

1 **Prosody conveys speaker's intentions: Acoustic cues for speech act**
2 **perception**

3 **Nele Hellbernd^{a*}, Daniela Sammler^a**

4 ^a Otto Hahn Group „Neural Bases of Intonation in Speech“, Max Planck Institute for Human Cognitive
5 and Brain Sciences, Leipzig, Germany

6 *** Correspondence:**

7 Nele Hellbernd

8 Otto Hahn Group „Neural Bases of Intonation in Speech“

9 Max Planck Institute for Human Cognitive and Brain Sciences

10 Stephanstraße 1a

11 04103 Leipzig, Germany

12 hellbernd@cbs.mpg.de

13 phone: +49 341 9940 2533

14 fax: +49 341 9940-2204

15 Daniela Sammler

16 sammler@cbs.mpg.de

17

18 **Keywords: Prosody, Intention, Speech acts, Acoustics, Pragmatics**

19 Number of Words: 8547

20 Number of Tables: 4

21 Number of Figures: 3

22 Number of appendices: 3

23 Supplementary Material: 6 Audio files

24 **Abstract**

25 Action-theoretic views of language posit that the recognition of others' intentions is key to
26 successful interpersonal communication. Yet, speakers do not always code their intentions
27 literally, raising the question of which mechanisms enable interlocutors to exchange
28 communicative intents. The present study investigated whether and how prosody—the vocal
29 tone—contributes to the identification of “unspoken” intentions. Single (non-)words were
30 spoken with six intonations representing different speech acts—as carriers of communicative
31 intentions. This corpus was acoustically analyzed (Experiment 1), and behaviorally evaluated in
32 two experiments (Experiments 2 and 3). The combined results show characteristic prosodic
33 feature configurations for different intentions that were reliably recognized by listeners.
34 Interestingly, identification of intentions was not contingent on context (single words), lexical
35 information (non-words), and recognition of the speaker's emotion (valence and arousal).
36 Overall, the data demonstrate that speakers' intentions are represented in the prosodic signal
37 which can, thus, determine the success of interpersonal communication.

38

39 Introduction

40 During conversations, humans regularly decode not only *what* is said but also *why*
41 (Bühler, 1934; Grice, 1957; Wittgenstein, 1953). Depending on the latter, we may understand the
42 same statement „It’s hard to be punctual in the morning“ as empathic concern, criticism, or
43 simply as a matter of facts. Pragmatic theory posits that it is particularly the *why*—the
44 communicative intention of the speaker—that drives the recipient’s behavior and is the motive of
45 communication. Yet, how intentions are (de)coded in interpersonal communication is still not
46 fully understood. Contemporary pragma-linguistic theories posit that listeners identify the
47 speaker’s goal via pragmatic inference (Wilson & Sperber, 2012), taking conversation context
48 and “common ground” (Clark & Carlson, 1981; Levinson, 2013; Stalnaker, 2002; Tomasello,
49 2005; Wichmann, 2002) into account. Alternatively, other studies seek to identify extralinguistic
50 cues that reveal a speaker’s intention, such as facial expressions (Fridlund, 1994; Frith, 2009;
51 Parkinson, 2005), properties of biological motion (Di Cesare, Di Dio, Marchi, & Rizzolatti,
52 2015), or gestures (Bucciarelli, Colle, & Bara, 2003; Enrici, Adenzato, Cappa, Bara, &
53 Tettamanti, 2011). The present study will focus on speech prosody—the tone of the voice—and
54 will weigh its potential to convey communicative intentions.

55 The question of how interlocutors decode the *why* of an utterance is grounded in *action-*
56 *theories of language*. In the middle of the 20th century, scholars like Karl Bühler (1934), Ludwig
57 Wittgenstein (1953), or Paul Grice (1975) recognized that language is more than strings of
58 symbols that are understood by retrieving their conventional, *coded* meaning. In their view,
59 language is an *intentional action* and gains meaning through its employment. Utterances become
60 instruments to influence the behavior of the interlocutor. The meaning of an utterance must be
61 found in its underlying intention. It was Grice (1957) who particularly promoted the central role

62 of intentions in communication. He advocated the idea that intentions drive speakers' behaviors
63 (e.g., utterances) whose sole function is to have an effect on the addressee by virtue of having
64 their intention recognized (cf. Levinson, 2006). Notably, the intention of the speaker—the
65 *speaker meaning* in Grice's terms—not necessarily surfaces in the overt lexical content of the
66 utterance, as shown in the example on punctuality above, but needs to be interpreted by the
67 listener.

68 This idea later became central to speech act theory by John L. Austin (1962) and John R.
69 Searle (1969) who considered utterances as actions—or *speech acts*—with specific interpersonal
70 goals such as promising, apologizing, or warning. Like Grice, they claimed that speakers convey
71 information on at least two levels: (1) the *propositional content* carrying the lexical meaning of
72 *what* is said, and (2) the *illocutionary force* representing the action and speaker's intention—the
73 *why*. As mentioned above, it is this second level—what the speaker is attempting to accomplish
74 with a remark—that is thought to predominantly drive the interlocutor's (conversational)
75 reaction. Notably, illocutionary force is often expressed implicitly (i.e. without the performative
76 verb) or even indirectly, hence requiring some sort of inference on the part of the listener
77 (Austin, 1962; Bach, 1994).

78 Interestingly, the notion of implicitness and indirectness conflicts with Grice's
79 *cooperative principle* (1975), which describes principles for effective communication in
80 conversation in four maxims. Following his *maxim of manner*, speakers ought to shape their
81 utterances in ways that support the purpose of the conversation. Hence, speakers should produce
82 unambiguous cues that make their intentions comprehensible to listeners. The fact that this seems
83 often not to be the case but listeners still efficiently recognize the speaker's intent has fueled
84 research on the cognitive and neural bases of comprehending communicative intentions. A great

85 deal of work has focused on implicit speech acts, i.e. utterances that express the speaker's
86 intention and illocutionary force without inclusion of the performative verb (e.g. "I will be
87 there." expressing a promise without including the verb "promise"). These studies demonstrated
88 the psychological reality of speech acts (Holtgraves, 2005), their automatic (Holtgraves, 2008a;
89 Liu, 2011) and early recognition during conversation turns (Egorova, Pulvermüller, & Shtyrov,
90 2014; Egorova, Shtyrov, & Pulvermüller, 2013; Gisladdottir, Chwilla, & Levinson, 2015), and
91 their importance for conversation memory (Holtgraves, 2008b). However, despite their
92 importance for understanding human communication, these studies remain incomplete in one
93 particular way: They often rely on written linguistic material and, thus, miss out on
94 extralinguistic cues that are usually available during natural spoken conversations. These cues
95 comprise signals expressed via additional communicative channels like eyes, face, body, or voice
96 and may render the speaker's intention less implicit and indirect than typically thought. The
97 present study will focus on vocal acoustic cues, i.e., prosody, as one non-verbal channel in
98 interpersonal conversation that may play an important role for speakers and listeners to express
99 and recognize communicative intentions.

100 The term prosody refers to variations in pitch, loudness, timing, or voice quality over the
101 course of an utterance (Warren, 1999) that can modify the communicative content of a message,
102 both linguistically and paralinguistically (Bolinger, 1986). Linguistically, prosody has direct
103 effects on the information structure of an utterance. It conveys, for example, semantic
104 relationships (Cutler, Dahan, & van Donselaar, 1997; Wagner & Watson, 2010), disambiguates
105 the syntactic constituent structure (Carlson, Frazier, & Clifton, 2009), and marks declarative vs.
106 interrogative sentence mode (Schneider, Lintfert, Dogil, & Möbius, 2006; Sammler, Grosbras,
107 Anwander, Bestelmeyer, & Belin, 2015; Srinivasan & Massaro, 2003). Paralinguistically, the

108 “manner of saying” conveys additional information that goes beyond the linguistic content.
109 Whether or not this includes intentions is a matter of debate (Bolinger, 1986) and will be topic of
110 the present research.

111 Until now, most studies on paralinguistic prosody either focused on the speaker’s
112 emotion (Banse & Scherer, 1996; Bänziger & Scherer, 2005; Frick, 1985; Simon-Thomas,
113 Keltner, Sauter, Sinicropi-Yao, & Abramson, 2009) or, more recently, on their attitude, for
114 example, the politeness, confidence, or sincerity of the speaker (Jiang & Pell, 2015; Monetta,
115 Cheang, & Pell, 2008; Rigoulot, Fish, & Pell, 2014) and often sought to determine links between
116 the acoustics of the prosodic signal and the listeners’ comprehension of the paralinguistic
117 message. Although opinions diverge on whether prosody as such can convey meaning, i.e.
118 without contextual information (see below) (Cutler, 1976; Wichmann, 2000, 2002), studies
119 revealed distinct acoustic properties for the prosodic expression of different emotions (Banse &
120 Scherer, 1996, Szameitat et al., 2009) and attitudes (Blanc & Dominey, 2003; Morlec, Bailly, &
121 Aubergé, 2001; Uldall, 1960). Similarly, on the perception side, researchers showed that
122 participants were able to identify the speaker’s attitude (Morlec et al., 2001; Uldall, 1960) and
123 emotion by prosodic differences alone, in verbal (Banse & Scherer, 1996; Morlec et al., 2001)
124 and non-verbal utterances (Monetta et al., 2008; Sauter, Eisner, Calder, & Scott, 2010a), in
125 laughter (Szameitat et al., 2009), and to some extent even cross-culturally (Sauter, Eisner,
126 Ekman, & Scott, 2010).

127 Compared to this active field of research, only little is known about the perceptual reality,
128 relevance and effectiveness of prosodic cues in conveying *intentions*. We consider
129 communicative intentions as the goals of interpersonal actions (e.g., language) that are meant to
130 be recognized by the interlocutor and to influence her (conversational) reactions. This

131 differentiates communicative intentions from basic emotions that do not necessarily need another
132 person to be displayed, and attitudes that are not necessarily meant to purposefully influence
133 conversation partners (Wichmann, 2000). Certainly, both emotions and attitudes can be
134 expressed for communicative purposes (Fridlund, 1994; Mead, 1934; Parkinson, 2005) and often
135 take an effect on the listener by virtue of their “expressive function” (Bühler, 1934). Yet, their
136 intended goal remains rather underspecified compared to the “specific intentions for specific
137 turns” (cf. Holtgraves, 2008a) proposed by action-theoretic accounts of language, particularly by
138 speech act theory (Austin, 1962; Searle, 1969).

139 To date, the role of prosody for the non-literal expression and recognition of different
140 intentions still lacks detailed investigation, although several findings from developmental studies
141 and psycholinguistics point to the relevance of extralinguistic vocal cues in intentional
142 communication. For quite some time, studies on intonational development have been focusing on
143 the emergence of illocutionary skills in infants, considering intonation patterns as primitive
144 devices that preverbal infants use to express their communicative intentions (Dore, 1975). For
145 example, 7- to 11-month-old babies were found to vocally distinguish between communicative
146 and investigative (Papaeliou, Minadakis, & Cavouras, 2002) or emotional functions when
147 babbling (Papaeliou & Trevarthen, 2006). This competence was proposed to regulate cooperative
148 interactions with their parents as a prerequisite for language acquisition. Furthermore, infants’
149 intonations of babble at the end of their first year (Esteve-Gibert & Prieto, 2013), or words in
150 their second year of life (e.g., Furrow, Podrouzek, & Moore, 1990; Marcos, 1987; Prieto,
151 Estrella, Thorson, & Vanrell, 2012) were found to differ between simple speech acts such as
152 complaining, requesting, or greeting. These combined findings were taken as evidence for a
153 prosodic choice that prelinguistic infants make to communicate their intentions (illocutions)

154 while their propositional (locutionary) abilities are still limited. One challenge that these studies
155 have to face, though, is their dependency on adult, post-hoc interpretations of infants' vocal
156 actions that are usually based on the context in which the vocalizations were produced. This
157 bears the risk that raters—although experts (e.g., mothers or phoneticians)—might overestimate
158 or misinterpret the children's (true) motives or draw conclusions from cues other than prosody.
159 Studies with adult speakers who can report on their intentions are necessary to corroborate the
160 link between prosody and communicative intentions, and to show its persistence in adulthood.

