1 **"Talk to You Later": Doing Social Robotics with Conversation Analysis. Towards**

2 **the Development of an Automatic System for the Prediction of Disengagement**

3 Nicolas Rollet

4 Télécom Paris, Institut polytechnique de Paris, I3 CNRS

5 Chloé Clavel

6 Télécom Paris, Institut polytechnique de Paris, Université Paris-Saclay

7

8 *Author Note*

9

10 Nicolas Rollet, Department of Economics and Social Sciences, Télécom Paris;

11 Chloé Clavel, LTCI, Télécom Paris, Université Paris-Saclay.

20 Correspondence concerning this article should be addressed to Nicolas Rollet

21 and Chloé Clavel, Télécom Paris, 19, place Marguerite Perey F-91120 Palaiseau,

22 France.

23 E-mail: nicolas.rollet@telecom-paris.fr; chloe.clavel@telecom-paris.fr

24                                              Abstract

25

26      This article presents an applied discussion of the possibility of integrating conversation

27      analysis (CA) methodology into that of machine learning. The aim is to improve the

28      detection of that which resembles disengagement in the interaction between a robot and

29      a human. We offer a novel analytical assemblage at the heart of the two disciplines, and

30      namely on the level of the annotation schemes provided by conversation analysis

31      transcription methods. First, we demonstrate that the need for a stable structure in

32      establishing an interaction scenario and in designing robot behaviours does not prevent

33      the emergence of ordinariness or creativity among the participants engaged in this

34      interaction. Secondly, based on an actual case, we emphasize the possibility of

35      systematicness in CA transcription to support the choice (a) of the categories targeted

36      by prediction methods and defined by the annotation scheme, and (b) of the verbal and

37      non-verbal features used to create prediction models.

38

39      Keywords: social robotics, engagement, machine learning, multimodal features,

40      annotation schemes, conversation analysis, transcription, closings

41 **"Talk to You Later": Doing Social Robotics with Conversation Analysis. Towards**

42 **the Development of an Automatic System for the Prediction of Disengagement**

43

44 In recent human-agent interaction studies, topics such as artificial agent's

45 sociality or engagement in interaction have given rise to a great deal of publications and

46 discussions. (Sidner, Lee, Kidd, & Rich, 2005; Dautenhahn, 2007; Pelachaud & Glas,

47 2015a; Clavel, Cafaro, Campano, & Pelachaud, 2016; Jones, 2017). With the

48 development of so-called social robotics, one observes a tension between *(i)* building

49 computational models using implementable traits to make an artificial agent sociable or

50 *socially interactive* (Fong, Nourbakhsh, & Dautenhahn, 2003; Breazeal, 2003); and *(ii)*

51 the recognition that this sociality is the product of a local organization, emerging from

52 the interaction (Suchman, 2007; Straub, 2016; Sabanovic, Chang, 2016). Such a

53 recognition suggests an increased focus on the mechanisms that organize an

54 environment in which interactions can take place, and where robot sociality and human

55 engagement can emerge. Such a recognition calls for interdisciplinarity.

56 As Pélachaud and Glas (2015a: 945) have stated, Sidner and Dzikovska (2002)

57 provide a definition of engagement, that is commonly used in the field of human-agent

58 interaction research, as a collaborative endeavour: "the process by which two (or more)

59 participants establish, maintain and end their perceived connection". For interactionists,

60 engagement does not necessarily involve verbal practices but rather any deployment of

61 orientation from one participant to another (or others) - in particular, the fact that

62 participants take into account the actions of others to produce, adjust their own

63 (Goffman, 1983). Engagement is a form of presence.

64    From a conversation analysis (CA) perspective, the question of engagement /

65    disengagement is linked to the observable deployment of interactional behaviours at the

66    scale of turns-at-talk and conversational sequences such as closings and pre-closings

67    (Sacks & Schegloff, 1973; Button, 1991), distinct orientations towards turn-taking

68    systems (Jefferson, 1978; Zimmerman, 2006), interruptions (Schegloff, 2002). In

69    addition, verbal behaviours can co-occur with other behaviours: gaze, body orientation,

70    etc. (Goodwin, 1981). Engagement can also refer to the multiple orientations of a

71    person who must respond to external events and reorganize, *in situ*, the practical

72    achievement of multi-activity, for example looking at his smartphone and driving

73    (Licoppe & Figeac, 2014), talking and playing an instrument (Rollet, 2010).

