



When the words are not everything: the use of laughter, fillers, back-channel, silence and overlapping speech in phone calls

Alessandro Vinciarelli^{1,*}, Paraskevi Chatziioannou¹, and Anna Esposito²

¹School of Computing Science, University of Glasgow, Glasgow, United Kingdom ²Psychology Department, Second University Naples, Caserta, Italy

Correspondence*: Alessandro Vinciarelli School of Computing Science, University of Glasgow, Sir A.Williams Building, Glasgow, G12 8QQ, United Kingdom, vincia@dcs.gla.ac.uk

2 ABSTRACT

This article presents an observational study on how some common conversational cues -3 laughter, fillers, back-channel, silence, and overlapping speech - are used during mobile phone 4 conversations. The observations are performed over the SSPNet Mobile Corpus, a collection 5 of 60 calls between pairs of unacquainted individuals (120 subjects for roughly 12 hours of 6 material in total). The results show that the temporal distribution of the social signals above is 7 not uniform, but it rather reflects the social meaning they carry and convey. In particular, the 8 9 results show significant use differences depending on factors such as gender, role (caller or receiver), topic, mode of interaction (agreement or disagreement), personality traits and conflict 10 handling style. 11

12 Keywords: Nonverbal Behaviour, Social Signals, Laughter, Back-Channel, Fillers, Pauses, Overlapping Speech, Corpus Analysis

1 INTRODUCTION

In general terms, nonverbal communication is the "process of one person stimulating meaning in the 13 mind of another person or persons by means of nonverbal messages" (Richmond et al., 1991). In 14 face-to-face conversations, people have at disposition a wide spectrum of cues - facial expressions, 15 gestures, mutual distances, posture, etc. - to accomplish nonverbal communication and enrich the words 16 being exchanged with multiple layers of meaning (social, psychological, emotional, etc.). However, the 17 situation changes dramatically in phone-mediated conversations where all the functions typically fulfilled 18 via nonverbal communication - e.g., conveying impressions, sending relational messages, expressing 19 emotions, etc. (Hecht et al., 1999) - must be constrained through the only available channel, i.e. speech. 20 The main difficulty in this case is that the same vocal apparatus must be used for both verbal and nonverbal 21 22 components of communication and, in some cases, one component can be used only at the expense of the 23 other (e.g., it is difficult to speak and laugh at the same time).

Given the above, the temporal distribution of nonverbal vocal cues should not be uniform, but rather correspond to the function and role of nonverbal communication in human-human interactions. For example, in the case of laughter, "the temporal segregation of speech and laughter on the single vocalization channel reveals the presence or absence of an underlying organizational principle" (**Provine**, 1993). More generally, "the circumstances in which an activity is performed and those in which it never

29 occurs [provide] clues as to what the behaviour pattern might be for (its function)" (Martin and Bateson, 2007). For these reasons, this article proposes an analysis of the temporal distribution of several nonverbal 30 vocal cues - laughter (audible contractions, typically rhythmical, of the diaphragm and other parts of the 31 respiratory system), fillers (expressions like "ehm" that fill the time intervals that should be occupied 32 by a word), back-channel (short voiced utterances like "*ah-ah*" that signal attention and encouragement 33 to continue to others), silence (time intervals during which nobody speaks or produces audible sounds) 34 35 and overlapping speech (time intervals during which at least two speakers talk at the same time) - in the SSPNet Mobile Corpus (Polychroniou et al., 2014), a collection of 60 phone calls between unacquainted 36 individuals (120 subjects in total). In particular, the observations show that the distribution of the cues 37 changes according to the following factors expected to account for the relational context: gender (male vs 38 *female*), role (*caller* vs *receiver*), topic of conversation (*task* vs *social*), mode of interaction (*agreement* vs 39 *disagreement*), Big-Five personality traits (McCrae, 2009), and conflict handling style (Rahim, 1983). 40

The *rationale* behind the choice of the cues above is that they tend to appear frequently in conversations (see Section 2) and this is probably an indication of their primacy in human-human communication. Furthermore, the five cues are the subject of extensive work in Social Signal Processing (**Vinciarelli et al.**, 2012), Computational Paralinguistics (**Schuller and Batliner**, 2013) and Human-Media Interaction (**Nijholt**, 2014), three computing domains involving automatic detection and interpretation of nonverbal behavioral cues. In this respect, the findings of this work can be helpful for automatic approaches aimed at automatically making sense of social interactions.

48 Overall, the observations of this work try to address the following three main questions:

- What are the physical, possibly machine detectable traces of socially relevant factors like gender, role,
 topic of conversation, mode of interaction, personality traits and conflict handling style?
- Is it possible to transfer observations made about face-to-face interactions to phone-mediated conversations?
- Does the use of phones introduce effects and biases that are not observed (or not applicable) in faceto-face interactions?

The results of the observations show that, far from distributing uniformly over time, nonverbal cues appear 55 with different frequency depending on the relational context factors. Therefore, the frequency of the cues 56 can be considered one of the physical traces that contextual factors leave. Furthermore, the results show 57 that several status and dominance effects observed in face-to-face interactions seem to apply in the case 58 59 of phone-mediated conversations as well. Hence, observations about co-located interactions appear to transfer, at least partially, to phone calls. Finally, the results show that the difference between calling or 60 receiving (peculiar of phone calls and not applicable to face-to-face encounters) tends to be perceived as 61 a difference in terms of status and dominance. Therefore, the use of phones induces peculiar effects that 62 are not observed in other interactional settings. 63

The rest of this article is organized as follows: Section 2 describes the Corpus and its scenario, Section 3 describes the methodology adopted in this work, Section 4 illustrates observations and findings, and Section 5 draws some conclusions.

2 THE SSPNET MOBILE CORPUS: SCENARIO AND CUES

The observations of this work are performed over the SSPNet Mobile Corpus (**Polychroniou et al.**, 2014), a collection of 60 phone calls between unacquainted individuals (120 subjects in total). The conversations, minutes and 24 seconds in total, revolve around the *Winter Survival Task* (see below) and are annotated in terms of the cues mentioned in Section 1, namely laughter, fillers, back-channel, silence and querlanging speech. The root of this section provides further details about both task and cues

and overlapping speech. The rest of this section provides further details about both task and cues.