161 The present study aimed to fill this gap by conceptualizing a speaker's intention in terms
162 of speech acts (Austin, 1962; Holtgraves, 2002; Searle, 1969) and investigating the role of
163 prosody in decoding illocutionary force. Note that it was not our goal to describe the prosody of
164 a complete set of speech acts or to investigate the reality of speech act theory. Rather, we aimed
165 to demonstrate that—in identical utterances pronounced according to a limited set of
166 intentions—speakers produce well-identifiable characteristic prosodic patterns, and that these
167 patterns can be reliably recognized by listeners. This adds to the debate whether prosody can
168 convey meaning on its own, i.e., may be conventionalized for different communicative concepts.
169 Alternative views regard prosody as a contrastive marker that does not carry meaning by itself
170 but signals the presence of “unspoken” meaning by deviating from normal prosody, and hence,
171 motivates listeners to infer the implied message by taking context information into account
172 (Cutler & Isard, 1980; Levinson, 2013). Here, we tested the hypothesis that prosodic patterns as
173 such can be sufficiently distinct, to a degree that listeners can recognize the broad
174 communicative concept and intention in the prosodic speech signal. Therefore, our stimulus set
175 comprised single words and non-words, i.e., tokens free of context and lexical meaning, that
176 were pronounced with six different intonations representing the speech acts criticism, doubt,

177 naming, suggestion, warning, and wish. In three experiments, we combined acoustic analyses of
178 these speech signals with perceptual judgments of listeners (for a similar approach, see Banse &
179 Scherer, 1996; Sauter, Eisner, Calder, et al., 2010). If prosody itself codes speakers' intentions,
180 different speakers should employ similar cue configurations when conveying the same intention,
181 and participants should be able to recognize the intention without contextual information (i.e., in
182 single words) and irrespective of whether the speech sound carries lexical meaning or not (i.e., in
183 words and non-words).

184 One important consideration for our investigations of communicative intentions in
185 prosody is the relation to emotional components in the speaker's tone of voice. Although we
186 advocated a conceptual differentiation of intentions and emotions above, we have to keep in
187 mind that emotions (e.g., fear) might drive intentions (e.g., to warn the interlocutor). Hence, both
188 may be intertwined in the production and perception of communicative utterances. In an attempt
189 to show that the comprehension of intentions is more than the recognition of emotions or affect
190 in the prosodic signal, we further assessed the valence and arousal of our speech stimuli
191 according to dimensional models of affect (Remington, Fabrigar, & Visser, 2000; Wundt,
192 1896). (We will use the term emotion throughout the text to refer to these affect measures).
193 These values were then used to correct the perceptual recognition of intentions for the
194 contribution of emotion (see below).

195 The present study took three steps: We started with analyses of the acoustics of speakers'
196 vocal expressions of speech acts by means of discriminant analyses (Experiment 1). If speech
197 acts as carriers of intentions are coded in characteristic prosodies (i.e. show some consistency of
198 the prosodic pattern across speakers and across tokens within speakers), it should be possible to
199 classify the different categories of speech acts based on their acoustic features alone, in words

200 and non-words alike. Second, we tested whether listeners are able to identify the correct intention
201 based on the prosodic pattern alone (Experiment 2) in a 6-alternative forced choice (6-AFC)
202 categorization task and ratings of the stimuli on every speech act scale (e.g., “How much does it
203 sound like criticism?”). If prosody conveys meaning in a partly conventionalized way, listeners
204 should be able to classify the intentions despite lack of context (i.e., in single words) and
205 irrespective of lexical meaning (i.e., similarly in words and non-words). Finally, we determined
206 which acoustic parameters contribute most to the perception of the respective intention
207 (Experiment 3). Therefore, we fed the acoustic parameters into multiple regression analyses to
208 predict the participants’ ratings on each speech act scale. Furthermore, to control for a possible
209 influence of emotion on intention recognition, the regression analyses were repeated once after
210 valence and arousal ratings of the stimuli had been regressed out.

211 In summary, the present study sought to demonstrate that prosody carries information
212 about the speaker’s communicative intention by (i) identifying characteristic prosodic feature
213 configurations of a set of speech acts that are (ii) reliably recognized by listeners, (iii) despite the
214 lack of context information (single words) and semantic content (non-words) and the control for
215 emotional processing of the stimuli.

216 **Experiment 1 – Acoustics**

217 The goal of Experiment 1 was to investigate whether speakers use characteristic acoustic
218 features to convey their intentions. If so, it should be possible to classify the speech stimuli into
219 the corresponding speech act categories based on their acoustic features alone and irrespective of
220 word meaning. Specifically, we focused on duration, intensity, pitch and spectral features that
221 have been analyzed in similar approaches in emotion research (e.g. Banse & Scherer, 1996;

222 Blanc & Dominey, 2003; Sauter, Eisner, Calder, & Scott, 2010b). In such studies, pitch cues
223 were predominant when emotions were expressed verbally, compared to a stronger weighting of
224 spectral features in non-verbal utterances, making it likely that pitch cues will play a major role
225 in the present experiment.

226 For the current study, four speakers produced single-word stimuli with varying prosodies
227 to express six different intentions, i.e., the speech acts criticism, doubt, naming, suggestion,
228 warning, and wish. To obtain stimuli that are representative for typical language use, all speakers
229 were non-actors, i.e., they relied on their intuition—not training in acting—to express the
230 intention in a way that could be understood by an imaginary interlocutor. For high stimulus
231 quality, all speakers were, however, familiar with sound recordings, i.e., working as voice
232 coaches or speech scientists. This choice of professional speakers with only minimal training in
233 acting is an attempt to face the criticism, first raised in emotion research, that actors' prosodic
234 patterns may deviate from those used in everyday conversations (Jürgens, Hammerschmidt, &
235 Fischer, 2011; see also General Discussion). Apart from that, it should be mentioned that
236 intentions are typically expressed more voluntarily than emotions and are, hence, less dependent
237 on the spontaneity of the utterance. Altogether, the present stimuli were recorded such to grant
238 generalizability of the results to natural language use.

239 *Materials and Methods*

240 *Ethics Approval.*

241 The ethics committee of the University of Leipzig, Germany approved the present and all
242 following experiments in this study.

243 *Stimulus recordings.*

244 Four trained native German speakers (voice coaches, 2 female) were invited to record the
245 German words “Bier” (*beer*) and “Bar” (*bar*) as well as the non-words “Diem” and “Dahm” (for
246 examples, see Supplementary Material). These (non-)words were intoned to express six different
247 communicative intentions or speech acts: criticism, doubt, naming, suggestion, warning, and
248 wish. The chosen speech acts were plausible for our stimulus words “beer” and “bar” and fit into
249 the broader speech act categories as defined by Searle and Vanderveken (1985). To elicit the
250 respective intentions in the speakers, they read short scenarios that described a situation in which
251 they interacted with an interlocutor (see Appendix A). They were allowed to utter an initial
252 sentence and to vocalize freely until they felt ready to articulate the intention shortened to the
253 single essential word. This recording approach, instead of using natural speech recordings, was
254 chosen to obtain clear portrayals of intentions in good sound quality. Recordings were conducted
255 in a soundproof room with the microphone (Rode NT55) approximately 20 cm in front of the
256 speaker and digitized at a 44.1 kHz sampling rate in a 16-bit mono format. The words and non-
257 words were repeated several times to obtain eight variants per stimulus in good quality. The
258 resulting stimulus set, thus, comprised 768 stimuli, with eight repetitions of four (non-)words
259 expressed as six speech acts by four speakers.

260 *Acoustic features.*

261 For investigating acoustic features of the speech acts, we obtained seven acoustic
262 measures that are commonly used in experiments on human voice and speech stimuli (e.g.,
263 Banse & Scherer, 1996; Sauter et al., 2010). Using Praat software (Boersma & Weenink, 2014)
264 we extracted the number of voiced frames as a measure of stimulus duration, mean intensity,
265 harmonics-to-noise ratio (HNR), mean fundamental frequency (f_0) as well as pitch rise,

266 measured as the difference between offset and onset f0. Furthermore, we extracted the spectral
267 center of gravity and the standard deviation of the spectrum. The mean acoustic characteristics as
268 measured with Praat are presented in Table C1 (Appendix C). Statistical analyses showed that
269 speakers had used very similar acoustic cues to intone speech acts in words and non-words. T-
270 tests for paired samples comparing the acoustics of words and non-words for each speech act
271 category were largely non-significant. Only exception were HNR and spectral center of gravity
272 that showed differences in some, but not all speech act categories (see Table C2 in Appendix C
273 for more details). These differences are, however, likely to be caused by the different consonants
274 (“r” in words vs. “m” in non-words) rather than by differences in prosody.

275 *Discriminant analyses.*

276 Discriminant analyses were performed for words and non-words separately, with the
277 seven acoustic features as independent variables and the speech act category (criticism, doubt,
278 etc.) as dependent variable. These analyses sought to identify linear functions of acoustic feature
279 combinations that maximize differences between speech act categories. In other words, these
280 analyses tested whether the acoustic features alone have sufficient discriminant power to reliably
281 group the stimuli that express the same intention. The discriminant analyses for both words and
282 non-words were cross-validated with a jack-knife procedure, and the distribution of the results
283 was validated with chi-square tests.

284 *Results*

285 The discriminant analyses classified the correct speech act category for 92% of all word
286 stimuli and 93% of all non-word stimuli. These results are highly above chance-level (17%) as
287 was tested with chi-square tests: $\chi^2(35) = 1717.6$, $p < .001$ for words and $\chi^2(35) = 1702.6$, $p <$

288 .001 for non-words. Classification results of the discriminant analyses for the different types of
289 speech acts are demonstrated as confusion matrices in Table 1. The highest results were found
290 for naming and warning (both 100% correct classification for words and non-words), while the
291 lowest results were obtained for criticism in words (76.6% correct classification), and non-words
292 (79.7%). Additional chi-square tests showed that the discriminant model classified our stimuli
293 better than chance (chance-level: 17%) for every type of speech act ($\chi^2(5) > 151$, p 's $< .001$).
294 Figure 1 shows the classification of the different speech acts by the first two discriminant
295 functions. The first function (x-axes) explained 49.6% of the variance for words and 54.2% for
296 non-words and was mainly based on the acoustic measure of pitch rise (offset – onset f_0). The
297 second discriminant function (y-axes) had an additional discriminant power of 38.6% for words
298 and 36.4% for non-words and was most related to the mean intensity and mean f_0 of the stimulus
299 (see Table C3, Appendix C). Additionally, a third function (not depicted in Figure 1 for reasons
300 of clarity) explained 10.5% of variance for words and 7.7% for non-words and showed highest
301 correlation with the duration of the stimuli. The last two discriminant functions from our
302 analyses explained only minor effects (function 4: 0.9% for words and 1.2% for non-words,
303 function 5: 0.4% for words and 0.5% for non-words) and were neglected from further
304 investigations.

305 Table 1

306 Figure 1

307 *Discussion*

308 The acoustic features of our stimulus set could be used to accurately classify the correct
309 speech act, for words and non-words alike. This demonstrates the distinctiveness of the prosodic

310 patterns that speakers deliberately applied to code their intentions in the tested speech act
311 categories and the relative independence of prosody from lexical content. Furthermore, the high
312 accuracy of the classification implies a reasonable consistency of the relevant prosodic cues
313 across speakers and utterances, and may point to the existence of feature configurations that
314 speakers consider conventional and appropriate for different communicative goals. For example,
315 the warning stimuli were loudest and had the most arched pitch contour with a salient peak in the
316 middle of the word as is appropriate for the urgent nature of a warning. In comparison, the
317 naming stimuli showed the least salient acoustic features with low mean pitch, flat pitch contour,
318 low intensity and little spectral variation in line with the neutral character of the expression. As
319 expected, pitch rise and mean f0, together with mean intensity were the most influential acoustic
320 features in these analyses, while spectral features had only weak discriminant power. In sum, the
321 data show that speakers can use prosody as a channel of communication to convey their
322 intentions. Note that we do not expect that speakers possess different prosodic patterns for all
323 possible intentions or speech acts. Yet, we believe that speakers choose salient, distinguishable
324 and probably culturally learned prosodic signatures to trigger cognitive processes in the
325 addressee to infer the communicative intent of the speaker beyond the overt lexical meaning.