74    Human-robot (HRI) and Human-agent (HAI) interaction studies using CA as a

75    data exploration and analysis method are recent – about two decades (see for example

76    Cassel et al., 1999)[1]. Nevertheless, the combination of CA and coding practices in

77    affective computing / machine learning is quite uncommon. Although quite new,

78    applications are growing and CA is even recently presented in a new handbook on HRI

79    (Bartnek et al., 2019). Moreover, one can distinguish between studies that use CA more

80    as a thematic basis or conceptual underpinning (Sadazuka, Kuno, Kawashima, &

81    Yamazaki, 2007; Yu et al., 2013; Pelachaud & Glas, 2015b) and those that, in addition

82    to this thematic use, conduct a meticulous analysis of the interactions themselves, for

83    themselves, and therefore often contain detailed transcriptions in their publications

84    (Pitsch  et al., 2009; Dickerson, Robins, & Dautenhahn, 2013; Pelikan & Broth, 2016;

[1] There are earlier discussions in the Computer Science community, for example Chapman, D. (1992), that claims that "an interactionist computational interpretation of the conversation analytical rules is possible".

85   Rollet, Jain, Licoppe, & Devillers, 2017; Porcheron, Fischer, Reeves, & Sharples,

86   2018).

87          We pursue the idea that interdisciplinarity in social robotics offers novel and

88   ambitious opportunities for design (Fong, Nourbakhsh, & Dautenhahn, 2003; Bartnek et

89   al., 2019), especially if it focuses on objects at the heart of the respective disciplines. In

90   this sense, we can discuss this interdisciplinarity by considering the nature and purpose

91   of collaboration between a so-called interactionist sociological approach and a

92   computational model of human-robot interaction. Such interdisciplinarity raises

93   multiple questions with regard to both the design and detection of behaviours –

94   especially how to address these in terms of segmentable and annotatable flows of

95   interaction. We address some of these questions through the subject of disengagement

96   in a human-robot interaction. Drawing on a study conducted by (Ben-Youssef et al.,

97   2017), which develops a system to predict engagement breakdown in the context of

98   face-to-face interaction with the robot Pepper (Softbank robotics), we present a

99   conversation analysis viewpoint of the methodology adopted.

100  CA is a sociological approach that addresses language and especially talk-in-interaction

101  as a social organization. Its general topic lies in the description of the details of this

102  organization through which social interaction is made possible in an orderly and

103  intelligible way (Sacks, Schegloff, & Jefferson, 1974; Levinson, 1983; Sacks, 1992).

104  One recognizes here the affiliation to ethnomethodology's perspective which relies on

105  the intelligibility of the methods (defined as 'accountability') and on the participants'

106  point of view to produce its scientific analyses (Garfinkel, 1967). In addition, CA

107  considers an utterance in conversation in its *sequential conditions of emergence*, i.e., as

108  a contribution retrospectively and prospectively referring to the temporal stream of

109    interaction locally managed by participants. In that sense, social action is context-

110    shaped and context-renewing (Heritage, 1991) – property defined as *reflexivity*. The

111    conversational approach provides a methodological and argumentative framework in

112    which social interaction itself constitutes a powerful resource for analyzing and

113    understanding the meaning of being engaged, adapting, collaborating, disengaging, and

114    sharing an experience between co-participants.

115

116         The computational models of human-robot interaction discussed in this article

117    are based on "supervised" machine learning (Mohri, Rostamizadeh, & Talwalkar,

118    2012). The supervision consists of using audiovisual recordings of human-robot

119    interactions that have been annotated into categories of behaviours to predict (here,

120    engagement-related categories), in order to learn the models associated with each of

121    these categories. The methodology is broken down into different stages that structure

122    this article.

123         The first stage consists of collecting audiovisual recorded data that will be used

124    for learning behavioural models. This requires an interaction scenario to be defined, that

125    will be followed by the robot during its interaction with the participant. In the first

126    section, we present the chosen scenario and correlate it with the notion of context as

127    understood in conversation analysis, in particular with regard to the notion of *relevance*.

128         The second stage consists in annotating the data into engagement categories that

129    will be used for learning the supervised model. The second section of this article

130    presents a conversational perspective on the methodology for creating affective

131    computing annotation schemes, and shows how conversation analysis can contribute to

132    identifying features relevant to the development of our human behaviour detection

133　system. Specifically, this approach emphasizes the details constituting social

134　behaviours, on the scale of turns-at-talk, sub-units composing turns (Turn

135　Constructional Units), and sequences. Based on an actual case of human-robot

136　interaction, this section discusses the problems of categories and of segmentation, and

137　proposes leads for an interdisciplinary *assemblage*.

138　　　Finally, Section 3 offers a summary and an extension simultaneously aiming at

139　short-term applications in the context of a project underway, and lines of reflection that

140　expand the horizon of possible collaboration between interactionist sociology, machine

141　learning, and Affective Computing.