Figure 1. The picture shows the experimental protocol. The subjects fill self-assessment questionnaires in the days before the call. The day of the call, they sit in one of the two offices used for the experiment (step 1), they read the protocol (step 2), they address the WST task individually (step 3), they receive a mobile phone (step 4), they negotiate a common solution during a call (step 5), they deliver a negotiated solution (step 6).

2.1 THE WINTER SURVIVAL TASK

The Winter Survival Task (WST) requires the participants to consider a list of 12 items (*steel wool, axe*, *pistol, butter can, newspaper, lighter without fuel, clothing, canvas, airmap, whisky, compass, chocolate*) and to identify those that can increase the chances of survival after a plane crash in Northern Canada (**Joshi et al.**, 2005). Before participating in the experiment, the participants have been asked to fill the *Big-Five Inventory* 10 (**Rammstedt and John**, 2007) and the *Rahim Organizational Conflict Inventory II* (**Rahim**, 1983), two questionnaires that measure personality traits (see Section 4.5 for details) and conflict handling style (see Section 4.6 for details), respectively.

Figure 1 illustrates the experimental protocol adopted for collecting the data. After having filled the questionnaires, the participants have been admitted to the experiment and the calls have been collected as follows:

- Step 1: The two subjects involved in the same call are conducted to two different rooms of the School of Computing Science at the University of Glasgow (the two subjects never enter in contact with one another before the call).
- Step 2: Once in their room, the participants receive the same document that explains the WST and are asked to read it carefully (the document includes the list of the 12 items at the core of the task).
- Step 3: Before starting the call, the subjects address the WST by filling a form where, for each of the 12 items, they have to tick a "*Yes*" or "*No*" box. A positive answer means that the item can increase the chances of survival and viceversa for the negative answer. The participants are asked to tick a box for each of the items (the call cannot start if any item is left blank).
- Step 4: The two subjects receive a mobile phone (the same model for both participants).
- Step 5: One of the two subjects, selected randomly, calls the other with the mobile phone provided by the experimenters.
- Step 6: During the call, the two subjects have to negotiate a common solution for the WST. Every time 95 they have ticked a different box about an item, one of the two participants has to shift to the decision

made by the other participant. The items have to be discussed one-by-one following the same order
 for all pairs; The call cannot be interrupted until a common decision has been reached for all items.

At the end of the call, the participants have received a payment that includes a fixed sum of $\pounds 6$ and a bonus calculated as follows: the WST has a golden standard that shows what are the items for which the box "*Yes*" should be ticked. Each time the participants tick the box "*Yes*" for one of these, they earn $\pounds 3$. However, if the participants tick the box *Yes* for an item for which the golden standard says "*No*", then they loose $\pounds 3$. If the bonus is negative (the false positives are more frequent than the true positives), the participants do not receive any extra bonus.

2.2 THE CUES

Figure 2 shows the distribution of the cues in terms of both occurrences and percentage of the Corpus
duration covered by each of them. The high frequency of all cues (16,235 occurrences in total) confirms
their primacy in human-human communication.

107 Laughter is "a common, species-typical human vocal act and auditory signal that is important in social 108 discourse" (Provine and Yong, 1991). Seminal findings about the temporal distribution of laughter in conversations have been proposed by **Provine** (1993), including the tendency of women to laugh 109 more than men, the tendency of listeners to laugh less than speakers, and the tendency to laugh only 110 when a sentence has been completed. More recently, laughter was found to signal topic changes in 111 spontaneous conversations (Bonin et al., 2014). This article confirms some of the previous observations 112 113 while proposing new effects that can emerge in the particular scenario of the SSPNet Mobile Corpus. 114 Figure 2 shows that the laughter occurrences in the Corpus are 1805 for a total duration of 1,114.8 seconds (2.6% of the total length of the corpus). When the speakers laugh together, the cue is counted twice. 115

Fillers are expressions like "*ehm*" and "*uhm*" that "*are characteristically associated with planning problems* [...] *planned for, formulated, and produced as parts of utterances just as any word is*" (**Clark and Fox Tree**, 2002). This means that speakers replace words with fillers when, e.g., they need time to look for the right term, they plan what to say next or they try to hold the floor. According to the distribution of Figure 2, the Corpus includes 3,912 filler occurrences that account for 1,815.9 seconds (4.2% of the total corpus time).

Another frequent event in human-human conversations is *back-channel*, i.e. the use of "*short utterances produced by one participant in a conversation while the other is talking*" (Ward and Tsukahara, 2000). In English, this corresponds to expressions like "*yeah*", "*aha-aha*", etc., that signal, in most cases, attention and agreement. Figure 2 shows that the speakers of the Corpus perform back-channel 1,015 times, for a total of 407.1 seconds (0.9% of the Corpus time).

Silence is the most frequent cue among those considered in this work: 6,091 occurrences for a total 127 128 of 4,670.6 seconds (10.9% of the corpus length). In some cases, silence accompanies the grammatical structure of the speech stream (e.g., a short silence can signal the end of a sentence), in others it manifests 129 130 hesitation in planning the next words or it is a latency time between questions and answers (Hall and Knapp, 1992). Furthermore, silence can serve communication purposes: "the main common link between 131 speech and silence is that the same interpretive processes apply to someone's remaining meaningfully 132 silent in discourse as to their speaking" (Jaworski, 1999). The observations of this work do not take into 133 account the differences mentioned above, but show that the frequency of silences changes according to 134 some relational context factors (see Section 4). 135

According to **Schegloff** (2000), "*Talk by more than one person at a time in the same conversation is one* of the two major departures that occur from what appears to be a basic design feature of conversation, [...] namely 'one at a time' (the other departure is silence, i.e. fewer than one at a time)". For this reason, the observations of this work take into account both silence (see above) and overlapping speech, i.e. the time intervals during which the two subjects involved in the same call talk simultaneously. The number



Figure 2. The left chart shows the number of occurrences for the cues considered in this article. The right chart shows the percentage of time covered by each cue in the Corpus.

of occurrences for this cue is 3,412 for a total of 2,000.5 seconds (4.7% of the corpus time). Unless there
is competition for the floor, simultaneous speakers resolve overlapping quickly to move back to the "*one at a time*" situation (the average duration of overlapping speech segments in the Corpus is 0.58 seconds).
Like in the case of the other cues, the observations of this work show how the frequency of overlapping
speech segments changes in different parts of a conversation.