326 **Experiment 2 – Behavior**

327 After finding consistent acoustic differences between prosodic speech act expressions, we
328 were interested in participants' perception of the stimuli. We investigated whether participants
329 would be able to identify the different intentions based on the prosodic information in a 6-
330 alternative forced-choice (6-AFC) categorization task. Participants, further, judged the valence
331 and arousal of every stimulus (Remington et al., 2000; Russell, 1980; Wundt, 1896) in the

332 second half of the experiment, which allowed us to assess in how far speech acts may be
333 classified based on their emotional tone.

334 *Materials and Methods*

335 *Participants.*

336 Ten participants were presented with the word stimuli (6 female, mean age \pm SD: 24.6 \pm
337 4.9), ten other volunteers performed the task with the non-word stimuli (4 females, mean age \pm
338 SD: 24.9 \pm 2.6). We tested separate groups of participants for words and non-words to avoid
339 transfer of the semantic meaning to the non-word stimuli. All participants reported normal
340 hearing ability, gave written informed consent and were paid 7€ per hour for their participation.

341 *Design and procedure.*

342 In the first half of the experimental session, participants were asked to assign each
343 stimulus to one of the six possible speech act categories (criticism, doubt, naming, suggestion,
344 warning, or wish). After having read short definitions for the different speech acts (Appendix B),
345 they heard each sound stimulus once via headphones and were instructed to press the keys 1–6
346 on a keyboard. The speech act labels with corresponding numbers were displayed on a computer
347 screen throughout the experiment. No feedback for the correctness of the response and no time
348 limits were given. The experiment was separated into four blocks—one for each speaker. Block
349 order and stimulus order within each block were pseudo-randomized by preallocating the speech
350 acts with balanced probabilities. Chi-square tests were performed to test for above-chance
351 classification across all speech acts and within single speech act categories.

352 In the second half of the experimental session, participants were asked to evaluate the
353 valence and arousal of each stimulus. Therefore, they listened to the same stimuli again, in the

354 same order as before. After each sound, they saw two visual analogue scales on the screen, first
355 for valence (positive/negative), then for arousal (calm/excited), and placed their ratings with a
356 continuous slider. The scales showed the outermost pictures of the Self-Assessment Manikin
357 (Bradley & Lang, 1994) at the margins. No time constraint was given for the answers. Friedman
358 tests were calculated to examine differences in the affect ratings among the speech act
359 categories.

360 *Results*

361 *Speech act categorization*

362 In the 6-AFC task, participants were able to identify the correct speech act category of
363 our stimuli with high accuracy for words (mean \pm SD: 82 \pm 13%) and for non-words (73 \pm
364 17%)—with no significant difference between the participant groups for words and non-words
365 ($t(18) = 1.26, p = .22$). Chi-square tests showed that participants' classification of every speech
366 act category was better than being predicted by chance (chance level: 17%), both for words
367 ($\chi^2(5) > 1082, p's < .001$) and non-words ($\chi^2(5) > 798, p's < .001$).

368 Confusion matrices for words and non-words are presented in Table 2. As with the
369 acoustic analyses, the identification of criticism was lowest among the six speech act categories.
370 For words as well as non-words, participants misclassified criticism most often as being doubt,
371 and to a lower extent as warning. Furthermore, common confusions of the non-word stimuli were
372 found for suggestion taken as doubt, wish, or criticism. To some extent, participants also
373 misclassified wish as suggestion, and naming as criticism.

374 Table 2

375 *Emotion ratings*

376 For the perception of emotion, mean ratings for valence and arousal differed significantly
377 between the speech act categories, for words (valence: $\chi^2(5) = 35$, arousal: $\chi^2(5) = 43$, p 's < .001)
378 and non-words (valence: $\chi^2(5) = 44$, arousal: $\chi^2(5) = 45$, p 's < .001). The results were very
379 similar for words and non-words in each speech act category (Fig. 2). On the valence scale, the
380 speech acts warning and criticism were perceived most negatively, whereas wish and suggestion
381 were associated with a more positive valence. Doubt and naming were rated neutrally with
382 regard to valence. The perception of the speakers' arousal was very calm for naming, wish, and
383 doubt, and very excited for warning and criticism. Suggestion stimuli were rated in the middle
384 range for arousal.

385 Figure 2

386 *Discussion*

387 Participants were well able to identify the speaker's communicative intention from the
388 prosody alone as indicated by the highly significant results in the 6-AFC categorization task.
389 Importantly, participants were able to make use of the prosodic signal with minimal context
390 descriptions and without lexical content (see below). These data show that prosody is a powerful
391 communicative channel that is used by listeners to decode the "unspoken" meaning and intention
392 of the speaker and that may determine their respective conversational reaction. Interestingly,
393 criticism was identified least reliably and was specifically confused with doubt, in line with the
394 very similar acoustic features of these two speech acts (Fig. 1, Table 1). It is well conceivable
395 that the acoustic similarity of criticism and doubt may amount from their conceptual similarity—

396 a rather depreciative stance towards an inner or outer event—a fact that would further illustrate
397 the intricate link between communicative intentions and prosody.

398 The valence and arousal ratings revealed distinctive affective properties of the different
399 speech act categories, which were consistent across words and non-words (Fig. 2). This suggests
400 that speaker’s intentions may have emotional connotations that listeners are able to detect in the
401 prosodic signal. A natural question that comes up is, how strongly the identification of intentions
402 in prosody *depends* on emotion recognition and whether intentions can be recognized without
403 taking emotion into account. We addressed this issue in Experiment 3.

404 **Experiment 3 – Behavior and Acoustics**

405 Experiment 1 and 2 demonstrated that prosodically coded intentions can be differentiated
406 (i) physically based on characteristic acoustic feature configurations as well as (ii) perceptually
407 in a 6-AFC task. What remains to be shown is, in how far the acoustic differences account for
408 participants’ ability to identify the speaker’s intention. If participants use the prosody’s acoustic
409 information for intention understanding, it should be possible to predict listeners’ perception
410 from the acoustic measures and, further, to identify the different feature combinations that evoke
411 specific speech act impressions. We addressed this question by feeding acoustic measures and
412 typicality ratings for every speech act into a multiple regression analysis. Moreover, to assess the
413 influence of emotion perception on intention recognition (see Experiment 2), we conducted an
414 additional regression analysis in which valence and arousal ratings were regressed out.

415 *Materials and Methods*

416 *Participants*

417 A new group of 20 healthy volunteers (10 females, mean age \pm SD: 24.8 \pm 4.1 years) for
418 the words and 20 participants for the non-words (10 females, mean age \pm SD: 24.6 \pm 3.2 years)
419 took part in a rating study. All participants reported normal hearing ability, gave written
420 informed consent and were paid 7€ per hour for their participation.

421 *Design and Procedure*

422 In this experiment, participants were asked to indicate to what extent each stimulus
423 sounded like a given speech act category (criticism, doubt, naming, suggestion, warning, or
424 wish). Compared to the 6-AFC categorization task, such speech act ratings provide a more
425 refined and less strategy dependent measure for the participants' perception and allowed for the
426 application of multiple regression analyses. In total, each stimulus was presented six times, once
427 for every speech act scale, in separate blocks. Responses were given with a slider on a visual
428 analogue scale from 0 to 100 ('intonation does not fit the intention at all' to 'intonation fits the
429 intention very well'). Each block comprised the same 192 stimuli—four tokens of two (non-
430)words expressed as six speech acts by four speakers—that were chosen from the full stimulus
431 pool of 384 stimuli. Stimuli and block order were again pseudo-randomized. The timing of the
432 experiment was self-paced and participants were able to take breaks between blocks. The results
433 of the ratings were analyzed by repeated-measures ANOVAs with Greenhouse-Geisser
434 correction for the factor SPEECH ACT for every speech act scale separately.

435 *Multiple regressions*

436 To elucidate which acoustic features guided the participants' ratings on the speech act
437 scales, we performed linear multiple regression analyses. Specifically, we used the acoustic
438 features as predictors (independent variable) for the subjective ratings of the 192 stimuli
439 (dependent variable), separately for words and non-words. Acoustic features were the same as in
440 Experiment 1 and were chosen such to include measures of duration, intensity, pitch, and
441 spectrum while keeping multicollinearity low (variance inflation factor words: < 3.526, non-
442 words: < 4.561).

443 Furthermore, to demonstrate that intention perception is not merely determined by
444 perceived emotional connotations in the speech signal, the emotion perception of the stimuli was
445 regressed out in two steps: Firstly, separate regressions were calculated with the valence and
446 arousal ratings as independent variables and the single speech act ratings as dependent variables.
447 This way, we bound all the variance in the perceived speech act that could be explained by
448 potentially perceived emotions. Thus, the residuals of these regressions should contain
449 information about the participants' intention perception devoid of the perceived valence and
450 arousal. Following this, new regressions were performed, now with the acoustic features as
451 independent and the standardized residuals of the speech act ratings as dependent variables.

452 *Results*

453 *Multiple regressions*

454 The mean ratings of the stimuli according to the six different speech act scales are shown
455 in Table 3. As can be seen in the diagonal, the highest ratings were obtained for the correct
456 speech act category. This was confirmed by a significant main effect of SPEECH ACT in

457 repeated-measures ANOVAs performed for every speech act scale separately (word stimuli: F 's
458 > 46.289 , p 's $< .001$; non-word stimuli: F 's > 29.326 , p 's $< .001$). Post-hoc paired comparisons
459 with Bonferroni correction showed that speech act stimuli were rated significantly higher on
460 their corresponding scale than any other speech act category with p 's $< .03$. Altogether, the
461 ratings replicate the findings in the 6-AFC categorization task, in that also this new group of
462 participants was well able to recognize and evaluate the speech acts correctly.

463 Table 3

464 To examine whether specific patterns of acoustic features can predict subjective
465 evaluation of the different speech act stimuli, the ratings together with the acoustic measures for
466 the single stimuli were entered into multiple regression analyses, separately for each speech act
467 rating scale. These regressions yielded highly significant results on all scales (see Table 4 for
468 detailed results). The variance explained by the regression models ranged from 11.6% for the
469 wish ratings to 52.7% for the warning ratings of the words and from 13.5% for the criticism
470 ratings to 62.1% for warning ratings of the non-words. The beta weights of the regression
471 functions indicate the degree to which the acoustic parameters predicted the ratings. Put
472 differently, high absolute values of the beta weights reflect the importance of the corresponding
473 acoustic feature for the prediction of the regression model. Almost all acoustic features
474 contributed significantly to the predictions of the speech act ratings (see Table 4). While the
475 spectral features (centre of gravity and standard deviation of the spectrum) as well as the HNR
476 yielded very low beta values in general (all < 0.2 , except for HNR in naming ratings in the non-
477 words), the acoustic measures of pitch, amplitude, and duration reached absolute beta values of
478 up to 0.524, suggesting that these parameters are key features for the comprehension of the
479 intentions.

480 Table 4

481 Figure 3

482 Figure 3 shows the beta weights of the main acoustic features and reveals specific
483 patterns of acoustic parameters for the prediction of the different speech act ratings: While high
484 ratings for criticism and doubt were mainly predicted for long stimuli with high mean pitch and a
485 rising pitch contour (positive beta weights for voiced frames, mean f0 and offset-onset f0), high
486 ratings for naming were associated with short and soft stimuli with low mean pitch and falling
487 pitch contour (negative beta weights for these measures). Suggestion ratings relied on short
488 stimuli with low mean pitch, but a strong pitch rise and high intensity. High ratings for warning
489 were predicted, if the stimuli had a high mean pitch and intensity with a negative f0 offset-onset
490 relation. Beta values in predictions for wish ratings showed the least clear pattern of acoustic
491 information. For these stimuli, a low mean pitch was the most informative parameter. In total,
492 the judgment of different speech act categories seems to be based on distinct acoustic patterns.
493 Crucially, predictions of criticism and doubt ratings were based on similar patterns that only
494 differed in the contribution of intensity. This is in line with the confusion of criticism and doubt
495 in the behavioral categorization and ratings (Tables 2 and 3). Overall, the pitch-related features
496 (mean f0 and offset-onset f0) had a high influence on the ratings in all speech act categories,
497 qualifying them as the most important acoustic features in these analyses. A subset of speech acts
498 was further influenced by amplitude and duration features, while spectral parameters only
499 seemed to play a minor role. Notably, the results of the multiple regressions were very similar for
500 the words and non-words, even though the analyzed data were not only based on a different
501 stimulus set, but also on the ratings of two independent groups of participants. Therefore, these

502 results validate the distinct acoustic patterns that shape the perception of different speech acts
503 and intentions in single-word utterances.