142　**1. Interaction Scenarios in Social Robotics and the Notion of Context in CA**

143　　　When seeking to develop methods to predict a participant's behaviour in his or

144　her interaction with a machine, it is essential to examine the situation or setting (*e.g.*

145　museum entrance hall (Campano, Clavel, & Pelachaud, 2015) or negociation game

146　(Langlet & Clavel, 2018) in which we want our prediction system to function— both to

147　define the interaction scenario used for data collection and to better understand the data

148　itself. Indeed, the participant's behaviours faced with the robot could depend heavily on

149　features of the situation of interaction.

150　　　The notion of context has been highly discussed and re-specified in conversation

151　analysis and sociolinguistics. We give a clarification below in order to better understand

152　the issues related to defining an interaction scenario and its impact on the type of

153　processing foreseen (Paragraph 1.a). We then give a contrastive view of CA and

154　affective computing approach regarding stability and emergence (Paragraph 1.b).

155 **a. Social robotics scenarios and the CA perspective**

156 The main goal of the interaction scenario that we have used, as defined in (Ben-

157 Youssef et al., 2017), is to collect the data on which the model to predict an engagement

158 breakdown will be built. Specifically, the goal is to collect data on a situation where a

159 participant is liable to exit the interaction before the end of the scenario. In this study,

160 the scenario defines the following aspects:

161 • The place of the interaction: the robot (in this case, Softbank's Pepper) is placed

162 in a hall where people frequently pass by.

163 • The mechanisms of entering into interaction: the robot starts speaking when it

164 detects the presence of a person, to invite him or her to interact, and the

165 participant is free to enter into the interaction or not.

166 • The mechanisms of exiting the interaction: the participant is free to exit the

167 interaction whenever he or she so wishes.

168 • The participant's engagement area: a space delimited by a semicircle with a

169 radius of 1.5 m.

170 • The phases of the scenario (note that the scenario was intentionally long, thus

171 including multiple phases, in order to trigger engagement breakdowns prior to

172 completion): the *welcome phase* (the robot introduces itself using very lively

173 animations, and gives instructions; the *dialog phase* (set of open questions that

174 the robot asks the participant about his or her tastes and personality); the

175 *cucumber phase* (with self-mockery and in the form of a game, the robot

176 presents its perceptive capabilities to show that, from its viewpoint, the

177 difference between a cucumber and a human is the human's face); and the *final*

178   *phase* (the robot concludes the interaction with questions intended to evaluate

179    the participant's interaction).

180

181   Now, one fundamental aspect of the *emic* perspective of conversation analysis

182 (that is, in which the orientation of the analysis is based on the participants' viewpoints

183 emerging in the interaction itself) regarding context, is characterized by *relevance*.

184   By following what participants make relevant themselves in the ongoing

185 interaction, it is possible to provide a characterization of the participants and of the

186 context. Such characterization provided *in situ* by the participants themselves is called

187 *internal setting relevancies* (Schegloff, 1987). In this sense, the selection of cues

188 relevant to the context from the multitude of available contextual elements, corresponds

189 to what is carried out visibly (accountably) by the interactants in the immediacy of the

190 interaction. These contextual aspects "internal to the interaction" heavily weigh on the

191 interaction, whilst they are not exclusive: background elements can also be important,

192 as well as the structure of the place, time, etc.[2]

193   The scenario consists of shaping a set of robot behaviours (of a finite number, by

194 definition) for interaction with humans. The robot's verbal and non-verbal behaviours

195 are defined in "interaction phases" (welcome phase, cucumber phase, etc.). From this

196 viewpoint, the scenario is designed asymmetrically: it consists of creating robot's

197 behaviours as ingredients of the different phases that make up a whole—the scenario—

198 but in which other ingredients of these phases, and namely human behaviours,

199 constitute hypothetical, fictional participation's opportunities.

---

[2] For a multi-dimensional definition of context, inspired by CA and linguistic anthropology, see (Duranti & Goodwin, 1992; Duranti, 1997)

200     With respect to the actions of the designer, this asymmetry consists in pre-

201     allocating turns-at-talk, which will then be experienced by a human participant.

202     Here are two examples of the same turn's occurrence:

203     (Extract 1)[3]

```
204   01    R     comment tu t'appelles ?
205               what's your name?
206   02    P     oui, Evelyne
207               yes, Evelyne
208

209         (Extract 2)
210   01    R     comment tu t'app@elles /
211               what's your name?
212   02    P                    @looks towards Robot's face
213   03          (1s)#(1,7s)
214   04    P        #leaves
215
```

216     In the two examples taken from the data, the robot (R) addresses a question

217     about the name (Line 1) to the participant that stands in front of it. Pepper produces

218     what CA calls a first part of an adjacency pair, namely a question-answer sequence

219     (Sacks & Schegloff, 1973; Schegloff, 2007). This is a very basic (ordinary) sequence in

220     which the first pair calls for possible seconds, that is: for the next speaker, the first pair

221     part is a context in which some relevant actions are expected, namely giving something

222     recognizable as an answer (i.e. a second pair part). That is what happens in the first

223     extract, but not in the second: the participant looks at the robot (L2) and after 1 second

224     pause, just leaves (L4). Contrary to what is expected in an interaction between two

---

[3] Some conventions of transcription are given below in 2.a.