3 METHODOLOGY

The goal of this work is to show whether the frequency of nonverbal cues changes according to six factors
expected to account for the relational context, namely gender (see Section 4.1), role (see Section 4.2), topic
of conversation (see Section 4.3), mode of interaction (see Section 4.4), personality traits (see Section 4.5)
and conflict handling style (see Section 4.6).

Each factor is modeled as a variable V that can take L values (numeric or nominal). For example, in the case of gender, the variable V can take 2 values, i.e. *male* and *female*. Given V, the Corpus can be segmented into intervals that correspond to one of the values of V. In the case of gender, this corresponds to segment the Corpus into intervals where the speaker is female and intervals where it is male. As a result, a fraction p_f of the corpus time corresponds to female speakers while a fraction p_m corresponds to male ones, with $p_f + p_m = 1, 0 < p_f < 1$ and $0 < p_m < 1$. In more general terms, if the variable V associated to a factor can take L values v_1, v_2, \ldots, v_L , the Corpus can be segmented into L subsets that account for fractions of the total time p_1, p_2, \ldots, p_L of the total time, where $0 < p_k < 1 \ \forall k$ and $\sum_{k=1}^{L} p_k = 1$.

If a cue (i.e., laughter, fillers, back-channel, silence or overlapping speech) occurs N times in the Corpus and its temporal distribution does not depend on the factor associated to V, the *expected* number of occurrences in the Corpus intervals where $V = v_k$ will be $E_k = N \cdot p_k$. For example, in the case of gender, the *expected* numbers of occurrences in correspondence of female and male speakers will be $E_f = N \cdot p_f$ and $E_m = N \cdot p_m$, respectively. However, the *observed* number of occurrences, i.e. the number of occurrences actually counted in the Corpus intervals where $V = v_k$, will be O_k . In the gender examples, O_f will be the number of times that female speakers actually display the cue while O_m will be the number of times that male ones do it. This allows one to define the following χ^2 variable:

$$\chi^2 = \sum_{k=1}^{L} \frac{(O_k - E_k)^2}{E_k},\tag{1}$$

where the number of Degrees of Freedom is L - 1. Such a variable can be used to test whether the null hypothesis is true (there is no statistically significant difference between observed and expected 168 distribution) or it must be rejected. In other words, the χ^2 variable above can tell us whether the frequency 169 of a given cue changes to a statistically significant extent depending on the value of V. In the case of 170 gender, if the null hypothesis can be rejected, it means that speakers of a given gender tend to display a 171 certain cue significantly more frequently than those of the other gender or viceversa.

172 The process for verifying whether a deviation with respect to the expected distribution is statistically 173 significant with confidence level α is as follows:

- 174 1. The value of the χ^2 variable resulting from the observed distribution is calculated.
- 175 2. The *p*-value corresponding to the χ^2 value is estimated;
- 176 3. If the *p*-value estimated at step 2 is lower than α/k , where α is the desired confidence level and k = 79 is the total number of statistical inferences made over the data, then the deviation is considered
- 178 statistically significant with confidence level α .

179 In other words, an effect is considered statistically significant with confidence level 0.01 when the *p*-value 180 is lower than 0.01/79 = 0.0001. Similarly, an effect is considered statistically significant with confidence 181 level 0.05 when the *p*-value is lower than 0.05/79 = 0.0006. Such a practice, known as Bonferroni 182 Correction, is typically applied when making a large number of statistical inferences over the same data 183 like it happens in this work. The Bonferroni Correction is subject to criticism because it reduces the 184 number of false positives at the cost of increasing significantly the number of false negatives (**Nakagawa**, 185 2004). However, it allows one to concentrate the analysis on the stronger effects observed in the Corpus.

4 CORPUS ANALYSIS

This section adopts the methodology described in Section 3 to test whether the frequency of nonverbalcues changes according to relational context factors, i.e. gender, role, topic, mode of interaction,personality and conflict handling style.

4.1 GENDER EFFECTS

The gender variable can take two values, male and female. In the SSPNet Mobile Corpus, male subjects 189 are 57 (47.5% of the total) and female ones are 63 (52.5% of the total). However, male subjects speak 190 54.5% of the time and this means that they tend to talk longer than female ones to a statistically significant 191 extent (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 test with Bonferroni Correction). A possible explanation is that the negotiation scenario adopted in the Corpus (**Polychroniou** 192 193 et al., 2014) activates gender stereotypes according to which "women are less assertive and agentic than 194 men" (Thompson et al., 2010). Therefore, women might tend to speak less, on average, than men. Further 195 196 confirmation comes from the duration of male-male calls that tend to be, on average, significantly longer than calls where at least one of the two speakers is female (see end of this section). 197

In absence of further gender effects, 54.5% of the occurrences of a cue should be displayed by male 198 subjects because these speak 54.5% of the total time (see Section 3). However, Figure 3 shows that there 199 are statistically significant deviations with respect to such an expectation. In particular, female subjects 200 201 tend to display laughter and back-channel significantly more frequently than male subjects (statistically significant with confidence level $\alpha = 0.01$ in both cases according to a χ^2 test with Bonferroni Correction). Furthermore, female subjects initiate overlapping speech significantly more frequently than 202 203 male ones (statistically significant with confidence level $\alpha = 0.01$ in both cases according to a χ^2 test 204 with Bonferroni Correction). Gender effects for the other cues, if any, are too weak to be observed in the 205 Corpus. 206



Figure 3. The charts show gender differences in the distribution of cues' occurrences. In particular, the left chart shows the distribution for female subjects while the right one shows it for male ones. The double asterisk means that the deviation is statistically significant with confidence level $\alpha = 0.01$ (according to a χ^2 test with Bonferroni Correction).