504 *Multiple regressions controlled for emotion*

505 The second regression approach was performed to control for all variance that could be
506 explained by the perception of valence and arousal. After an initial regression to predict speech
507 act ratings as dependent variables from emotion ratings as independent variables, we conducted a
508 second regression to explain residual speech act information by the acoustic measures. These
509 regressions were still highly significant on all speech act scales (p 's < .001, Table 4) and
510 explained variance in the range from 15.3% for criticism ratings to 30.5% for suggestion ratings
511 of the words, and from 9.2% for criticism ratings to 30.5% for doubt ratings of the non-words.
512 Compared to the original regressions, there was a noticeable decrease in explained variance for
513 naming and warning in words and non-words, which indicates that some of the variance could be
514 explained by the emotion perception of these speech acts. Prediction of the other speech acts was
515 virtually unchanged.

516 *Discussion*

517 The current experiment confirmed a link between the acoustics and perception of the
518 speech act stimuli by multiple regression analyses in which distinct acoustic feature
519 configurations significantly (but not fully) predicted the listeners' perception of the speech acts.
520 The amount of variance explained (ranging between 12% and 62% depending on speech act
521 type) was overall comparable to estimates found in previous studies on emotional prosody
522 (Banse & Scherer, 1996; Sauter, Eisner, Calder, et al., 2010a), validating our approach. In
523 general, pitch features (mean pitch and offset-onset f_0) were most influential for the perception

524 of different intentions. Further important cues could be derived from intensity and duration
525 measures, whereas spectral features contributed least to the intention predictions.

526 The amount of variance explained by the regressions was significant, but the values for
527 some speech acts (e.g. 12% for wish in words or 13% for criticism in non-words) suggest that the
528 acoustic features chosen for the analyses are not the only basis for intention perception. The
529 inclusion of additional acoustic features might further increase the precision of the regression
530 models. On the other hand, higher cognitive processes, such as social inference, may contribute
531 to the recognition of the communicative intention (see below; Wichmann, 2002, 2000; Szameitat
532 et al., 2010). Still, the fact that different acoustic patterns can explain the perception of different
533 speech acts, generally leads to the assertion that prosody carries information about the speaker's
534 intended meaning.

535 *Emotions*

536 Importantly, perception of intentions was not solely based on recognition of the speaker's
537 emotion as shown by the additional regression analyses, taking valence and arousal of stimuli
538 into account. As mentioned in the introduction, we do not exclude that communicative intentions
539 are partly based on the emotional state of the speaker. Indeed, listeners could classify the stimuli
540 in terms of valence and arousal (Fig. 2; Experiment 2). Nevertheless, regression analyses still
541 explained a significant amount of variance of the speech act ratings and most speech act
542 predictions were virtually unchanged after these affective components had been regressed out
543 (Table 4). Only warning and naming showed a considerable decrease in the prediction rate which
544 might be explained by their extreme positions on the arousal scale (Fig. 2). On the other hand,
545 participants might have first identified the speaker's intention and then assigned the
546 corresponding valence and arousal because they were asked to do so in the experiment

547 (Experiment 2). Overall, although emotional connotations may be important for the recognition
548 of some speech acts, our results give no reason to assume a systematic influence of emotions on
549 the recognition of communicative intentions.

550 *Ratings vs. 6-AFC*

551 Finally, it is of note that the ratings replicated the results of the 6-AFC task used in
552 Experiment 2. Importantly, ratings are a more sensitive measure than forced-choice
553 categorization tasks because they not only allow participants to reject predefined response
554 categories but also to flexibly adjust their responses on every (visual-analog) speech act scale.
555 The fact that both typicality ratings in Experiment 3 and 6-AFC judgments in Experiment 2
556 (conducted in separate participant groups) yielded very similar results demonstrates the
557 robustness of our findings.

558 **General Discussion**

559 Action-theoretic views of language (Austin, 1962; Bühler, 1934; Grice, 1957; Searle,
560 1969) propose that speakers' intentions are the main core and driver of interpersonal
561 communication. Yet, speakers rarely express their intentions literally in the propositional content
562 of an utterance, raising the question of how the speaker's meaning is transmitted from sender to
563 receiver. Here, we conceptualized intentions in terms of speech acts and provide evidence that
564 prosody serves as an extralinguistic channel to convey intentions non-verbally. Acoustically,
565 speakers used distinct prosodic feature configurations for different speech acts. Behaviorally,
566 listeners were well able to differentiate these intentions from voice tone alone, even when no
567 semantic meaning (non-words) or situational context was available (single words). Further, a
568 direct link between acoustics and perception was demonstrated, in that acoustic features reliably

569 (although not fully) accounted for the listeners' perception of the stimulus—even when the
570 emotional connotation of the stimuli was controlled for.

571 Notably, our results were consistent across all three experiments. For example, in all
572 measures from acoustics to perception, warning was classified with highest and criticism with
573 lowest accuracy. Moreover, in both the stimulus-based discriminant analyses (Experiment 1) and
574 the multiple regressions (Experiment 3), pitch rise, mean f0, as well as mean intensity and
575 duration were the most and spectral features the least relevant cues for correct speech act
576 recognition. This is consistent with a special role of pitch features observed in similar studies on
577 verbal emotion (Banse & Scherer, 1996) and attitude recognition (Blanc & Dominey, 2003).
578 Overall, the consistency of our results across experiments and participant groups lends strong
579 support for the relevance of prosody in conveying communicative intentions.

580 *Conventional Prosodic Expressions*

581 Our results invite the assumption that speakers' intents are expressed in conventionalized
582 prosodic forms. This view is supported (i) by the consistency of the prosodic patterns across four
583 independent speakers for each of the six speech acts, and (ii) by the robustness of listeners'
584 performance in identifying the expressed intentions, despite absence of contextual or semantic
585 information. Arguably, prosodic patterns do not refer to communicative intentions as
586 unambiguously as words refer to objects in the world. Rather we propose that they represent
587 “communicative complexes” that connote a set of conceptually related pragmatic categories
588 (e.g., speech acts), whose distributions of relevant acoustic cues partly overlap. This acoustic and
589 conceptual overlap may account for the confusion of criticism and doubt in our experiments and
590 predicts a rather loose labeling of speakers' intentions in open choice tasks, licensing our use of

591 forced-choice task and typicality ratings (see below). Notably, our data suggest that the acoustic
592 characteristics of these “complexes” are conventionalized to the extent that listeners can infer the
593 relevant communicative concept by matching the perceived prosodic pattern with an internalized
594 probabilistic distribution of acoustic cue configurations for different intentions.

595 Such a direct recognition of speakers’ intent from prosody is reminiscent of previous
596 work on written speech acts (Holtgraves, 2008a) suggesting that the default interpretation of
597 illocutionary force can be based on generalized rather than particularized implicatures, i.e. can be
598 directly understood without contextual information, similar to most idioms (e.g., to call it a day)
599 or metaphors (e.g., He is a walking dictionary) (Glucksberg, Gildea, & Bookin, 1982;
600 Glucksberg, 2003; Keysar, 1989). The relevance of context for the classification of speakers’
601 prosodic intentions or attitudes has been a matter of debate for a while (Cutler, 1976; Wichmann,
602 2000, 2002). Some accounts posit that prosody mainly acts in a contrastive way, without
603 conveying meaning by itself (e.g. Attardo, Eisterhold, Hay, & Poggi, 2003; Bryant & Fox Tree,
604 2005). By deviating from its “default”, prosody is thought to motivate the listener to look for
605 “unspoken” meanings in the utterance, i.e. to infer implicit speech actions from literal meaning
606 by taking context information into account (Cutler & Isard, 1980; Levinson, 2013). However, as
607 Wichmann (2002) rightly pointed out, this view requires knowledge about prosodic “defaults”.
608 We argue that this knowledge is best characterized as experience-dependent inventory of
609 situationally distinct acoustic patterns that allows listeners to recognize broad communicative
610 concepts based on prosody. Such a distinguished role for prosody in intention transmission is
611 supported by the fact that these communicative concepts could be conveyed despite absence of
612 contextual information and without knowledge of the lexical content in non-words (Experiments
613 2 and 3).

614 Note that we do not claim that context plays no role at all. Very much like lexical and
615 syntactic processing is not based on acoustics alone but varies with context (for example in case
616 of homophones such as “meet” vs. “meat” or ambiguous word category as in “report”; for
617 review, see Piantadosi, Tily, & Gibson, 2012), also the prosodic recognition of speakers’ intents
618 can be shaped by context (Tanenhaus, Kurumada, & Brown, 2015). First, context predicts what
619 interpretations are likely and may, thus, resolve perceptual ambiguity between overlapping
620 distributions within the “communicative complex”, e.g., allowing listeners to better discriminate
621 between doubt and criticism. Second, context provides a sample of the speaker’s prosodic “style”
622 that allows listeners to flexibly adapt (even reverse) their prosodic interpretations accordingly
623 (Tanenhaus et al., 2015). Altogether, we conclude that (paralinguistic) prosody is a signal that is
624 able to convey a broad communicative concept on its own but becomes cognitively interlinked
625 and specified with complementary contextual information, if available.

626 *Prosody’s Initial Relevance for Social Communication*

627 Overall, the transfer of intentions via prosody might be a capability that forms the initial,
628 non-linguistic foundation of interpersonal communication (Bates, Camaioni, & Volterra, 1975;
629 Dore, 1975) that becomes gradually complemented and refined—yet not erased—by growing
630 verbal capacities, over the course of ontogeny and perhaps even phylogeny (Oller & Griebel,
631 2014). For example, primate calls have been found to signal the producer’s interactive stance
632 intentionally (Schel, Townsend, Machanda, Zuberbühler, & Slocombe, 2013) via distinct
633 acoustic structures (Crockford & Boesch, 2003; Seyfarth & Cheney, 2014), even if they lack
634 lexical (referential) meaning (Wheeler & Fischer, 2012). Developmentally, young infants start to
635 produce acoustically distinct prosodic patterns in the middle of their first year of life that are

636 initially used in communicative as opposed to self-centred emotional or exploratory contexts
637 (Papaeliou et al., 2002; Papaeliou & Trevarthen, 2006), later express specific “primitive intents”
638 (Esteve-Gibert & Prieto, 2013; Prieto et al., 2012) and endow pointing gestures with
639 communicative goals (Grünloh & Liszkowski, 2015). Notably, interactive prosodic patterns
640 emerge earlier than verbal skills and become meaningful communicative instruments, most likely
641 because parents differentiate their responses based on the acoustics of the child’s vocalizations
642 (cf. Lester et al., 1995; Oller & Griebel, 2014). The present data show that prosody continues to
643 be indicative of speakers’ intents in adulthood, despite mature verbal skills. More than that, the
644 data suggest that the use of prosodic cues evolves further beyond infancy to express more
645 complex intentions than those infants would ever produce (e.g., criticism or doubt). Whether
646 speakers resort more strongly to these (early) prosodic building blocks of communication when
647 verbal capacities may get lost or are nonexistent as in conditions of non-fluent aphasia (Barrett,
648 Crucian, Raymer, & Heilman, 1999; Warren, Warren, Fox, & Warrington, 2003) or foreign
649 languages is an interesting topic for future research.

650 *Prosody in Natural Language Use*

651 Single-word utterances are part of our everyday life and humans start to use prosody to
652 code for different pragmatic intentions in single words in early infancy (Dore, 1975; Prieto et al.,
653 2012). Yet, compared to longer sentences with additional semantic information, the brevity of
654 the present context-free stimuli may have led speakers to emphasize the relevant prosodic
655 features. Listeners, in turn, may be more used to decode intentions in sentential contexts that
656 often resolve ambiguities (even if ambiguities were mitigated by the 6-AFC task and typicality
657 ratings in the present study). Future studies can help to generalize our results by using a wider set

658 of recordings (as suggested by Banse and Scherer, 1996), for example, including sentence-level
659 stimuli, more variable tokens (i.e. more words/sentences), more speech acts, and more speakers.