225    humans, the strong pressure exerted by the first part of the pair on the possible second

226    part, in this situation, is visibly not addressed by the participant through his unilateral

227    disengagement. This particular observation raises an interesting point regarding the

228    categorization of the robot as a machine rather than as a social partner. For the latter

229    case, mitigation marks would generally be produced before leaving (Goffman, 1973;

230    Sacks & Schegloff, 1973).

231    Indicating that "the participant is free to exit the interaction" in the protocol illustrates

232    the rather logical asymmetry in the process of its assembling involving the design of

233    robot behaviours: the scenario is designed by projecting a hypothetical participant.

234    Whilst, in ordinary social interaction each apparently identical turn occurs in particular

235    circumstances, often as *another first time*. Hence, the naming of phases such as the

236    "welcome phase" or expectations such as "the participant is free to exit the interaction"

237    doesn't give any details on how this is factually, in a particular moment of the

238    interaction, both relevant and experienced[4]. Reversing the reasoning, we acknowledge

239    that the design of robot behaviours is based on the designer's ordinary interactional

240    knowledge: he/she assumes that these hypothetic interactions should start with a

241    welcome phase, and that saying hello will trigger a hello.

242          In addition, it is not given in advance that a framework externally considered

243    artificial or experimental implies that the participants treat it as such sometimes, or even

244    never. Context influences practices, but practices actualize and coproduce context as

245    well. In ethnomethodology, this refers to *reflexivity*: social practices pre-suppose

246    (context-shape) and constitute (context-renew) the framework of the interaction

---

[4] Although the interaction strategies defined within the phases of the scenario take into account the participant's responses/reactions, they are still part of a planning process: they do not predict the particular circumstances under which actions will occur (Suchman, 2007).

247 embedded within them (Heritage, 1991). There is nothing to prevent a particularly

248 constrained framework, such as a scenario based on question-answer games, from

249 seeing the emergence of unexpected, creative, natural behaviours.

250       **b.  Emergence, spontaneity, stability**

251       The emerging dimension (i.e. the situated character of actions) is at the heart of

252 the organization of social interactions (human-human). Conversation analysis'

253 analytical approach puts emphasis both on the *accountability* of actions (that which the

254 participants make visible themselves for themselves to coordinate their actions and

255 structure an activity) and on their normative aspect. Moreover, if the human treats the

256 robot like a conversation partner (and not an answering machine), then despite being

257 "scripted", the interaction will nonetheless adopt an emerging quality (fully on the

258 human side, and in the form of a tree diagram on the robot side). The fact that the

259 context is not only a set of imposed external characteristics but also a set of resources to

260 organize the interaction affords the participants the opportunity to demonstrate

261 creativity and spontaneity based on this framework. This question of spontaneity is

262 addressed from another viewpoint below (Section 3) regarding affective computing

263 annotation schemes in contrast with conversation analysis transcription practices. In

264 both cases, the problem of the reification of emerging behaviours is raised.

265       Now, the problem is that the contrast between the CA approach and affective

266 computing using supervised machine learning lies at the intersection of a problem of

267 stabilizing cues and the fundamentally emerging nature of social interaction. On the one

268 hand, on the CA side, we address the versatility of the context, or to put it another way,

269 the situated character of actions, as a strong resource to analyze the meaning of these

270     same actions. Actions are to be analyzed in their sequential deployment. (Where

271     versatility does not mean that there is no stability: it is punctual, circumstantial.). Few

272     data are processed generally, CA researchers work on singular cases and collections as

273     well but they are in relatively limited number. On the other hand, on the supervised

274     Machine Learning side, there is a technically justifiable need for stability of features

275     and identification of large classes that must be detectable and in limited numbers to

276     optimize the mass annotation work. Thus, from the viewpoint of a system to

277     automatically analyze participant behaviours, the question of spontaneity is addressed as

278     follows: how can defining constraints limiting the interaction steer the automatic

279     participant behaviour analysis system. Large amount of data is to be processed in this

280     approach.