207 This pattern is compatible with a large body of work showing that "men and women are generally 208 perceived as differing in status (importance, dominance, power, etc.) and also that they often feel themselves to differ in this way" (Leffler et al., 1982). In other words, even if the scenario of the 209 Corpus does not introduce a status difference between subjects and there is no status difference between 210 male and female subjects (Polychroniou et al., 2014), males are still more likely to adopt behaviors 211 typical of higher-status individuals, including speaking longer (see above), laughing less (**Provine**, 1993; 212 Leffler et al., 1982) and showing back-channel less frequently (Hall et al., 2005). The only contradictory 213 evidence is that female subjects tend to initiate overlapping speech significantly more often than male 214 ones. A possible explanation is that female subjects initiate overlapping speech more often to avoid a 215 216 stereotype threat, i.e. the risk to confirm negative stereotypes about a category someone belongs to (Steele 217 and Aronson, 1995). In this case, female subjects might interrupt more to contradict the sterotype that depicts women are less assertive and agentic (Thompson et al., 2010). 218

When it comes to gender composition, the SSPNet Mobile Corpus includes 17 female-female calls 219 220 (28.3% of the total time), 14 male-male calls (23.3% of the total time) and 31 female-male calls (48.4%) of the total time). The average duration of female-female, male-male and female-male calls is 595 s, 899 s 221 and 639 s, respectively. Therefore, male-male pairs seem to need significantly more time to complete a call 222 (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 test with Bonferroni Correction). 223 In absence of further gender composition effects, nonverbal cues should distribute over the three types of 224 225 call according to the same proportions indicated above, namely 28.3% (female-female), 23.3% (malemale) and 48.4% (female-male). However, the observed distribution is significantly different from the 226 227 expected one for fillers and silences (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 test with Bonferroni Correction). 228

229 A possible explanation of the latter observations is that male subjects tend to compete more to hold the 230 floor. This can explain the higher frequency of fillers - one of the functions of such a cue is to keep the floor while planning what to say next or addressing any other communication performance problem (Clark and 231 Fox Tree, 2002; Hall and Knapp, 1992) - as well as lower frequency of silences. Higher competition in 232 holding the floor might contribute to explain the longer duration of male-male calls as well. In fact, 233 234 competition to hold the floor is typically associated to higher levels of conflict (Schegloff, 2000; Smith-Lovin and Brody, 1989) that result into longer negotiations before reaching a consensual solution for the 235 Winter Survival Task (see Section 4.4). 236



Figure 4. The charts show the effects of roles (*callers* on the left and *receivers* on the right) measured in terms of frequency of the cues under examination. The double asterisk means that the deviation is statistically significant with confidence level $\alpha = 0.01$ (according to a χ^2 test with Bonferroni Correction).

4.2 ROLE EFFECTS: CALLING VS RECEIVING

The scenario of the SSPNet Mobile Corpus does not introduce any difference between two subjects involved in the same call (**Polychroniou et al.**, 2014). However, given that the conversations take place over the phone, one subject plays the role of the *caller* (the person that makes the call) while the other one plays the role of the *receiver* (the person that receives the call). For every pair, the two roles were assigned randomly. This section adopts the methodology of Section 3 (the variable takes the values *caller* and *receiver*) to test whether the role has any effect on the frequency of nonverbal cues.

By design, 50% of the subjects are callers while the other 50% are receivers. The former speak 49.9% 243 of the time and the latter 50.1%. According to a χ^2 test with Bonferroni Correction, the difference is 244 not significant and the effect of role on speaking time, if any, is too weak to be observed in the Corpus. 245 If the same applies to the cues under examination, 49.9% of their occurrences should be displayed by 246 callers and the remaining 50.1% by receivers. However, Figure 4 shows that, to a statistically significant 247 extent, callers tend to display fillers more often while receivers tend to initiate overlapping speech more 248 frequently (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 test with Bonferroni 249 Correction). 250

Initiating overlapping speech is typically associated with dominance (Anderson and Leaper, 1998) and 251 higher-status (Leffler et al., 1982). Therefore, a possible reason why receivers initiate overlapping speech 252 significantly more frequently than callers is that they tend to be perceived, and perceive themselves, as 253 higher-status individuals. As a possible confirmation, previous results obtained over the SSPNet Mobile 254 Corpus (Vinciarelli et al., 2014) show that receivers persuade callers 70% of the times (statistically 255 significant with confidence level $\alpha = 0.01$ according to a χ^2 test with Bonferroni Correction) when 256 there is disagreement about one item of the Winter Survival Task (see Section 4.4): being perceived as a 257 higher-status individual is one of the two factors that "stand out as particularly powerful determinants of 258 a person's ability to influence others" (Oldmeadow et al., 2003). The tendency of callers to display fillers 259 more frequently goes in the same direction, given that "when communicating with a higher status person, 260 the lower status person [...] has more filled and unfilled pauses than normal" (Richmond et al., 1991), 261 where the "filled pauses" correspond to the fillers of this work. Overall, while perceived status has not 262 been measured and the Corpus scenario does not involve any status difference, role related effects seem 263 to be compatible with a situation where the receiver is perceived to be higher in status. 264



Figure 5. The charts show how the distribution of the cues changes depending on whether the subjects establish social contact (left chart) or address the Winter Survival Task (right chart). The double asterisk means that the deviation is statistically significant with confidence level $\alpha = 0.01$ (according to a χ^2 test with Bonferroni Correction).

4.3 TOPIC EFFECTS

The calls of the SSPNet Mobile Corpus revolve around the Winter Survival Task (**Polychroniou et al.**, 2014). The two subjects involved in each call are asked to identify objects that are likely to increase the chances of survival in a polar environment (**Joshi et al.**, 2005). The subjects spend only 90.3% of the total Corpus time in addressing the task. The remaining 9.7% is dedicated to mutual introductions, small-talk, greetings, comments about the experiment and other activities that, in general, aim at establishing a social contact between fully unacquainted subjects. This allows one to apply the methodology of Section 3 with a variable that takes the values *task* and *social*.

Figure 5 shows that laughter, silence and overlapping speech are significantly more frequent than expected when the subjects do not address the task (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 test with Bonferroni Correction) and viceversa for back-channel (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 test with Bonferroni Correction).