660 Apart from that, another point of discussion is in how far prosodies produced in the sound
661 lab using fictional scenarios correspond to prosodies produced in natural conversations.
662 Although a direct empirical investigation is still pending, there are several reasons that grant the
663 ecological validity of our sound stimuli. First, cues for expressing intentions are typically
664 produced voluntarily during an interaction. Therefore, they have a posed character by nature and
665 may not suffer from artificial recording situations to the extent as emotions do (Jürgens, Grass,
666 Drolet, & Fischer, 2015; Jürgens et al., 2011). Second, our speakers—although trained in
667 producing clear and artifact-free speech—were non-actors. Hence, they relied on their everyday
668 speech experience to express the intention in a way they would naturally do to be understood by
669 an interlocutor. Last, studies on non-prosodic cues for speech acts (Bucciarelli et al., 2003;
670 Reeder, 1980) and voluntary vocal expressions of social affect (Rilliard, Shochi, Martin,
671 Erickson, & Auberge, 2009) suggest that cues for expressing intentions are not innate but
672 culturally learned. On this assumption, the fact that our speakers and listeners used and
673 understood the specific prosodic cues suggests that these cues must occur in natural
674 conversations.

675 *Future Research on Intentional Prosody*

676 An interesting question with regard to speaker's intent in natural communication is, then,
677 how prosodic cues are weighed and cognitively interlinked with other paralinguistic cues such as
678 facial expressions. Notably, the latter have been shown to serve explicit interpersonal functions
679 that reach beyond the inadvertent display of basic emotions (Ekman, 1992), for example when

680 (voluntarily) communicating *social motives* (e.g., in case of compassion or empathy for pain)
681 (Fridlund, 1994; Parkinson, 2005). Concerning audio-visual integration, recent motion-capture
682 and neuroimaging studies revealed interactions between linguistic/emotional prosody and facial
683 expressions, in speakers (Cvejic, Kim, & Davis, 2012; Kitamura, Guellaï, & Kim, 2014) and in
684 listeners, respectively (Brück, Kreifelts, & Wildgruber, 2011; Watson et al., 2014). Yet, whether
685 and how prosody and facial cues are fused in the transfer of speaker's meaning is currently not
686 known and an interesting topic for future research.

687 Another point that deserves further examination is our observation that acoustic
688 information predicted participant's speech act recognition successfully, yet not fully. This raises
689 the interesting hypothesis that the comprehension of speaker's intentions from prosody relies on
690 a weighted contribution of auditory-prosodic and other, socio-cognitive processes whose exact
691 nature and ways of interaction still need to be clarified. On the socio-cognitive side, recent
692 neuroimaging work lends initial evidence for inferential processes, i.e. involving theory of mind
693 areas, during the comprehension of speech acts (Egorova et al., 2014; Egorova, Shtyrov, &
694 Pulvermüller, 2015) and speaker meaning (Bašnáková, Weber, Petersson, van Berkum, &
695 Hagoort, 2014; Jang et al., 2013), as well as motor system involvement during the processing of
696 directive speech acts (Egorova et al., 2014) and indirect requests (van Ackeren, Casasanto,
697 Bekkering, Hagoort, & Rueschemeyer, 2012). Yet, none of these studies involved prosody,
698 leaving the fundamental question unresolved how prosody potentially interlinks with these socio-
699 cognitive systems. Future neurocognitive investigations with the present stimuli may help to
700 elucidate this question and are currently underway.

701 **Conclusion**

702 Speakers rarely code their intentions in the lexical content of an utterance. Yet, listeners
703 easily recognize the speaker’s communicative goals. The present study shows that
704 conversationalists are able to use prosody as extralinguistic cue to specify communicative
705 intentions—an early capacity that complements adults’ mature verbal abilities. Interlocutors
706 produce and understand prosodic cues independently of the semantic meaning, contextual
707 information, and emotional coloring of the utterance. These results argue in favor of
708 conventionalized acoustic feature configurations that connote communicative concepts, although
709 their acoustic and conceptual distributions may partly overlap. The present study leads towards
710 future research on the interaction between auditory-prosodic cues, conversation context, and
711 socio-cognitive processes serving the transfer of speaker meaning as the foundation of successful
712 interpersonal communication.

713 **Acknowledgments**

714 We thank Prof. Dr. J. D. Jescheniak and Prof. Dr. K. von Kriegstein for their helpful
715 feedback on this project. The study was funded by the Otto Hahn award of the Max Planck
716 Society to DS.

717

718 **References**

- 719 Attardo, S., Eisterhold, J., Hay, J., & Poggi, I. (2003). Multimodal markers of irony and sarcasm.
720 *Humor - International Journal of Humor Research*, 16(2), 243–260.
- 721 Austin, J. L. (1962). *How to do things with words*. Cambridge: Harvard University Press.
- 722 Bach, K. (1994). Conversational implicature. *Mind & Language*, 9(2), 124–162.
- 723 Banse, R., & Scherer, K. R. (1996). Acoustic profiles in vocal emotion expression. *Journal of*
724 *Personality and Social Psychology*, 70(3), 614–36.
- 725 Bänziger, T., & Scherer, K. R. (2005). The role of intonation in emotional expressions. *Speech*
726 *Communication*, 46(3-4), 252–267.
- 727 Barrett, M. D., Crucian, G. P., Raymer, A. M., & Heilman, K. M. (1999). Spared Comprehension
728 of Emotional Prosody in a Patient with Global Aphasia. *Neuropsychiatry, Neuropsychology,*
729 *and Behavioral Neurology*, 12(2), 117–120.
- 730 Bašnáková, J., Weber, K., Petersson, K. M., van Berkum, J., & Hagoort, P. (2014). Beyond the
731 Language Given: The Neural Correlates of Inferring Speaker Meaning. *Cerebral Cortex*,
732 24(10), 2572–2578.
- 733 Bates, E., Camaioni, L., & Volterra, V. (1975). The acquisition of performatives prior to speech.
734 *Merrill-Palmer Quarterly*, 21(3), 205–226.
- 735 Blanc, J., & Dominey, P. (2003). Identification of prosodic attitudes by a temporal recurrent
736 network. *Cognitive Brain Research*, 17, 693–699.
- 737 Boersma, P., & Weenink, D. (2014). Praat: doing phonetics by computer [Computer program].
738 Version 5.3.80, retrieved 29 June 2014 from <http://www.praat.org/>.
- 739 Bolinger, D. (1986). *Intonation and Its Parts: Melody in Spoken English*. Stanford: Stanford
740 University Press.
- 741 Bradley, M., & Lang, P. J. (1994). Measuring Emotion: The Self-Assessment Manikin and the
742 Semantic Differential. *Journal of Behavior Therapy and Experimental Psychiatry*, 25(1),
743 49–59.
- 744 Brück, C., Kreifelts, B., & Wildgruber, D. (2011). Emotional voices in context: a neurobiological
745 model of multimodal affective information processing. *Physics of Life Reviews*, 8(4), 383–
746 403.
- 747 Bryant, G. a., & Fox Tree, J. E. (2005). Is there an ironic tone of voice? *Language and Speech*,
748 48(Pt 3), 257–277.
- 749 Bucciarelli, M., Colle, L., & Bara, B. G. (2003). How children comprehend speech acts and
750 communicative gestures. *Journal of Pragmatics*, 35(2), 207–241.
- 751 Bühler, K. (1934). *Sprachtheorie: Die Darstellungsfunktion der Sprache*. Jena: Gustav Fischer.
- 752 Carlson, K., Frazier, L., & Clifton, C. J. (2009). How prosody constrains comprehension: A
753 limited effect of prosodic packaging. *Lingua*, 119(7), 1066–1082.

- 754 doi:10.1016/j.biotechadv.2011.08.021.Secreted
- 755 Clark, H. H., & Carlson, T. B. (1981). Context for Comprehension. In J. Ling & A. Baddeley
756 (Eds.), *Attention and Performance IX* (pp. 313–330). Hillsdale, New Jersey: Lawrence
757 Erlbaum Associates.
- 758 Crockford, C., & Boesch, C. (2003). Context-specific calls in wild chimpanzees, *Pan troglodytes*
759 *verus*: analysis of barks. *Animal Behaviour*, *66*(1), 115–125.
- 760 Cutler, A. (1976). The Context-Dependence of “Intonational Meanings,” (2), 104–115.
- 761 Cutler, A., Dahan, D., & van Donselaar, W. (1997). Prosody in the comprehension of spoken
762 language: a literature review. *Language and Speech*, *40*(2), 141–201.
- 763 Cutler, A., & Isard, S. D. (1980). The Production of Prosody. In B. Butterworth (Ed.), *Language*
764 *Production: Speech and Talk* (Volume 1., pp. 245–270). Academic Press.
- 765 Cvejic, E., Kim, J., & Davis, C. (2012). Recognizing prosody across modalities, face areas and
766 speakers: Examining perceivers’ sensitivity to variable realizations of visual prosody.
767 *Cognition*, *122*(3), 442–453.
- 768 Di Cesare, G., Di Dio, C., Marchi, M., & Rizzolatti, G. (2015). Expressing our internal states and
769 understanding those of others. *Proceedings of the National Academy of Sciences*, *112*(33),
770 10331–10335.
- 771 Dore, J. (1975). Holophrases, speech acts and language universals. *Journal of Child Language*,
772 *2*(01), 21–40.
- 773 Egorova, N., Pulvermüller, F., & Shtyrov, Y. (2014). Neural dynamics of speech act
774 comprehension: an MEG study of naming and requesting. *Brain Topography*, *27*(3), 375–
775 92.
- 776 Egorova, N., Shtyrov, Y., & Pulvermüller, F. (2013). Early and parallel processing of pragmatic
777 and semantic information in speech acts: neurophysiological evidence. *Frontiers in Human*
778 *Neuroscience*, *7*(86), 1–13.
- 779 Egorova, N., Shtyrov, Y., & Pulvermüller, F. (2015). Brain basis of communicative actions in
780 language. *NeuroImage*, *125*, 857–867.
- 781 Ekman, P. (1992). An argument for basic emotions. *Cognition & Emotion*.
- 782 Enrici, I., Adenzato, M., Cappa, S., Bara, B. G., & Tettamanti, M. (2011). Intention processing in
783 communication: a common brain network for language and gestures. *Journal of Cognitive*
784 *Neuroscience*, *23*(9), 2415–2431.
- 785 Esteve-Gibert, N., & Prieto, P. (2013). Prosody signals the emergence of intentional
786 communication in the first year of life: evidence from Catalan-babbling infants. *Journal of*
787 *Child Language*.
- 788 Frick, R. W. (1985). Communicating emotion: The role of prosodic features. *Psychological*
789 *Bulletin*, *97*(3), 412–429.
- 790 Fridlund, A. J. (1994). *Human Facial Expression: An Evolutionary View*. Academic Press.
- 791 Frith, C. (2009). Role of facial expressions in social interactions. *Philosophical Transactions of*