281        To summarize, we gather the elements of this contrast in Figure 1.

282

| | Conversation Analysis | Supervised Machine Learning for Affective Computing |
|---|---|---|
| Actors | sociologists, anthropologists, linguists | linguists, computer scientists, annotators |
| Goals | analyze the intimacy of interaction, the practical reasonings, the improvised choreography, account for the intelligibility of actions, the sequential conducts | define classes that can be learned by a system, define cues for these classes, detect human behaviour, present a robust system |
| Steps in scientific production | audio-video recordings ; transcription ; text (analysis) | audio-video recordings ; annotation ; programming (model training) |
| Data Scale | small corpus (collections) | Big data |

283        Figure 1. Variety of shared and distinct aspects of conversation analysis and

284        affective computing approaches

285

286  The underlying idea is that by defining that which is potentially stable (the context-

287  shaped induced by the interaction scenario) and by integrating it into our behaviour

288  prediction models, we can improve the performance of prediction systems.

289  Our goal is to explore the possibility of improving the acuity of the detectable

290  cues, while guaranteeing some stability and not excluding the fact that this stability can

291  be temporary (on the scale of a turn or an adjacent pair, for example).


292  **2. Affective Computing Annotation of Recordings in Social Robotics vs. CA**

293  **Transcription**

294  In this section, we present a comparative study of two productions

295  independently obtained with the same recording data from interactions with the robot

296  Pepper: an affective computing annotation to develop behaviour prediction (Figure 2)

297  and a CA transcription (Transcript 1). The goal of this section is to provide an in-depth

298  comparison of what is produced by both approaches. The first paragraph contrasts the

299  categories produced by the affective computing annotations with the principle of

300  emergence in conversation analysis from a methodological point of view. The second

301  paragraph aims to compare productions of each approach, showing: *i)* similar results

302  that have emerged from both affective computing annotation and transcription processes

303  and *ii)* complementary results that show how both processes could benefit from each

304  other.  The third paragraph discusses the affective computing segmentation processes

305  from a conversation analysis perspective. We use as a guideline the framework

306  concerning the development of the system to predict participant engagement

307  breakdowns, and illustrate our study with examples of affective computing annotation
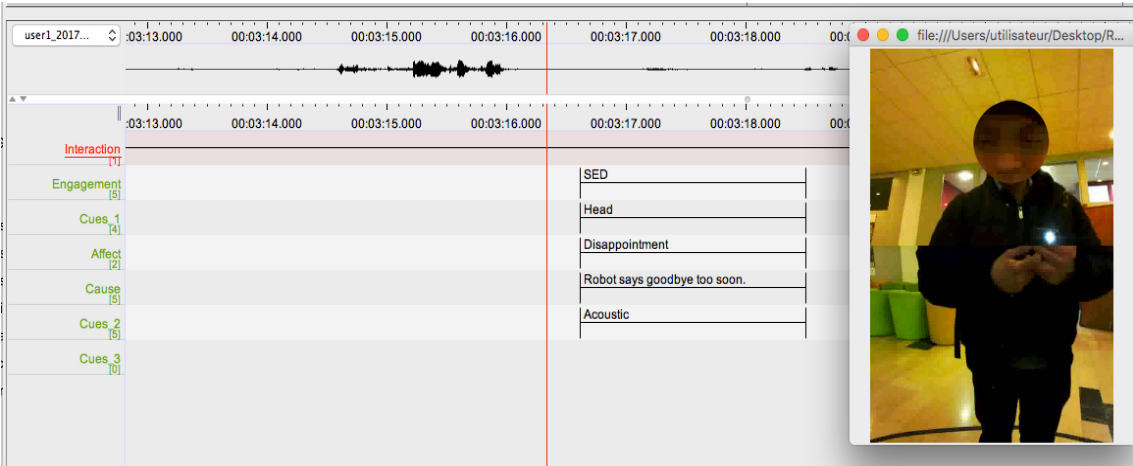
308  and CA transcription.

309

310  *Figure 2*. Affective computing annotation of interaction recordings conducted with the

311  Elan software[5] corresponding to lines 1-4 in the transcript 1 below.

312  Transcript 1 user1_2017-03-03

313  ```
01    P      ablas esp[agnol
```
314         *hablas español*

315  ```
02    R              [une autre fois\
```
316             *another time*

317  ```
03    P      <oh: ((look at smartphone)) (0,5s) > (0,5s) ok (..)je ne
```
318  ```
04           sais pas qué: qu'est-ce qué tou (1,1s) dire\
```
319         *oh, ok I don't know what, what do you..say*

320  *Transcript 1*[6]. In this excerpt from an interaction between the robot Pepper and a

321  human, turns are delimited to the left by a letter for each participant (R for robot and P

322  for human participant), and to the right by the end of the line or several lines below (e.g.