In the case of laughter, a possible explanation is that the cue is "*important in social discourse*" (**Provine** and Yong, 1991). Therefore, it probably tends to appear more frequently when interactions are socially rather than task oriented. For what concerns fillers and silence, one possible explanation is that these cues can account for communication difficulties between fully unacquainted individuals speaking to one another for the first time. In particular, silences can reflect a difficulty in planning what to say next in absence of a predefined topic of conversation (**Hall and Knapp**, 1992). Furthermore, overlapping speech might account for lack of coordination in turn-taking (**Schegloff**, 2000).

4.4 MODE OF INTERACTION EFFECTS: AGREEMENT VS DISAGREEMENT

Before participating in the experiment, the two subjects involved in the same call are asked to look at 283 a list of 12 items and decide, for each of them, whether it increases the chances of survival in a polar 284 285 environment or not. In this way, it is possible to know whether the two subjects agree (they have made the same decision) or disagree (they have made a different decision) about an item, given that agreement 286 can be defined as "a relation of identity, similarity or congruence between the opinions of two or more 287 *persons*" (**Poggi et al.**, 2011). During the call, the two subjects are asked to discuss the items sequentially, 288 289 one at a time, and to reach a consensual decision for each of them. As a result, the Corpus can be segmented into 720 discussions (12 items \times 60 calls) about individual items and, for each discussion, 290 it is possible to know whether the subjects agree or disagree. This allows one to adopt te methodology of 291



Figure 6. The charts show the percentage of total cues' occurrences displayed during agreement and disagreement, respectively. The double asterisk means that the deviation is statistically significant (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 test with Bonferroni Correction).

292 Section 3 to test whether the mode of interaction has an effect on the frequency of nonverbal cues. The 293 variable used to segment the Corpus takes the values *agreement* and *disagreement*.

Disagreement is less frequent than agreement (283 out discussion of the total 720), but it accounts 294 295 for 61.7% of the total time spent on the task in the Corpus. The reason is that it takes more time to reach a consensual decision when the subjects have different opinions about a given item. If the mode of 296 297 interaction has no effect, 61.7% of a cue's occurrences (within statistical fluctuations) should be displayed during disagreement discussions. Figure 6 shows how the occurrences distribute over agreement and 298 disagreement. The observed distribution is statistically significantly different from the expected one for 299 silence and overlapping speech (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 300 test with Bonferroni Correction). 301

In the case of silences, the lower frequency during disagreement can have at least two possible reasons. 302 303 The first is that people tend to react immediately to interlocutors they disagree with (**Bilmes**, 1988), thus reducing, if not eliminating, the latency time in responding. The second is that conversation participants 304 tend to hold the floor during competitive interactions (Smith-Lovin and Brody, 1989) and, therefore, the 305 306 chances of observing silence are reduced. The lower frequency of overlapping speech seems to contradict 307 previous observations showing that the cue is associated with competitive interactions (Schegloff, 2000; Smith-Lovin and Brody, 1989). However, it should be noted that the conversations take place over 308 309 mobile phones and people cannot hear one another when they speak simultaneously more than a few hundreds of second. The need of mutual monitoring while negotiating a common solution to the task 310 might therefore lead people to reduce overlapping (McGinn and Croson, 2004). 311

4.5 PERSONALITY EFFECTS

312 Every subject of the SSPNet Mobile Corpus has filled the *Big-Five Inventory 10* (Rammstedt and 313 John, 2007), a questionnaire aimed at assessing personality in terms of the *Big-Five* traits (Saucier 314 and Goldberg, 1996): *Openness* (tendency to be intellectually curious, to have wide interests, etc.), 315 *Conscientiousness* (tendency to be planful, reliable, thorough, etc.), *Extraversion* (tendency to be 316 assertive, energetic, outgoing, etc.), *Agreeableness* (tendency to be kind, sympathetic, generous, etc.), and 317 *Neuroticism* (tendency to be anxious, self-pitying, touchy, etc.). The questionnaire allows one to calculate 318 five integer scores that measure how well an individual fits the tendencies associated to the Big-Five traits.



Figure 7. The bubble plot shows the deviation of the observed distribution of the cues with respect to the expected ones as a function of the personality traits. The larger the bubble, the larger the deviation. When the bubble is red, the deviation is negative (less occurrences than expected), when the bubble is blue the deviation is positive (more occurrences than expected). The double asterisk means that the deviation is statistically significant with confidence level $\alpha = 0.01$ (according to a χ^2 test with Bonferroni Correction).

319 The scores range in the interval [-4, 4] and, for each trait, it is possible to define a variable V that has 320 value low when the score is in the interval [-4, -2], middle when it is in [-1, 1], and high when it is in [2, 4]. This allows one to apply the methodology of Section 3. Figure 7 shows the deviations of the 321 observed frequencies with respect to the expected ones. In particular, the size of a circle is proportional 322 to the ratio (O - E)/E, where O is the number of times a cue actually occurs and E is the number of 323 times the cue is expected to occur. The circle is blue when O > E and red otherwise. The stars are plotted 324 in correspondence of deviations statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 325 test with Bonferroni Correction. 326

327 In the case of laughter, there are statistically significant deviations for Extraversion and Conscientiousness (statistically significant with confidence level $\alpha = 0.01$ according to a χ^2 test with 328 Bonferroni Correction). Previous findings suggest that "the tendency to laugh is a characteristic of the 329 extraverted person, albeit the empirical basis for this assertion is somewhat meagre" (Ruch and Deckers, 330 1993). The observations of the Corpus seem to be coherent with such an indication and, in particular, show 331 that subjects scoring *low* in Extraversion laugh significantly less often than the others (one of the largest 332 deviations in Figure 7). However, the deviation is negative in the case of subjects scoring *high* as well 333 334 and only people scoring *middle* appear to laugh significantly more than expected. In this respect, the 335 observations of the Corpus confirm that the evidence of an association between laughter and Extraversion is "meager". Similar considerations apply to Conscientiousness, with the subjects scoring middle that 336 laugh more frequently than the others to a statistically significant extent. A possible explanation is that the 337 subjects scoring *low* do not feel comfortable and/or motivated in addressing the task and, therefore, tend 338 to laugh less. At the opposite extreme, subjects scoring *high* tend to remain concentrated on the task and, 339 in line with the findings of Section 4.3, reduce the laughter frequency. 340

Figure 7 shows that subjects scoring *low* and *high* in Openness tend to initiate silences more frequently than expected to a statistically significant extent. To the best of our knowledge, the literature does not provide indications that can help to explain or interpret such an observation.

