target words	110.4	81.0, 115.9	
residual	197.8	193.9, 205.4	

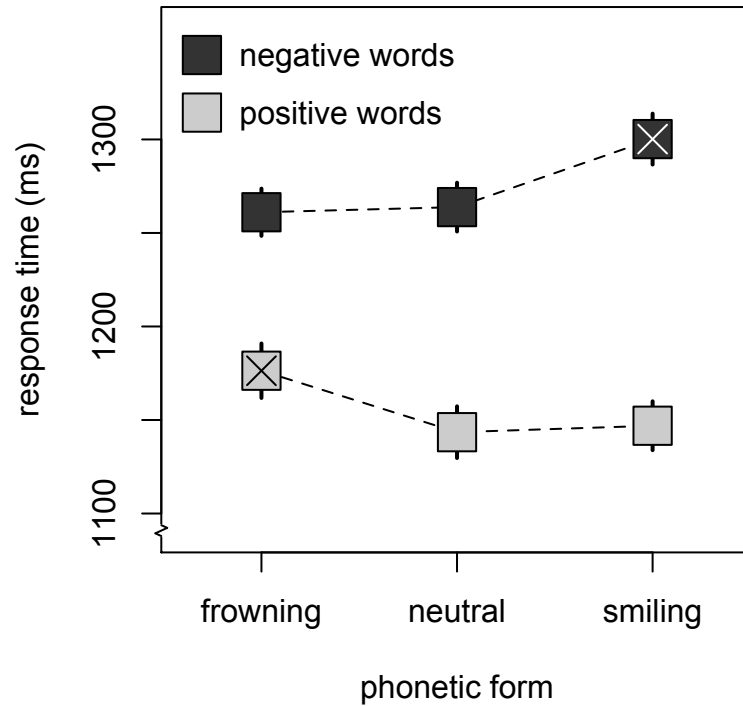


Figure 2. Mean response times for negative and positive target words, broken down by phonetic form. Error bars indicate standard error of the mean. Dashed lines are for visual guidance only. Conditions with incongruent phonetic forms are crossmarked.

4. Discussion

In the first pretest, listeners' phonetic judgements of affective stimulus words were influenced by the semantic valence of those words, yielding more "smiling" ratings for semantically positive words and more "frowning" ratings for negative words, across phonetic manipulation conditions. This interference effect of semantic valence on phonetic ratings supports the main hypothesis in this study.

Conversely, in the main experiment, listeners showed a significant impairment in their comprehension of these spoken words, as measured by a semantic classification task, if the phonetic form of a word was incongruent with its semantic valence. Although some neurological and behavioral studies have suggested a general dissociation between emotional auditory processing and linguistic processing (e.g. [37]), our results indicate that vocal expressions of affect are integrated with linguistic properties of an utterance, thus adding affective redundancy. When redundancy is not provided (i.e., in incongruent conditions), the utterance is more difficult to understand.

Previous work in the literature (e.g. [38,24]) investigated perceptual interference by varying affective congruency of prosody in natural speech tokens. Articulatory interference of acting speakers who simulated the intended emotions was a possible account for the effect. Namely, human actors may have more difficulty in producing natural tokens of affective words with incongruent emotions, as opposed to words with congruent emotions. To address this limitation, here synthetic speech was used. Therefore, the present interference effect cannot be attributed to articulatory interference in acted speech, nor can it be ascribed to visual affective expressions, because visual cues were absent, nor to priming effects [39,18] because effects were measured on the affected words themselves. Our results are most likely due to immediate resonance between the valence of the target word and the affective phonetic properties of that spoken word.

The present results are in line with embodied theories of language comprehension suggesting that listeners resonate to both acoustic and semantic information during spoken language comprehension (cf. [14]). The affective resonance observed in the present experiment corresponds with the recently proposed Simulation of Smiles Model

[28, cf. 40,41,17]) which claims that the meaning of a smile is conveyed by means of motor mimicry of the observed smile – extended here to include not only visually but also phonetically observed smiles – somewhat analogous to the Motor Theory of Speech Perception [42]. Similarly, imitation (i.e., speech mimicry) of a foreign accent improves comprehension of other utterances spoken in that accent [43], and listeners' facial expressions contribute to their recognition of auditorily presented emotions [44]. All these findings suggest that motor mimicry may contribute not only to speech comprehension, but also to affect perception. In our view, a perceived smile may elicit a smiling gesture (albeit a weak one) in the listener, which in turn interferes with speech comprehension through motor resonance [14] and/or through affective resonance [15,16].