323  the first robot's turn starts in L02, while the first participant's turn starts in L01, then he

324  produces another turn that starts in L03 and ends in L04). Pauses are indicated in

325  seconds. The font used (courier) allows for ideal vertical alignment. Such a layout is

---

[5] https://tla.mpi.nl/tools/tla-tools/elan
[6] We use an adaptation of ICOR transcription conventions:
http://icar.cnrs.fr/projets/corinte/documents/2013_Conv_ICOR_250313.pdf

326    used to retain the outline of the course of the interaction over time by seeking to

327    reproduce the details of verbal behaviours, overlaps (marked by square brackets []),

328    intonations indicating the end of a turn (going up or down with the signs / or \), but also

329    bodily behaviours (in this case looking at the smartphone, indicated in double

330    parentheses in L03 with the signs "< >" that delineate the co-occurrence of this

331    behaviour with the verbal conduct). Intended to be as neutral as possible, such a

332    transcription is then subject to analysis that can lead to it being refined, for instance by

333    detailing the timing of the participant's orientation towards his smartphone with respect

334    to his turn in L03.

335    **a.   Categories vs. emergence principle: methodological comparison**

336    The goal of an affective computing annotation scheme is to define the macro-

337    categories that can be learned by the system. These categories must be sufficiently

338    represented in the data and relatively easy to annotate. The quality of the models learned

339    will depend heavily on the quality and quantity of the annotations obtained. All the

340    difficulty lies in defining an annotation protocol to establish convergence between the

341    annotations of multiple annotators, knowing that socio-emotional behaviours are highly

342    subjective phenomena that are difficult to define in an annotation (Cowie & Cornelius,

343    2003). The performances of the models learned are also evaluated using these

344    annotations as a reference, acknowledging that it is sometimes difficult to determine

345    what annotation is the most relevant between the automatic annotation of the system

346    and human annotation (Clavel & Callejas, 2016).

347        In the case of the study of disengagement, we conducted an annotation of the

348    different videos collected, an example of which is given in Figure 2. Four categories

349    were defined upstream to delimit the phenomena to annotate:

350       • BD (engagement BreakDown): phase when the participant leaves the

351         interaction;

352       • EBD (Early sign of future engagement BreakDown): the first precursor sign of

353         an engagement breakdown (necessarily results in an engagement breakdown,

354         and is therefore different from a SED and a TD);

355       • SED (Sign of Engagement Decrease); and

356       • TD (Temporary disengagement) (TD): phase during which the participant

357         interrupts the interaction before returning (this is a disengagement related to an

358         external interruption, such as a third person).

359        Modelled on the work in (Clavel, Vasilescu, & Devillers, 2011), a sub-

360    characterization of these macro-categories was also proposed and consisted of the

361    following tasks:

362     1. Defining the main verbal and non-verbal cues that characterize the annotated

363        phenomena (no sub-segmentation provided): speech, facial expressions, or

364        gestures (see *Cues1 field: Head* in Figure 2).

365     2. Specifying if emotions were expressed in these annotation segments (no sub-

366        segmentation provided) and identifying them in the following list of negative

367        emotions: frustration, boredom, nervousness, disappointment, anger, submission

368        (see *Affect Field: Disappointment* in Figure 2).

369     3. Providing an interpretation of the participant's disengagement (see *Cause field:*

370        *"Robot says goodbye to soon"* in Figure 2).

371  4.  Identifying secondary cues (see *Cues2 field: Acoustic* in Figure 2)

372  The objective of this subcategorization was to provide cues to understand how

373  the system functions and to interpret the reasons (explanatory cues) for which the

374  system detected the emergence of this category of phenomena. Indeed, when analyzing

375  the performance of the machine learning models of the marco-categories, the

376  subcategories can explain the behaviours of the system. For example, if the errors of the

377  automatic detection of an engagement breakdown are always located in segments where

378  boredom is expressed, it may think that the system missed to model this type of

379  expression of engagement breakdown.

380  The tool ELAN was used for these affective computing annotations (see Figure

381  2). It is an annotation tool for multimodal dialogue. This tool allows us to define our

382  own annotation scheme. For example, the segment annotated on Figure 2 is constructed

383  as a so-called "parent" category (SED) followed by associated cues or comments. In

384  this case, there is a bodily cue called Cue 1 ("head"), a non-lexical cue called Cue 2

385  ("acoustic"), an emotional cue ("disappointment") and, last of all, a Cause comment

386  ("Robot says goodbye too soon"). Note here that the choice of these categories is guided

387  by the task, that is, by what the system must detect. Affective computing annotation can

388  be considered top-down given that the categories are what guide the annotation of

389  explanatory cues—which greatly contrasts with the transcription mentality of

390  conversation analysis.