Overall, the results suggest that scenario and setting adopted for the Corpus induce low "*Relevance* (*i.e.*, the environment must allow the person to express the trait) and Availability (*i.e.*, the trait must be perceptible to others" (**Wright**, 2014). In other words, it appears that addressing the Winter Survival Task over the phone does not let the traits emerge with evidence, at least through the lens of the nonverbal cues analysed in this work.



Figure 8.

The bubble plot shows how the observed frequency of nonverbal cues deviates from the expected one according to the conflict handling style of individuals. When the bubble is red, the deviation is negative (less occurrences than expected), when the bubble is blue the deviation is positive (more occurrences than expected). The double asterisk means that the deviation is statistically significant with $\alpha = 0.01$ (according to a χ^2 test with Bonferroni Correction).

4.6 CONFLICT HANDLING STYLE EFFECTS

All subjects of the SSPNet Mobile Corpus have filled the Rahim Organizational Conflict Inventory 349 350 II (Rahim, 1983), a questionnaire aimed at measuring the attitude towards conflict and disagreement in 351 terms of five dimensions: *Compromising* (tendency to find solutions where all parties loose something to 352 reach mutually acceptable solutions), Avoiding (tendency to withdraw and sidestep rather than finding solutions), Obliging (tendency to focus on commonalities to satisfy the concerns of other parties), 353 Dominating (tendency to look for win-lose solutions), and Integrating (tendency to look for solutions 354 355 acceptable to all parties). The questionnaire allows one to calculate five integer scores that measure how 356 well an individual fits the tendencies associated to the five conflict handling styles above.

The scores range in the interval [-14, 14] and, for each trait, it is possible to define a variable V that has value *low* when the score is in the interval [-14, -5], *middle* when it is in [-4, 4], and *high* when it is in [5, 14]. This allows one to apply the methodology of Section 3. Figure 8 shows the deviations of the observed frequencies with respect to the expected ones (see Section 4.5 for more details on how deviations are calculated). The number of statistically significant deviations is higher than in the case of personality (see Figure 8). The probable reason is that the scenario adopted in the Corpus (**Polychroniou et al.**, 2014; **Vinciarelli et al.**, 2014) lets the Conflict Handling Style to emerge more clearly than the personality traits.

In the case of laughter, Figure 8 shows that people scoring low and middle along the Obliging style tend 365 to laugh significantly less than expected. To the best of our knowledge, the literature does not provide 366 indications that can explain such an observation. However, it has been shown that individuals that score 367 lower along the Obliging style tend to show less empathy and lower social skills (**Rahim et al.**, 2002). This 368 might result into lower tendency to laugh as well given the highly pro-social value of such a cue (Provine 369 and Yong, 1991; Provine, 1993). Similar explanations can apply to the tendency of subjects scoring high 370 along the Compromising dimension to laugh more than expected. In fact, people with such a style tend to 371 show concern for the others and, therefore, tend to adopt pro-social behaviors like laughter (Rahim et al., 372 373 2002).

374 Subjects that score *low* and *middle* along the avoiding style tend to display back-channel less frequently 375 than expected while those scoring *high* tend to display it more frequently than expected. A possible 376 explanation is that one of the main functions of back-channel is to encourage others to hold the floor and continue speaking (Hall and Knapp, 1992; Richmond et al., 1991). Such a type of behavior is
compatible with the tendencies associated to the avoiding style, namely to sidestep, to leave others to
address the problems, etc. In a similar vein, Figure 8 shows that people that score *middle* and *high* in
Avoiding tend to initiate silences more frequently than expected (the effect size is small, but statistically
significant). In this case as well, the cue appears to be compatible with the tendencies associated to the
conflict handling style.

The same tendency to initiate silence more frequently than expected can be observed for people scoring *middle* along the Integrating style. Given that the main tendency associated to this way of handling conflict is to find solutions acceptable to all parties, higher frequency of silence might correspond to the tendency to leave others talk and express their points of view (**Rahim et al.**, 2002). In this sense, the observations of Figure 8 seem to be compatible with the attitude the Integrating style accounts for.

In the case of overlapping speech, significant effects can be observed for all styles except Compromising. 388 389 This is not surprising because the cue has been extensively shown to be associated with conflict, both in human sciences (Schegloff, 2000; Smith-Lovin and Brody, 1989) and computing (Grezes et al., 2013; 390 Kim et al., 2014). The overall pattern of association (see Figure 8) suggests that the subjects that tend to 391 satisfy concerns for others tend to initiate overlapping less frequently than expected (subjects that score 392 *low* and *middle* in Dominating or *high* in Integrating and Obliging), while those that tend to privilege 393 394 concerns for the self tend to initiate overlapping more frequently than expected (*high* in Dominating or middle in Obliging and Integrating). Not surprisingly, such a pattern does not apply to Avoiding because 395 such a handling style accounts for attitudes that do not privilege neither concerns for the self or concerns 396 for the other. 397

5 CONCLUSIONS

The article has presented a detailed analysis of the temporal distribution of nonverbal cues (laughter, fillers, back-channel, silence and overlapping speech) in the SSPNet Mobile Corpus, a collection of 60 phone calls between unacquainted individuals (120 subjects in total). In particular, the analysis shows how the frequency changes according to six factors expected to account for the relational context, namely gender, role, topic of conversation, mode of interaction, personality and conflict handling style of the interactants.