5. Conclusion

In sum, comprehension of spoken language is not merely based on *what* is said, but also on *how* it is said – namely the affective facial expression coinciding with speech production. Speech comprehension is therefore an integrated process that benefits from the affective expressions that modulate how we say what we say. If speakers smile or frown while talking, then the audible effects of these affective cues influence our comprehension of spoken words.

Acknowledgements

We thank Sieb Nooteboom, Frank Wijnen, Marcel Schmeets and Karin Wagenaar for helpful comments and suggestions, and Theo Veenker and Iris Mulders for technical and logistical assistance.

References

- [1] Williams, C. E., & Stevens, K. N. (1972). Emotions and speech: Some acoustical correlates. *Journal of the Acoustical Society of America*, *52(4B)*, 1238-1250.
- [2] Frick, R. W. (1985). Communicating emotion: The role of prosodic features. *Psychological Bulletin*, *97(3)*, 412-429.
- [3] Neumann, R., & Strack, F. (2000). "Mood contagion": The automatic transfer of mood between persons. *Journal of Personality and Social Psychology*, *79(2)*, 211-223.
- [4] Scherer, K. R. (2003). Vocal communication of emotion: A review of research paradigms. *Speech Communication*, *40(1-2)*, 227-256.
- [5] Batliner, A., Fischer, K., Huber, R., Spilker, J., & Noth, E. (2003). How to find trouble in communication. *Speech Communication*, *40(1-2)*, 117-143.
- [6] Ohala, J. J. (1980). The acoustic origin of the smile. *Journal of the Acoustical Society of America*, *53*. Retrieved from <http://linguistics.berkeley.edu/phonlab/users/ohala/papers/smile.pdf>.
- [7] Ohala, J. J. (1983). Cross-language use of pitch: an ethological view. *Phonetica*, *40(1)*, 1-18.
- [8] Tartter, V. C., & Braun, D. (1994). Hearing smiles and frowns in normal and whisper registers. *Journal of the Acoustical Society of America*, *96(4)*, 2101-2107.
- [9] Drahotová, A., Costall, A., & Reddy, V. (2008). The vocal communication of different kinds of smile. *Speech Communication*, *50(4)*, 278-287.
- [10] Pell, M. D., Paulmann, S., Dara, C., Allasseri, A., & Kotz, S. A. (2009). Factors in the recognition of vocally expressed emotions: A comparison of four languages. *Journal of Phonetics*, *37(4)*, 417-435.

- [11] Gallese, V., Fadiga, L., Fogassi, L., & Rizzolatti, G. (1996). Action recognition in the premotor cortex. *Brain*, *119*, 593-609.
- [12] Kohler, E., Keysers, C., Umiltà, M. A., Fogassi, L., Gallese, V., & Rizzolatti, G. (2002). Hearing sounds, understanding actions: Action representation in mirror neurons. *Science*, *297*(5582), 846-848.
- [13] Wilson, S. M. W., Saygin, A. P., Sereno, M. I., & Iacoboni, M. (2004). Listening to speech activates motor areas involved in speech production. *Nature Neuroscience*, *7*, 701-702.
- [14] Zwaan, R. A., & Taylor, L. J. (2006). Seeing, acting, understanding: Motor resonance in language comprehension. *Journal of Experimental Psychology: General*, *135*(1), 1-11.
- [15] Gallese, V. (2003). The roots of empathy: The shared manifold hypothesis and the neural basis of intersubjectivity. *Psychopathology*, *36*(4), 171-180.
- [16] Gallese, V. (2009). Motor abstraction: a neuroscientific account of how action goals and intentions are mapped and understood. *Psychological Research*, *73*(4), 486-498.
- [17] Niedenthal, P. M. (2007). Embodying emotion. *Science*, *316*(5827), 1002-1005.
- [18] Foroni, F., & Semin, G. R. (2009). Language that puts you in touch with your bodily feelings: The multimodal responsiveness of affective expressions. *Psychological Science*, *20*(8), 974-980.
- [19] Stroop, J. R. (1935). Studies of interference in serial verbal reaction. *Journal of Experimental Psychology*, *18*, 643-662.
- [20] Schirmer, A., & Kotz, S. A. (2003). ERP evidence for a sex-specific Stroop effect in emotional speech. *Journal of Cognitive Neuroscience*, *15*(8), 1135-1148.
- [21] Mehrabian, A., & Wiener, M. (1967). Decoding of inconsistent communications. *Journal of Personality and Social Psychology*, *6*(1), 109-114.