- 792 *the Royal Society of London. Series B, Biological Sciences*, 364(1535), 3453–3458.
- 793 Furrow, D., Podrouzek, W., & Moore, C. (1990). The acoustical analysis of children’s use of
794 prosody in assertive and directive contexts. *First Language*, (10), 37–49.
- 795 Gisladdottir, R. S., Chwilla, D. J., & Levinson, S. C. (2015). Conversation Electrified: ERP
796 Correlates of Speech Act Recognition in Underspecified Utterances. *PloS ONE*, 10(3), 1–
797 24.
- 798 Glucksberg, S. (2003). The psycholinguistics of metaphor. *Trends in Cognitive Sciences*, 7(2),
799 92–96.
- 800 Glucksberg, S., Gildea, P., & Bookin, H. B. (1982). On Understanding Nonliteral Speech: Can
801 People Ignore Metaphors? *Journal of Verbal Learning and Verbal Behavior*, 21, 85–98.
- 802 Grice, H. P. (1957). Meaning. *The Philosophical Review*, 66(3), 377–388.
- 803 Grice, H. P. (1975). Logic and Conversation. In P. Cole & J. Morgan (Eds.), *Syntax and*
804 *Semantics: 3. Speech acts* (pp. 41–58). New York: New York: Academic Press.
- 805 Grünloh, T., & Liszkowski, U. (2015). Prelinguistic Vocalizations Distinguish Pointing Acts.
806 *Journal of Child Language*, 49(0), 1–48.
- 807 Holtgraves, T. (2002). *Language as Social Action*. London: Lawrence Erlbaum Associates.
- 808 Holtgraves, T. (2005). The production and perception of implicit performatives. *Journal of*
809 *Pragmatics*, 37(12), 2024–2043.
- 810 Holtgraves, T. (2008a). Automatic intention recognition in conversation processing. *Journal of*
811 *Memory and Language*, 58(3), 627–645.
- 812 Holtgraves, T. (2008b). Conversation, speech acts, and memory. *Memory & Cognition*, 36(2),
813 361–374.
- 814 Jang, G., Yoon, S. A., Lee, S. E., Park, H., Kim, J., Ko, J. H., & Park, H. J. (2013). Everyday
815 conversation requires cognitive inference: Neural bases of comprehending implicated
816 meanings in conversations. *NeuroImage*, 81, 61–72.
- 817 Jiang, X., & Pell, M. D. (2015). On how the brain decodes vocal cues about speaker confidence.
818 *Cortex*, 66, 9–34.
- 819 Jürgens, R., Grass, A., Drolet, M., & Fischer, J. (2015). Effect of Acting Experience on Emotion
820 Expression and Recognition in Voice: Non-Actors Provide Better Stimuli than Expected.
821 *Journal of Nonverbal Behavior*, 195–214.
- 822 Jürgens, R., Hammerschmidt, K., & Fischer, J. (2011). Authentic and play-acted vocal emotion
823 expressions reveal acoustic differences. *Frontiers in Psychology*, 2(JUL), 1–11.
- 824 Keysar, B. (1989). On the Functional Equivalence of Literal and Metaphorical Interpretations in
825 Discourse. *Journal of Memory and Language*, 28, 375–385.
- 826 Kitamura, C., Guellai, B., & Kim, J. (2014). Motherese by Eye and Ear: Infants Perceive Visual
827 Prosody in Point-Line Displays of Talking Heads. *PLoS ONE*, 9(10), e111467.
- 828 Lester, B. M., Boukydis, C. F. Zachariah Garcia-Coll, C. T., Peucker, M., McGrath, M. M.,

- 829 Vohr, B. R., Brem, F., & Oh, W. (1995). Developmental Outcome as a Function of the
830 Goodness of Fit Between the Infant's Cry Characteristics and the Mother's Perception of
831 Her Infant's Cry. *Pediatrics*, 95(4), 516–521.
- 832 Levinson, S. C. (2006). On the “Human Interaction Engine.” In N. J. Enfield & S. C. Levinson
833 (Eds.), *Roots of Human Sociality: Culture, Cognition, and Interaction* (pp. 39–69). Oxford:
834 Berg.
- 835 Levinson, S. C. (2013). Action Formation and Ascription. In T. Stivers & J. Sidnell (Eds.), *The*
836 *Handbook of Conversation Analysis* (pp. 103–130). Wiley-Blackwell.
- 837 Liu, S. (2011). An experimental study of the classification and recognition of Chinese speech
838 acts. *Journal of Pragmatics*, 43(6), 1801–1817.
- 839 Marcos, H. (1987). Communicative functions of pitch range and pitch direction in infants.
840 *Journal of Child Language*, 14(02), 255.
- 841 Mead, G. H. (1934). *Mind, Self, and Society: From the Standpoint of a Social Behaviorist*. (C.
842 W. Morris, Ed.). University of Chicago Press.
- 843 Monetta, L., Cheang, H. S., & Pell, M. D. (2008). Understanding speaker attitudes from prosody
844 by adults with Parkinson's disease. *Journal of Neuropsychology*, 2(2), 415–430.
- 845 Morlec, Y., Bailly, G., & Aubergé, V. (2001). Generating prosodic attitudes in French: Data,
846 model and evaluation. *Speech Communication*, 33(4), 357–371.
- 847 Oller, D. K., & Griebel, U. (2014). On Quantitative Comparative Research in Communication
848 and Language Evolution. *Biological Theory*, 9, 296–308.
- 849 Papaeliou, C., Minadakis, G., & Cavouras, D. (2002). Acoustic Patterns of Infant Vocalizations
850 Expressing Emotions and Communicative Functions. *Journal of Speech, Language, and*
851 *Hearing Research*, 45, 311–317.
- 852 Papaeliou, C., & Trevarthen, C. (2006). Prelinguistic pitch patterns expressing “communication”
853 and “apprehension”. *Journal of Child Language*, 33(1), 163–178.
- 854 Parkinson, B. (2005). Do facial movements express emotions or communicate motives?
855 *Personality and Social Psychology Review : An Official Journal of the Society for*
856 *Personality and Social Psychology, Inc*, 9(4), 278–311.
- 857 Piantadosi, S. T., Tily, H., & Gibson, E. (2012). The communicative function of ambiguity in
858 language. *Cognition*, 122(3), 280–291.
- 859 Prieto, P., Estrella, A., Thorson, J., & Vanrell, M. D. M. (2012). Is prosodic development
860 correlated with grammatical and lexical development? Evidence from emerging intonation
861 in Catalan and Spanish. *Journal of Child Language*, 39(02), 221–257.
- 862 Reeder, K. (1980). The emergence of illocutionary skills. *Journal of Child Language*, 7(1), 13–
863 28.
- 864 Remington, N. a, Fabrigar, L. R., & Visser, P. S. (2000). Reexamining the circumplex model
865 of affect. *Journal of Personality and Social Psychology*, 79(2), 286–300.
- 866 Rigoulot, S., Fish, K., & Pell, M. D. (2014). Neural correlates of inferring speaker sincerity from

867 white lies: an event-related potential source localization study. *Brain Research*, 1565, 48–
868 62.

869 Rilliard, a., Shochi, T., Martin, J.-C., Erickson, D., & Auberge, V. (2009). Multimodal Indices
870 to Japanese and French Prosodically Expressed Social Affects. *Language and Speech*, 52(2-
871 3), 223–243.

872 Russell, J. a. (1980). A circumplex model of affect. *Journal of Personality & Social Psychology*.

873 Sammler, D., Grosbras, M.-H., Anwander, A., Bestelmeyer, P. E. G., & Belin, P. (2015). Dorsal
874 and Ventral Pathways for Prosody. *Current Biology*, 1–7.

875 Sauter, D. A., Eisner, F., Calder, A. J., & Scott, S. K. (2010a). Perceptual cues in nonverbal
876 vocal expressions of emotion. *Quarterly Journal of Experimental Psychology (2006)*,
877 63(11), 2251–72.

878 Sauter, D. A., Eisner, F., Calder, A. J., & Scott, S. K. (2010b). Perceptual cues in nonverbal
879 vocal expressions of emotion. *Quarterly Journal of Experimental Psychology*, 63(11),
880 2251–72.

881 Sauter, D. A., Eisner, F., Ekman, P., & Scott, S. K. (2010). Cross-cultural recognition of basic
882 emotions through nonverbal emotional vocalizations. *Proceedings of the National Academy
883 of Sciences of the United States of America*, 107(6), 2408–12.

884 Schel, A. M., Townsend, S. W., Machanda, Z., Zuberbühler, K., & Slocombe, K. E. (2013).
885 Chimpanzee Alarm Call Production Meets Key Criteria for Intentionality. *PLoS ONE*,
886 8(10).

887 Schneider, K., Lintfert, B., Dogil, G., & Möbius, B. (2006). Phonetic Grounding of Prosodic
888 Categories. In *Methods in empirical prosody research* (pp. 335–362). Berlin: De Gruyter.

889 Searle, J. R. (1969). *Speech acts: An essay in the philosophy of language*. Cambridge University
890 Press.

891 Searle, J. R., & Vanderveken, D. (1985). *Foundation of Illocutionary Logic*. Cambridge
892 University Press.

893 Seyfarth, R. M., & Cheney, D. L. (2014). The evolution of language from social cognition.
894 *Current Opinion in Neurobiology*, 28, 5–9.

895 Simon-Thomas, E. R., Keltner, D. J., Sauter, D., Sinicropi-Yao, L., & Abramson, A. (2009). The
896 voice conveys specific emotions: evidence from vocal burst displays. *Emotion (Washington,
897 D.C.)*, 9(6), 838–846.

898 Srinivasan, R. J., & Massaro, D. W. (2003). Perceiving prosody from the face and voice:
899 distinguishing statements from echoic questions in English. *Language and Speech*, 46(Pt 1),
900 1–22.

901 Stalnaker, R. (2002). Common Ground. *Linguistics and Philosophy*, 25, 701–721.

902 Szameitat, D. P., Alter, K., Szameitat, A. J., Darwin, C. J., Wildgruber, D., Dietrich, S., & Sterr,
903 A. (2009). Differentiation of emotions in laughter at the behavioral level. *Emotion
904 (Washington, D.C.)*, 9(3), 397–405.

905 Szameitat, D. P., Alter, K., Szameitat, A. J., Wildgruber, D., Sterr, A., & Darwin, C. J. (2009).
906 Acoustic profiles of distinct emotional expressions in laughter. *The Journal of the*
907 *Acoustical Society of America*, 126(1), 354–66.

908 Szameitat, D. P., Kreifelts, B., Alter, K., Szameitat, A. J., Sterr, A., Grodd, W., & Wildgruber, D.
909 (2010). It is not always tickling: distinct cerebral responses during perception of different
910 laughter types. *NeuroImage*, 53(4), 1264–71.

911 Tanenhaus, M. K., Kurumada, C., & Brown, M. (2015). Prosody and Intention Recognition. In
912 L. Frazier & E. Gibson (Eds.), *Explicit and Implicit Prosody in Sentence Processing* (pp.
913 99–118). Springer.

914 Tomasello, M. (2005). Intention Reading and Imitative Learning. In S. Hurley & N. Chater
915 (Eds.), *Perspectives on Imitation: From Neuroscience to Social Science* (pp. 133–148). MA:
916 MIT Press.

917 Uldall, E. (1960). Attitudinal meanings conveyed by intonation contours. *Language and Speech*,
918 3, 223–234.

919 van Ackeren, M. J., Casasanto, D., Bekkering, H., Hagoort, P., & Rueschemeyer, S.-A. (2012).
920 Pragmatics in Action: Indirect Requests Engage Theory of Mind Areas and the Cortical
921 Motor Network. *Journal of Cognitive Neuroscience*, 24(11), 2237–2247.

922 Wagner, M., & Watson, D. G. (2010). Experimental and theoretical advances in prosody: A
923 review. *Language and Cognitive Processes*, 25(7-9), 905–945.

924 Warren, J. D., Warren, J. E., Fox, N. C., & Warrington, E. K. (2003). Nothing to say, something
925 to sing: primary progressive dynamic aphasia. *Neurocase*, 9(2), 140–55.

926 Warren, P. (1999). Prosody and language processing. In *Language Processing* (pp. 155–188).
927 Psychology Press Ltd.

928 Watson, R., Latinus, M., Noguchi, T., Garrod, O., Crabbe, F., & Belin, P. (2014). Crossmodal
929 Adaptation in Right Posterior Superior Temporal Sulcus during Face-Voice Emotional
930 Integration. *Journal of Neuroscience*, 34(20), 6813–6821.

931 Wheeler, B. C., & Fischer, J. (2012). Functionally referential signals: A promising paradigm
932 whose time has passed. *Evolutionary Anthropology: Issues, News, and Reviews*, 21(5), 195–
933 205.

934 Wichmann, A. (2000). The Attitudinal Effects Of Prosody, And How They Relate To Emotion.
935 In *ITRW on Speech and Emotion*.

936 Wichmann, A. (2002). Attitudinal intonation and the inferential process. In *Speech Prosody*
937 *2002, International Conference*.

938 Wilson, D., & Sperber, D. (2012). *Meaning and Relevance*. Cambridge University Press.

939 Wittgenstein, L. (1953). *Philosophische Untersuchungen*. Suhrkamp.

940 Wundt, W. (1896). *Grundriss der Psychologie*. Leipzig: Engelmann.

941

942

943 **TABLES**

944 **Table 1.** Results of cross-validated / jackknife discriminant analysis for classification of speech
 945 acts from the acoustic features (in %).