404 The results show that the nonverbal cues do not distribute uniformly over time, but appear more or less frequently according to one or more of the abovementioned factors. In particular, male subjects and people 405 playing the role of receiver appear to display more frequently nonverbal cues associated to dominance 406 407 and/or higher social status. This happens even if the scenario adopted in the Corpus does not introduce any status or power difference between the two subjects involved in the same call. In the case of gender, 408 this is coherent with previous results showing that people tend to perceive male subjects as higher in status 409 (see Section 4.1). To the best of our knowledge, this is the first time that a similar effect is observed for 410 callers and receivers. 411

In the case of conflict handling style, nonverbal cues appear to change frequency according to the 412 413 tendencies associated to the various styles while, in the case of personality, statistically significant deviations with respect to the expected distributions take place only in a limited number of cases. This 414 is not surprising given that the scenario of the Corpus includes negotiation and disagreement aspects 415 that allow the conflict handling styles to emerge more clearly in terms of behavioral cues. This finds 416 417 confirmation in the changes observed when the relational context factor accounts for the mode of interaction (agreement vs disagreement). Finally, several cues change of frequency to a statistically 418 significant extent depending on whether the subjects are addressing the task at the core of the scenario or 419 not. In this case as well, the observations are compatible with previous work in the literature. 420

421 Overall, the findings suggest that the subjects manage to convey the same socially relevant information 422 as in face-to-face encounters even if they have to constrain their expressiveness through the phone. In other words, the lack of visual feedback is not an obstacle towards manifesting dominance, power differences and/or social verticality and reproduce, to a substantial extent, the patterns observed in the cases where the WST or other negotiation tasks are addressed in co-located settings. At the same time, the use of the phones appears to introduce at least one specific bias, i.e. the tendency to associate the role of receiver with behaviours typical of dominance and higher status. Such a finding might depend on the particular scenario adopted in the Corpus, but still shows that communication technologies can actually influence human-human communication and are not a mere passive channel.

430 From a technological point of view, the main interest of the findings above is that social and 431 psychological phenomena that cannot be observed and accessed directly can still be inferred from physical traces - the nonverbal cues and their frequency - that can be sensed and detected automatically. In this 432 respect, the analysis presented in this work provides a solid ground for domains like Social Signal 433 Processing (Vinciarelli et al., 2012), Computational Paralinguistics (Schuller and Batliner, 2013) or 434 Human Media Interaction (Nijholt, 2014) that aim at making machines socially intelligent, i.e. capable to 435 understand social interactions in the same terms as humans do. In particular, the observations suggest that 436 it is possible to develop automatic approaches for the inference of the factors the cues account for (e.g., 437 438 mode of interaction, topic of conversation, conflict handling style, etc.). However, while the inference of certain factors can be expected to achieve satisfactory performance because the number of statistically 439 significant effects is high (e.g., the conflict handling style), the inference of other factors might be difficult 440 or not possible because the corresponding physical traces are too weak (e.g., the personality traits), at 441 least for what concerns the cues analyzed in this work. 442

443 The development of the approaches above will contribute to further improve the state-of-the-art of conversational technologies (Renals et al., 2014). These include, e.g., the analysis of agent-customer 444 interactions at call centres¹ with the goal of improving the quality of services (Galanis et al., 2013), the 445 development of dialogue systems capable to interact naturally with human users (Keizer et al., 2014), the 446 improvement of tutoring systems aimed at supporting students in collective learning processes (Scherer 447 et al., 2012), the creation of speech synthesizers² that convey both verbal and nonverbal aspects of 448 a text (Schroeder, 2009), the enrichment of multimedia indexing systems with social and affective 449 450 information (Andre, 2013), etc.

In light of the above, the continuation of this work can take two parallel, but intertwined directions. 451 The first is the development of automatic approaches that perform the tasks mentioned above. The second 452 is the analysis of interaction effects between multiple factors. The findings described above focus on 453 individual factors because this makes it possible to observe larger number of events and, hence, to collect 454 more reliable statistics. However, the analysis of interaction effects can show further, more subtle effects 455 456 like, e.g., possible changes in the frequency of certain cues for subjects that have the same gender but different conflict handling styles. This, in turn, can help to further enhance the performance of automatic 457 458 systems that aim at inferring the relational factors from the frequency of the cues analyzed in this study.

ACKNOWLEDGEMENT

The authors are indebted with Dr Anna Polychroniou and Dr Hugues Salamin for their invaluable contribution to the SSPNet Mobile Corpus.

Funding: This work was possible thanks to the support of the European Commission (SSPNet, Grant
Agreement 231287), of the Swiss National Science Foundation (National Centre of Competence in
Research IM2) and the Finnish Ministry for Technological Innovation (TEKES).

This is a provisional file, not the final typeset article

 $^{^1}$ See http://www.cogitocorp.com for a company working on the analysis of call centre conversations.

 $^{^2~{\}rm See}~{\rm https://www.cereproc.com}$ for a company active in the field.