391  In conversation analysis transcription is a textual translation of repeated

392  observations from audio-visual data. In this sense, it constitutes a particular

393  configuration of the recorded reality, which is itself simply a particular configuration of

394  overall reality. From a methodological point of view, it is a research support that is not

395    sufficient in itself: the analysis is always carried out, transcription at hand, with repeated

396    visualizations of the corresponding audiovisual data. Transcription is both a necessary

397    and a reifying tool: it leads to decision-making and materializes a graphic layout (Ochs,

398    1979; Mondada, 2008), as we can see below in Figures 3 and 4, which show such

399    examples of transcription.

400        To summarize:

401      • Each type of transcription corresponds to a position and a goal for the researcher

402        or the transcriber.

403      • Each type of transcription corresponds to a status ascribed to the verbal and non-

404        verbal, and to the relationship maintained between the two.

405      • Transcription involves theoretical assumptions on the part of the transcriber,

406        which have a configuring effect on this transcription.

407

408        Transcription and *a fortiori* affective computing annotation result in de-

409    contextualization, that is, extraction from the singular context of production and

410    transformations (for example, certain phenomena that appear anecdotal in the field

411    become worthy of interest and fixed during transcription / affective computing

412    annotation practices).

413        Nevertheless, the status of transcriptions in conversation analysis is radically

414    different from that of the computational approach's annotation diagrams—even though

415    both "describe" an undertaking to categorize interactional behaviours, whether by the

416    conversation analysis researcher or multiple affective computing annotators. In affective

417    computing annotation schemes, macro-categories (that can be learned by the system)

418    are associated with cues (which are relatively limited), and the practical reasoning used

419    by annotators is to some extent invisible. By contrast, conversation analysis

420    transcription is an intermediary that attempts to be as neutral as possible between the

421    raw data and the researcher's analysis. This tendency of relative neutrality in the

422    production of transcripts refers to an orientation of the researcher towards the analysis

423    of participants' practical reasoning in the here and now of their social interactions. Note

424    that CA community is familiar with ELAN interface as a mean to visualize and account

425    for multimodal phenomena during collective work processes ('data sessions') and as

426    screenshots for publications (Mondada, 2006, 2008).

427        Nonetheless, conversation analysis is not fundamentally prevented from seeking

428    out systematicity (Sacks, Schegloff, & Jefferson, 1974; Stivers, 2015). Even though

429    debate exists within the CA community, it is not unreasonable to want to improve or

430    challenge the macro-categories resulting from the description modes of the

431    computational approach. In this approach, using the analysis of micro-phenomena can

432    feed a diversity of cues associated with these macro-categories in annotation schemes.

433    **b.  Categories vs. emergence principles: comparison of two types of**

434        **production**

435        To demonstrate this, let us contrast the affective computing annotation presented

436    in Figure 2 with the transcription of the same excerpt, taken from a common corpus of

437    interactions between the robot Pepper and humans. Note that both processes have been

438    carried out independently. This affective computing annotation is characterized by the

439    annotated segment as a "sign of engagement decrease" (SED). The complete

440    transcription of this segment and of what goes after is presented in Transcript 2.

441    Transcript 2 user1_2017-03-03

442   01   P      `ablas esp[agnol`

443                 *hablas español*

444   02   R            `[une autre fois\`

445                     *another time*

446   03   P      `<oh: ((look at smartphone)) (0,5s) > (0,5s) ok (..)je ne`

447   04         `sais pas qué: qu'est-ce qué tou (1,1s) dire\`

448                 *oh, ok I don't know what, what do you..say*

449   05         `(3,6s)`

450   06   P      `<((with greeting gesture)) au revoir/>`

451                             *goodbye*

452   07         `(6,2s)`

453   08   P      `<((with greeting gesture and body torq)) au revoir Pepper\>`

454                          *goodbye Pepper*

455   *Transcript 2*. Transcription associated with Figure 3 and 4

456        In order to visualize the correspondence that could be found between the two

457   productions, we manually integrated the segments of transcriptions corresponding to the

458   annotation environment in Figure 3 and 4. These figures show that in CA, analyzing the

459   excerpt reveals much more detailed information than that annotated, and may provide

460   explanatory indications of signs of disengagement: CA approach could give some hints

461   about how an identified closing (Figure 3) could be analyzed by scrutinizing also what