Stimulus type	Classification					
	Criticism	Doubt	Naming	Suggestion	Warning	Wish
<i>Words</i>						
Criticism	76.6	4.7	0	10.9	7.8	0
Doubt	1.6	92.2	0	1.6	0	4.7
Naming	0	0	100	0	0	0
Suggestion	3.1	1.6	0	95.3	0	0
Warning	0	0	0	0	100	0
Wish	0	0	9.4	0	0	90.6
<i>Non-Words</i>						
Criticism	79.7	6.2	0	7.8	4.7	1.6
Doubt	7.8	84.4	3.1	0	0	4.7
Naming	0	0	100	0	0	0
Suggestion	1.6	0	0	98.4	0	0
Warning	0	0	0	0	100	0
Wish	0	0	4.7	0	0	95.3

946

947

948 **Table 2.** Behavioral categorization of speech acts (in %).

Stimulus type	Participants' responses					
	Criticism	Doubt	Naming	Suggestion	Warning	Wish
<i>Words</i>						
Criticism	62.0	24.0	0.2	4.9	6.9	2.0
Doubt	5.3	83.4	3.1	4.4	0.6	3.3
Naming	1.7	0.3	90.0	1.3	0.8	5.9
Suggestion	4.5	9.2	3.3	80.3	0.5	2.2
Warning	4.1	0	0.3	0.5	89.5	5.6
Wish	1.9	0.8	8.3	4.4	1.1	83.6
<i>Non-Words</i>						
Criticism	52.4	29.7	1.1	3.9	9.4	3.4
Doubt	2.7	82.6	2.2	8.5	0	4.1
Naming	17.7	1.9	74.1	2.3	0.3	3.8
Suggestion	9.7	12.8	3.3	63.7	0.5	10.0
Warning	3.8	0.5	0.0	0.2	94.2	1.4
Wish	5.3	1.6	8.9	14.2	0.6	69.4

949

950 **Table 3.** Participants' ratings of speech acts (min = 0, max = 100).

Stimulus type	Speech act scale					
	Criticism	Doubt	Naming	Suggestion	Warning	Wish
<i>Words</i>						
Criticism	70.0	56.2	9.4	18.5	27.5	16.5
Doubt	45.7	82.5	13.7	16.7	10.7	9.5
Naming	12.2	13.2	85.3	10.5	9.3	17.4
Suggestion	15.7	28.0	17.7	77.3	7.7	26.9
Warning	22.8	18.6	8.2	16.0	86.0	29.9
Wish	6.5	10.1	21.3	16.3	6.0	80.4
<i>Non-Words</i>						
Criticism	63.3	45.5	11.3	21.3	31.3	13.3
Doubt	23.0	77.5	17.2	22.7	6.1	20.6
Naming	14.9	13.1	78.0	17.7	8.7	21.7
Suggestion	20.5	35.3	24.9	70.6	10.9	18.3
Warning	26.0	5.4	7.5	8.5	94.1	9.1
Wish	10.1	11.4	25.7	22.8	5.4	75.9

951

952

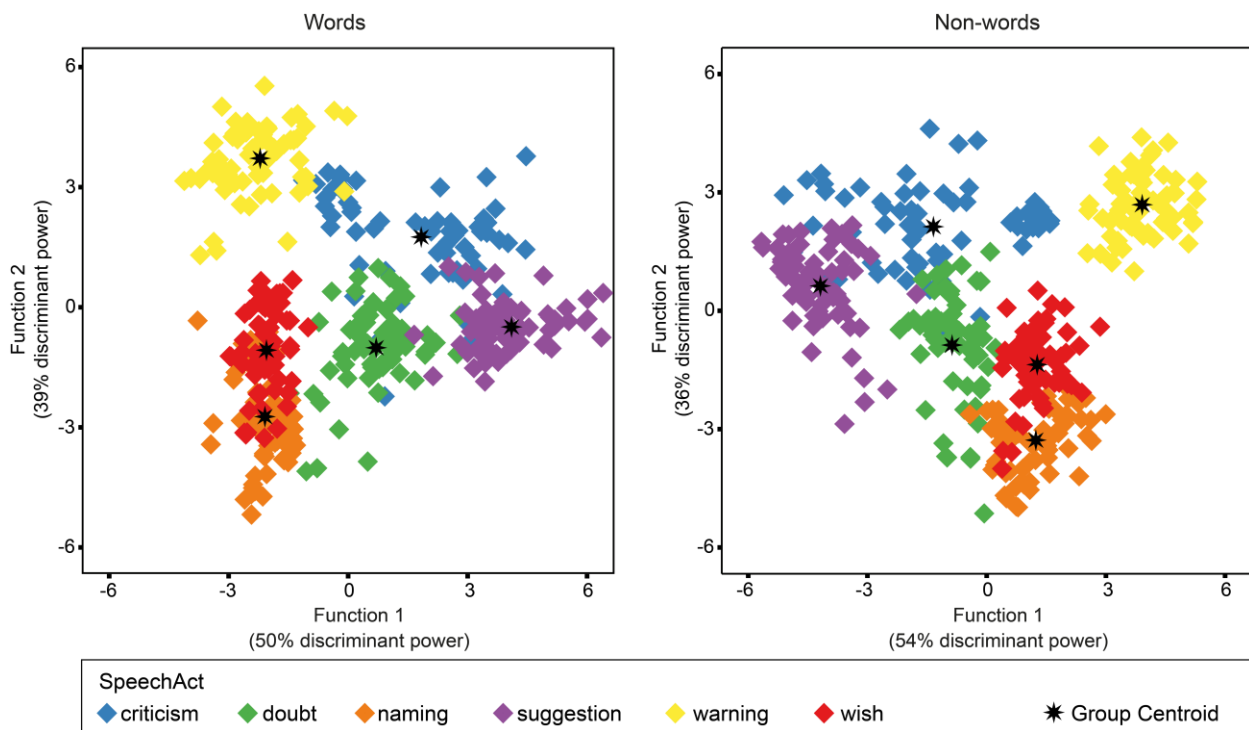
953 **Table 4.** Multiple regression analyses of acoustic features and speech act ratings (beta-weights).

Acoustic parameter	Speech act ratings											
	Criticism		Doubt		Naming		Suggestion		Warning		Wish	
<i>Words</i>												
Voiced Frames	0.276	***	0.348	***	-0.420	***	-0.227	***	0.073	***	0.083	***
Mean F0	0.371	***	0.402	***	-0.326	***	-0.211	***	0.385	***	-0.294	***
Offset-onset F0	0.224	***	0.315	***	-0.277	***	0.524	***	-0.387	***	-0.135	***
Mean Intensity	-0.015		-0.250	***	-0.266	***	0.368	***	0.343	***	0.118	***
Mean HNR	-0.171	***	-0.106	***	0.141	***	-0.008		-0.127	***	0.044	
Centre of Gravity	-0.083	**	-0.106	***	0.097	***	-0.059	*	-0.039	*	0.042	
SD Spectrum	-0.085	***	-0.082	***	0.014		0.035	*	-0.107	***	0.149	***
Adj R²	0.171	***	0.273	***	0.364	***	0.349	***	0.527	***	0.116	***
Adj R² (emo-corr)	0.153	***	0.285	***	0.228	***	0.305	***	0.162	***	0.175	***
<i>Non-Words</i>												
Voiced Frames	0.183	***	0.237	***	-0.388	***	-0.184	***	0.056	***	0.141	***
Mean F0	0.121	***	0.331	***	-0.376	***	-0.258	***	0.289	***	-0.204	***
Offset-onset F0	0.259	***	0.416	***	-0.206	***	0.407	***	-0.387	***	-0.062	**
Mean Intensity	0.168	***	-0.315	***	-0.291	***	0.250	***	0.479	***	-0.100	***
Mean HNR	-0.052	*	0.033	***	0.249	***	-0.093	***	-0.135	***	-0.056	*
Centre of Gravity	0.052	*	-0.074	***	0.199	***	-0.182	***	0.015		-0.077	**
SD Spectrum	0.105	***	0.034	***	-0.033	***	-0.027	***	-0.145	***	0.175	***
Adj R²	0.135	***	0.269	***	0.302	***	0.260	***	0.621	***	0.176	***
Adj R² (emo-corr)	0.092	***	0.305	***	0.132	***	0.159	***	0.270	***	0.142	***

954

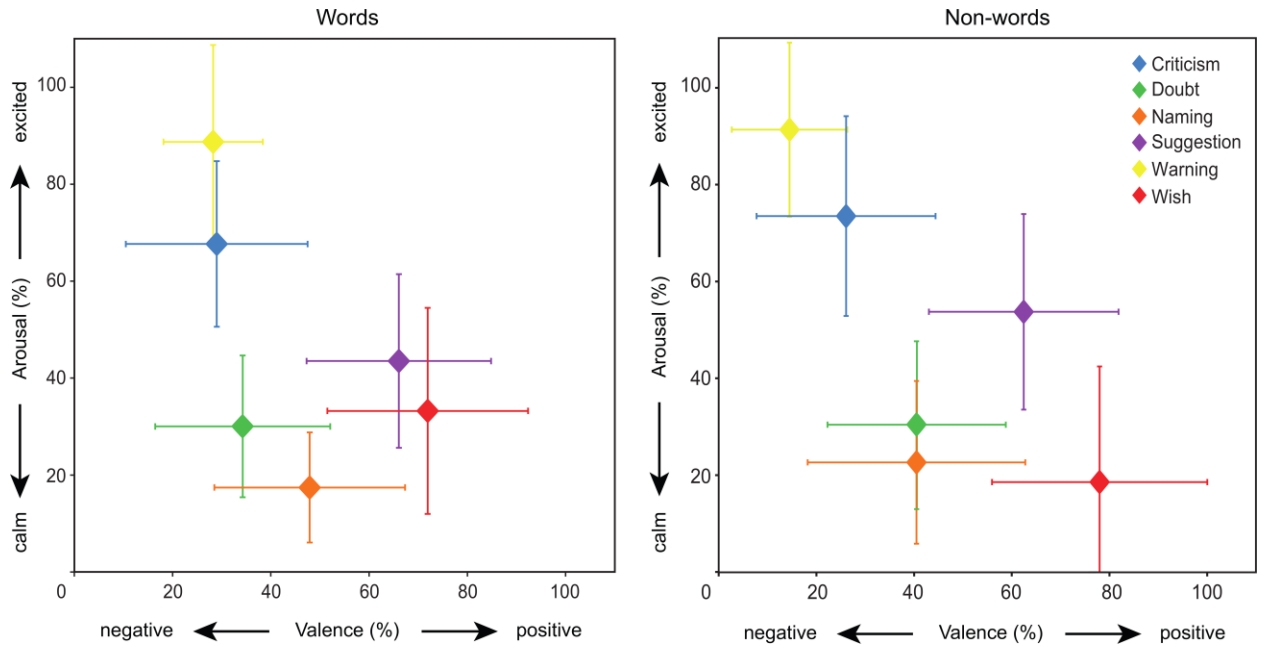
955 Beta weights and adjusted R² are depicted for multiple regressions using acoustic features as
 956 predictors and speech act ratings as dependent variables. Additionally, adjusted R² values are
 957 depicted after controlling for emotion perception (emo-corr). F0 = Fundamental frequency; HNR
 958 = Harmonics-to-noise ratio; SD = standard deviation; Adj = Adjusted; emo-corr = overall
 959 performance of the multiple regressions after affective ratings had been regressed out; * p < .05,
 960 ** p < .01, *** p < .001

961 **FIGURES**



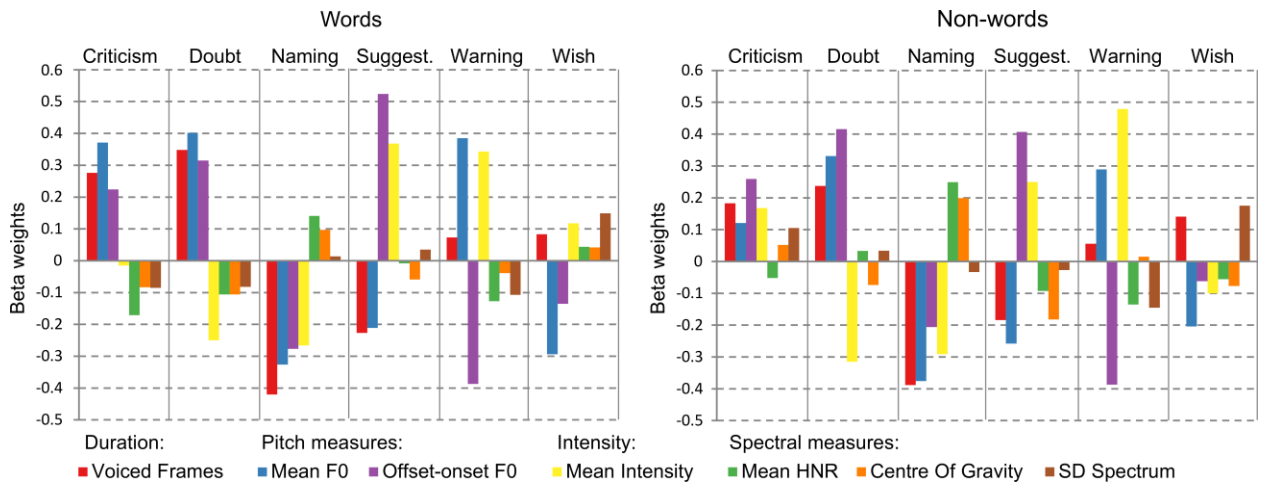
963 **Figure 1.** Results of the discriminant analyses plotted along the first and second
964 discriminant function. Each color corresponds to stimuli of one speech act category.