REFERENCES

- Anderson, K. and Leaper, C. (1998), Meta-analyses of gender effects on conversational interruption: Who,
 what, when, where, and how, *Sex Roles*, 39, 3-4, 225–252
- Andre, E. (2013), Exploiting unconscious user signals in multimodal human-computer interaction, ACM
 Transactions on Multimedia Computing, Communications, and Applications, 9, 1s, 48
- Bilmes, J. (1988), The concept of preference in conversation analysis, *Language in society*, 17, 02, 161– 181
- 470 Bonin, F., Campbell, N., and Vogel, C. (2014), Time for laughter, *Knowledge-Based Systems*, 71, 15–24
- 471 Clark, H. H. and Fox Tree, J. E. (2002), Using "uh" and "um" in spontaneous speaking, *Cognition*, 84, 1,
 472 73–111
- Galanis, D., Karabetsos, S., Koutsombogera, M., Papageorgiou, H., Esposito, A., and Riviello, M.-T.
 (2013), Classification of emotional speech units in call centre interactions, in Proceedings of IEEE
 International Conference on Cognitive Infocommunications, 403–406
- Grezes, F., Richards, J., and Rosenberg, A. (2013), Let me finish: Automatic conflict detection using
 speaker overlap, in Proceedings of Interspeech
- Hall, J. and Knapp, M. (1992), Nonverbal communication in human interaction (Harcourt Brace College
 Publishers)
- Hall, J. A., Coats, E., and Smith LeBeau, L. (2005), Nonverbal behavior and the vertical dimension of
 social relations: a meta-analysis., *Psychological bulletin*, 131, 6, 898–924
- Hecht, M., De Vito, J., and Guerrero, L. (1999), Perspectives on nonverbal communication: Codes,
 functions, and contexts, in The nonverbal communication reader (Waveland Press), 3–18
- Jaworski, A. (1999), The power of silence in communication, in L. Guerrero, J. De Vito, and M. Hecht,
 eds., The nonverbal communication reader (Waveland Press), 156–162
- Joshi, M., Davis, E., Kathuria, R., and Weidner, C. (2005), Experiential learning process: Exploring
 teaching and learning of strategic management framework through the winter survival exercise, *Journal of Management Education*, 29, 5, 672–695
- Keizer, S., Foster, M., Wang, Z., and Lemon, O. (2014), Machine learning for social multi-party human robot interaction, ACM Transactions on Intelligent Interactive Systems, 4, 3, 14:1–14:32
- Kim, S., Filippone, M., Valente, F., and Vinciarelli, A. (2014), Predicting continuous conflict perception
 with bayesian gaussian processes, *IEEE Transactions on Affective Computing*, 5, 2, 187–200
- Leffler, A., Gillespie, D., and Conaty, J. (1982), The effects of status differentiation on nonverbal behavior,
 Social Psychology Quarterly, 45, 3, 153–161
- 495 Martin, P. and Bateson, P. (2007), Measuring Behaviour (Cambridge University Press)
- McCrae, R. (2009), The Five-Factor Model of personality, in P. Corr and G. Matthews, eds., The
 Cambridge handbook of personality psychology (Cambridge University Press), 148–161
- McGinn, K. and Croson, R. (2004), What do communication media mean for negotiations? A question of
 social awareness, in M. Gelfand and J. Brett, eds., The handbook of negotiation and culture (Stanford
 University Press), 334–339
- Nakagawa, S. (2004), A farewell to Bonferroni: the problems of low statistical power and publication
 Behavioral Ecology, 15, 6, 1044–1045
- 503 Nijholt, A. (2014), Breaking fresh ground in human-media interaction research, Frontiers in ICT, 1, 4
- Oldmeadow, J. A., Platow, M. J., Foddy, M., and Anderson, D. (2003), Self-categorization, status, and
 social influence, *Social Psychology Quarterly*, 66, 2, 138–152
- Poggi, I., D'Errico, F., and Vincze, L. (2011), Agreement and its multimodal communication in debates:
 A qualitative analysis, *Cognitive Computation*, 3, 3, 466–479
- Polychroniou, A., Salamin, H., and Vinciarelli, A. (2014), The SSPNet Mobile Corpus: Social signal
 processing over mobile phones, in Proceedings Language Resources and Evaluation Conference, 1492–
 1498
- 511 Provine, R. (1993), Laughter punctuates speech: Linguistic, social and gender context of laughter,
 512 *Ethology*, 95, 4, 291–298
- 513 Provine, R. and Yong, Y. (1991), Laughter: A stereotyped human vocalization, *Ethology*, 89, 2, 115–124

- 514 Rahim, M. (1983), A measure of styles of handling interpersonal conflict, Academy of Management Journal, 26, 2, 368-376 515
- 516 Rahim, M., Psenicka, C., Polychroniou, P., Zhao, J., Yu, C., Chan, K., et al. (2002), A model of emotional intelligence and conflict management strategies: a study in seven countries, International Journal of 517 518 Organizational Analysis, 10, 4, 302–326
- Rammstedt, B. and John, O. (2007), Measuring personality in one minute or less: A 10-item short version 519 520 of the Big Five Inventory in English and German, Journal of Research in Personality, 41, 1, 203–212
- Renals, S., Carletta, J., Edwards, K., Bourlard, H., Garner, P., Popescu-Belis, A., et al. (2014), ROCKIT: 521 Roadmap for conversational interaction technologies, in Proceedings of the Workshop on Roadmapping 522 the Future of Multimodal Interaction Research Including Business Opportunities and Challenges, 39-523 42
- 524
- 525 Richmond, V., McCroskey, J., and Payne, S. (1991), Nonverbal behavior in interpersonal relations 526 (Prentice Hall)
- 527 Ruch, W. and Deckers, L. (1993), Do extraverts like to laugh?: An analysis of the Situational Humor Response Questionnaire (SHRQ), European Journal of Personality, 7, 4, 211–220 528
- Saucier, G. and Goldberg, L. (1996), The language of personality: Lexical perspectives on the five-factor 529 model, in J. Wiggins, ed., The Five-Factor Model of Personality (Guilford Press), 21-50 530
- Schegloff, E. (2000), Overlapping talk and the organization of turn-taking for conversation, Language in 531 society, 29, 01, 1-63 532
- Scherer, S., Weibel, N., Morency, L., and Oviatt, S. (2012), Multimodal prediction of expertise and 533 leadership in learning groups, in Proceedings of the International Workshop on Multimodal Learning 534 535 Analytics
- Schroeder, M. (2009), Expressive speech synthesis: Past, present, and possible futures, in Affective 536 537 Information Processing (Springer), 111–126
- Schuller, B. and Batliner, A. (2013), Computational paralinguistics: emotion, affect and personality in 538 speech and language processing (John Wiley & Sons) 539
- Smith-Lovin, L. and Brody, C. (1989), Interruptions in Group Discussions: The Effects of Gender and 540 541 Group Composition, American Sociological Review, 54, 3, 424–435
- Steele, C. and Aronson, J. (1995), Stereotype threat and the intellectual test performance of african 542 americans., Journal of personality and social psychology, 69, 5, 797-811 543
- Thompson, L., Wang, J., and Gunia, B. (2010), Negotiation, Annual review of psychology, 61, 491-515 544
- Vinciarelli, A., Pantic, M., Heylen, D., Pelachaud, C., Poggi, I., D'Errico, F., et al. (2012), Bridging 545 546 the Gap Between Social Animal and Unsocial Machine: A Survey of Social Signal Processing, IEEE Transactions on Affective Computing, 3, 1, 69–87 547
- Vinciarelli, A., Salamin, H., and Polychroniou, A. (2014), Negotiating over mobile phones: Calling or 548 549 being called can make the difference, *Cognitive Computation (to appear)*
- Ward, N. and Tsukahara, W. (2000), Prosodic features which cue back-channel responses in english and 550 japanese, Journal of Pragmatics, 32, 8, 1177-1207 551
- Wright, A. (2014), Current directions in personality science and the potential for advances through 552 553 computing, IEEE Transactions on Affective Computing, 5, 3, 292–296