- [22] Grimshaw, G. M. (1998). Integration and interference in the cerebral hemispheres: Relations with hemispheric specialization. *Brain and Cognition*, 36(2), 108-127.
- [23] Schirmer, A., Kotz, S. A., & Friederici, A. D. (2002). Sex differentiates the role of emotional prosody during word processing. *Cognitive Brain Research*, 14, 228-233.
- [24] Nygaard, L. C., & Queen, J. S. (2008). Communicating emotion: Linking affective prosody and word meaning. *Journal of Experimental Psychology: Human Perception and Performance*, 34(4), 1017-1030.
- [25] Paulmann, S., Titone, D., & Pell, M. D. (2012). How emotional prosody guides your way: Evidence from eye movements. *Speech Communication*, 54(1), 92-107.
- [26] Van Berkum, J. J. A., van den Brink, D., Tesink, C. M. J. Y., Kos, M., & Hagoort, P. (2007). The neural integration of speaker and message. *Journal of Cognitive Neuroscience*, 20(4), 580-591.
- [27] Tesink, C. M. J. Y., Petersson, K. M., van Berkum, J. J. A., van den Brink, D., Buitelaar, J. K., & Hagoort, P. (2008). Unification of speaker and meaning in language comprehension: An fMRI study. *Journal of Cognitive Neuroscience*, 21(11), 2085-2099.
- [28] Niedenthal, P. M., Mermillod, M., Maringer, M., & Hess, U. (2010). The Simulation of Smiles (SIMS) model: Embodied simulation and the meaning of facial expression. *Behavioral and Brain Sciences*, 33(6), 417-433.
- [29] Schröder, M. (2006). Expressing degree of activation in synthetic speech. *IEEE Transactions on Audio and Speech Processing*, 14(4), 1128-1136.
- [30] Chuenwattanapranithi, S., Xu, Y., Thipakorn, B., & Maneewongvatana, S. (2008). Encoding emotions in speech with the size code. *Phonetica*, 65(4), 210-230.
- [31] Lasarcyk, E., & Trouvain, J. (2008). Spread lips + raised larynx + higher F0 = smiled

speech? An articulatory synthesis approach. Paper presented at the 8th International Speech Production Seminar (ISSP), Strasbourg.

[32] Xu, Y., & Kelly, A. (2010). Perception of anger and happiness from resynthesized speech with size-related manipulations. Paper presented at the Speech Prosody 2010 Conference, Chicago. Retrieved from http://www.isca-speech.org/archive/sp2010/papers/sp10_027.pdf

[33] Klatt, D. H. (1980). Software for a cascade/parallel formant synthesizer. *Journal of the Acoustical Society of America*, 67(3), 971-995.

[34] Boersma, P., & Weenink, D. (2011). Praat: Doing phonetics by computer [Computer software]. Retrieved from <http://www.praat.org>.

[35] Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language*, 59(4), 390-412.

[36] Quené, H., & Van den Bergh, H. (2008). Examples of mixed-effects modeling with crossed random effects and with binomial data. *Journal of Memory and Language*, 59(4), 413-425.

[37] Scott, S. K., Young, A. W., Calder, A. J., Hellawell, D. J., Aggleton, J. P., & Johnsons, M. (1997). Impaired auditory recognition of fear and anger following bilateral amygdala lesions. *Nature*, 385(6613), 254-257.

[38] Kitayama, S., & Ishii, K. (2002). Word and voice: Spontaneous attention to emotional utterances in two languages. *Cognition & Emotion*, 16(1), 29-59.

[39] Klauer, K. C., & Musch, J. (2003). Affective priming: Findings and theories. In J. Musch & K. C. Klauer (Eds.), *The Psychology of Evaluation: Affective Processes in Cognition and Emotion* (pp. 7-49). Mahwah, NJ: Lawrence Erlbaum Ass.

[40] Hietanen, J. K., Surakka, V., & Linnankoski, I. (1998). Facial electromyographic responses

to vocal affect expressions. *Psychophysiology*, 35(5), 530-536.

[41] Dimberg, U., Thunberg, M., & Elmehed, K. (2000). Unconscious facial reactions to emotional facial expressions. *Psychological Science*, 11(1), 86-89.

[42] Galantucci, B., Fowler, C. A., & Turvey, M. T. (2006). The motor theory of speech perception reviewed. *Psychonomic Bulletin & Review*, 13(3), 361-377.

[43] Adank, P., Hagoort, P., & Bekkering, H. (2010). Imitation improves language comprehension. *Psychological Science*, 21(12), 1903-1909.

[44] Hawk, S. T., Fischer, A. H., & Van Kleef, G. A. (2012). Face the noise: Embodied responses to nonverbal vocalizations of discrete emotions. *Journal of Personality and Social Psychology*, 102(4), 796-814.