965



966

967 **Figure 2.** Emotion ratings. Average scores of the valence and arousal ratings for each
 968 speech act category.



969

970 **Figure 3.** Results of the multiple regressions for each speech act scale (columns). The
 971 bars represent the beta-weights for the seven acoustic features indicating how strongly they
 972 predicted the speech act rating.

973

974 **APPENDIX A**

975 Situation descriptions presented to the speakers for speech act recordings (English
976 translations). All speakers read the scenarios and were asked to place themselves in the
977 interpersonal situation, for example, by uttering the example sentences (in bold), before uttering
978 the single words (or non-words) with the corresponding prosody. Scenarios contained either
979 words or non-words (in square brackets below).

980 *Criticism*

981 Your colleague Tom and you will present your first big job in an important meeting this
982 afternoon. Therefore, you are extremely nervous and do not want to disappoint your boss. You
983 sit at your desk and go through the presentation one more time. Suddenly, there is a knock on the
984 door and Tom peeks into your office. He asks you whether you would like to join him for a beer
985 [diem] in the bar [dahm], although he knows how important the upcoming meeting is. You think
986 it would be best for him to prepare the joint talk and ask disapprovingly:

987 **BIER [DIEM]: “(Are you serious? A) beer [diem], (now)?”**

988 **BAR [DAHM]: “(Are you serious? The) bar [dahm], (now)?”**

989 *Doubt*

990 You arrive home after a hard day at work and it is quite late. Your mobile rings when you
991 have just hung your coat up. It is your friend, Eva, and she suggests having a beer [diem] at the
992 bar [dahm]. You would actually like to meet her as you have not seen her for a long time, but are

993 very tired and need to leave for work early the next day. You do not know whether this is a good
994 idea. Therefore, you ask doubtfully:

995 BEER [DIEM]: **“(A) beer [diem], (now)?”**

996 BAR [DAHM]: **“(The) bar [dahm], (now)?”**

997 *Naming*

998 Please say the words: beer / bar / diem / dahm with a neutral intonation, for example, as
999 in the sentence:

1000 **“(I’m going to have a) beer [diem] (tonight).”**

1001 **“(I’m going to a) bar [dahm] (tonight).”**

1002 *Suggestion*

1003 It is Thursday evening and you have almost finished your work. You achieved a lot today
1004 and are satisfied with your work. You really deserve to go to the bar [dahm] for a beer [diem]
1005 now. You think you can perhaps convince your colleague, Anne, to join as you sometimes go to
1006 your favorite bar together after work. In pleasant anticipation of a nice evening, you peek into
1007 Anne’s office and ask invitingly:

1008 BEER [DIEM]: **“(Are you up for a) beer [diem]?”**

1009 BAR [DAHM]: **“(Do you want to go to a) bar [dahm]?”**

1010

1011 *Warning*

1012 You invited a friend to have a beer at your apartment. He talks excitedly about his last
1013 football match and vividly tries to imitate one of his maneuvers. He spins around wildly, back
1014 and forth, left and right, and you start getting worried about your furniture. He suddenly starts
1015 running and does not see your mini bar [dahm] where he put his glass of beer [diem]. You try to
1016 warn him:

1017 BEER [DIEM]: “**(Watch out, your) beer [diem]!**”

1018 BAR [DAHM]: “**(Watch out, the) bar [dahm]!**”

1019 *Wish*

1020 It is a hot summer day and you descend after a hard, but wonderful mountain hike. After
1021 all those kilometers and the great view, you are pleasantly exhausted, hungry, and thirsty—the
1022 hotel is almost within sight. You only have one thought on your mind: You would like a nice
1023 cool beer [diem] at the hotel bar [dahm] to make this day truly perfect. You say longingly:

1024 BEER [DIEM]: “**(Now for a) beer [diem]!**”

1025 BAR [DAHM]: “**(Now to the) bar [dahm]!**”

1026

1027 **APPENDIX B**

1028 Definitions of speech acts presented to participants before the behavioral tests (English
1029 translations).

1030 *Criticism*

1031 The speaker, a friend of yours, is disapprovingly expressing criticism, for example, about
1032 one of your suggestions.

1033 *Doubt*

1034 The speaker is deliberately expressing doubt, for example about whether to accept a
1035 proposal you made.

1036 *Naming*

1037 The speaker is saying something for no specific purpose, for example, to name an object.

1038 *Suggestion*

1039 The speaker is invitingly suggesting something, for example, to undertake something
1040 together.

1041 *Warning*

1042 The speaker is warning you of a possible accident, for example, not to fall over an object.

1043 *Wish*

1044 The speaker is longingly expressing a wish for something, for example, a relaxing
1045 evening after a successful working day.

1046 **APPENDIX C**1047 **Table C1.** Mean acoustic features per speech act category.

Speech Act	Acoustic feature						
	Number of Voiced Frames	Mean F0 (Hz)	Offset-onset F0 (Hz)	Mean Intensity (dB)	Mean HNR (dB)	Spectral Centre of Gravity (Hz)	SD Spectrum (Hz)
<i>Words</i>							
Criticism	450.1 ± 61.1	230.7 ± 48.4	81.5 ± 91.6	65.1 ± 3.9	12.1 ± 3.2	712.1 ± 263.8	728.6 ± 243.4
Doubt	482.8 ± 68.7	188.9 ± 41.6	72.6 ± 30.9	57.2 ± 3.6	14.2 ± 3.4	511.4 ± 186.0	713.6 ± 204.7
Naming	341.1 ± 64.5	13.3 ± 42.6	-51.6 ± 29.0	56.9 ± 4.2	13.3 ± 2.3	587.7 ± 248.1	554.4 ± 139.1
Suggestion	320.0 ± 56.4	206.1 ± 37.2	184.7 ± 55.2	63.3 ± 2.0	13.2 ± 2.7	617.1 ± 224.0	611.7 ± 136.0
Warning	428.4 ± 97.9	268.5 ± 49.4	-122.7 ± 27.0	71.8 ± 2.3	13.9 ± 2.8	897.1 ± 209.3	762.2 ± 244.4
Wish	485.7 ± 62.9	148.6 ± 39.5	-61.7 ± 16.4	59.1 ± 3.1	12.7 ± 2.1	591.1 ± 238.2	781.1 ± 319.6
Average	418.0 ± 95.3	196.9 ± 62.3	17.2 ± 115.5	62.2 ± 6.2	13.2 ± 2.9	652.8 ± 259.8	652.8 ± 237.0
<i>Non-Words</i>							
Criticism	473.0 ± 81.1	250.1 ± 60.7	118.9 ± 112.6	64.8 ± 3.6	15.5 ± 4.7	645.5 ± 279.2	806.2 ± 227.4
Doubt	510.7 ± 57.5	185.8 ± 47.5	70.9 ± 30.8	57.0 ± 4.1	18.9 ± 4.1	386.8 ± 118.9	641.9 ± 230.0
Naming	427.3 ± 72.9	134.1 ± 41.7	-56.8 ± 28.4	55.2 ± 4.2	17.2 ± 2.8	504.2 ± 235.0	597.5 ± 174.1
Suggestion	342.0 ± 50.8	217.4 ± 41.0	213.2 ± 49.5	62.2 ± 3.1	16.1 ± 3.6	455.0 ± 169.0	567.9 ± 154.4
Warning	474.2 ± 124.9	280.5 ± 46.7	-125.3 ± 23.3	71.4 ± 1.6	18.1 ± 2.6	755.1 ± 288.8	713.7 ± 184.0
Wish	560.8 ± 87.1	142.3 ± 38.0	-53.2 ± 19.0	57.9 ± 3.5	15.9 ± 3.9	494.5 ± 188.3	760.9 ± 205.4
Average	464.7 ± 106.8	201.7 ± 70.7	27.9 ± 128.7	61.4 ± 6.5	17.0 ± 3.9	540.2 ± 252.5	681.4 ± 214.6

1048
1049 Values depict mean ± *SD*. HNR = harmonics-to-noise ratio; *SD* = standard deviation. All values were extracted using PRAAT 5.3.01
1050 (<http://www.praat.org>).

1051

1052 **Table C2.** Statistical comparison of acoustic features between words and non-words.

Acoustic parameter	Speech act ratings											
	Criticism		Doubt		Naming		Suggestion		Warning		Wish	
	t(6)	p	t(6)	p	t(6)	p	t(6)	p	t(6)	p	t(6)	p
Voiced Frames	-0.453	0.666	-1.376	0.218	-1.831	0.117	-0.710	0.504	-0.512	0.627	-1.679	0.144
Mean F0	-0.427	0.684	-0.027	0.980	0.124	0.906	-0.191	0.854	-0.272	0.795	0.251	0.810
Offset-onset F0	-0.299	0.775	0.277	0.791	0.630	0.552	-0.618	0.559	-0.198	0.849	-1.652	0.150
Mean Intensity	0.159	0.879	-0.191	0.855	0.605	0.567	0.837	0.435	0.105	0.920	0.903	0.401
Mean HNR	-1.735	0.134	-7.098	0.000	-2.749	0.033	-2.317	0.060	-2.659	0.038	-2.309	0.060
Centre of Gravity	0.605	0.567	2.619	0.040	1.193	0.278	4.292	0.005	2.011	0.091	2.562	0.043
SD Spectrum	-0.962	0.373	0.611	0.564	-0.415	0.693	0.689	0.516	0.443	0.673	-0.008	0.994

1053

1054 Acoustic features of words and non-words were compared with paired *t*-tests for each speech act category (columns). Significant

1055 results ($p < 0.05$) are marked in bold. F0 = Fundamental frequency; HNR = Harmonics-to-noise ratio; SD = standard deviation.

1056 **Table C3.** Results of the discriminant analyses (Experiment 1).

	Words	Non-Words
<i>Function 1</i>		
Offset-onset f0	0.881	0.836
<i>Function 2</i>		
Mean Intensity	0.721	0.716
Mean f0	0.467	0.538

1057
1058 Within-group correlations between acoustic
1059 measures and standardized canonical
1060 discriminant functions. The table includes
1061 all values of the first two functions above a
1062 threshold of $r = 0.3$.
1063

1064 **SUPPLEMENTARY MATERIAL**

1065 Supplementary material available via

1066 <https://www.dropbox.com/sh/vu6q0g8f9stvm35/AACbDBBIdyHKFnHVIDyrajHNa?dl=0>

1067 **Audio S1_criticism.** Example stimulus “criticism”.

1068 **Audio S2_doubt.** Example stimulus “doubt”.

1069 **Audio S3_naming.** Example stimulus “naming”.

1070 **Audio S4_suggestion.** Example stimulus “suggestion”.

1071 **Audio S5_warning.** Example stimulus “warning”.

1072 **Audio S6_wish.** Example stimulus “wish”.