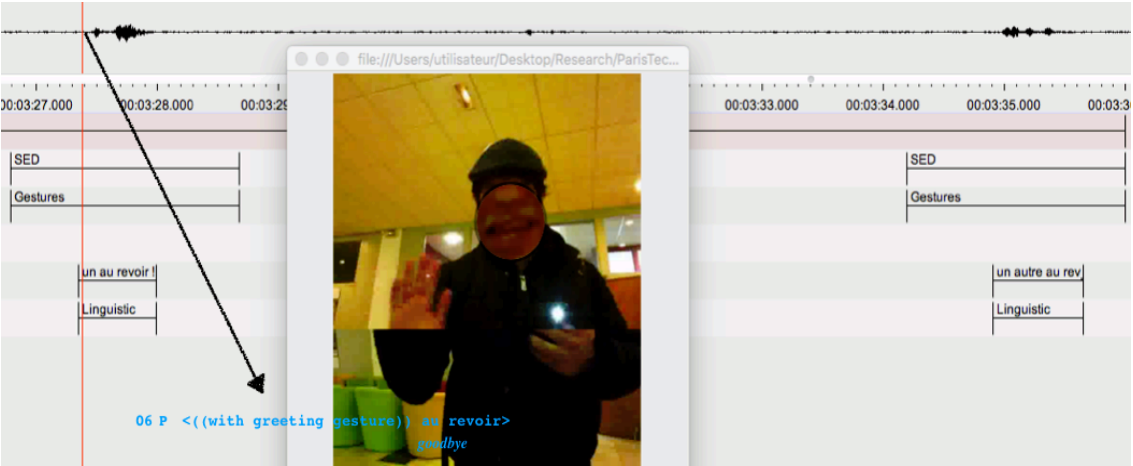462   happens right before in the interaction (Figure 4).

463

464    *Figure 3*. Assemblage of CA transcription – affective computing annotation



465

466    *Figure 4*. Assemblage of CA transcription – affective computing annotation (continued)

467            A first example of interesting details provided by CA transcription is given in

468    Figure 3. Within the segment annotated signs of engagement decreased (SED), the

469    affective computing annotation indicates globally that it is characterized by linguistic

22

470    ("un au revoir") and gestural cues, while the CA transcription details the type and the

471    exact timing of gesture that accompanied the "au revoir".

472         A second example is given in Figure 4. The affective computing annotation

473    indicates a decrease in engagement by using the annotation of the SED category, and

474    specifies that it consists of acoustic cues and cues related to head movements (Head).

475    Here, the CA transcription reveals different cues with their interpretation allowing one

476    to anticipate the engagement breakdown that are not given by the affective computing

477    annotation: *i)* in  (L01-02), the CA transcription allows us to identify a potential cause

478    of participant's engagement decrease that occurs before the SED affective computing

479    segment: an overlap between participant turns (the robot (R) produces a turn that causes

480    an overlap (L01-02)) and a violation of adjacency principle by the robot (the

481    participant's turn is a question addressed at R concerning a language skill but R

482    produces a turn that is topically inconsistent with P's question) ; *ii)* the "oh" produced

483    quietly and transcribed in the CA transcription  (L03) follows this overlap and marks

484    thus a reaction to this transgression. This cue occurs during SED segment but is just

485    signaled as acoustic cues by affective computing annotation; *iii)* the head cue annotated

486    in the affective computing annotation in Figure 4 is more precise in CA transcript. It is

487    detailed as an accompaniment of the verbal behaviour "oh" with the orientation of P

488    towards his smartphone (L03); *iv)* the CA transcription indicates in L03-L04 linguistic

489    cues denoting the re-engagement of the participant in the interaction. Reengagement is

490    not so far included in the affective computing annotation categories. In CA

491    transcription, we note that the "ok" seems to mark a re-engagement, a *springboard*

492    (Beach, 1993; Rollet, 2013) towards a new orientation: that of moving towards the end

493    of the interaction. This orientation is made visible by the production of a unit ("je ne

494    sais pas qué: qu'est-ce qué tou (1,1s) dire\") that topicalizes an unresolvable non-

495    understanding; *v)* the pre-closing phenomenon, that is, an interactional behaviour that

496    sequentially precedes and serves to project the closing as such (Sacks & Schegloff,

497    1973). This phenomenon is illustrated here by the re-engagement cues to move towards

498    the interactions described above. In this case, an "interactional blank cheque" is being

499    given by P, which is followed by a 3.6s slot granted to R which is therefore transcribed

500    on a separate line (L05).

501

502        Another interesting aspect of comparing the two forms of notation of this extract

503    concerns the relationship between turns L02, L03 and L04. An initial analysis offered

504    by the transcription is that in L03-04, P is marking disengagement according to two

505    mechanisms (gestural—he looks at his cell phone—and verbal). And, still according to

506    this analysis, the robot's turn L02 can be described as a transgression of the "one

507    speaker at a time" rule (Sacks, Schegloff, & Jefferson,1974), made visible by the

508    reaction that this transgression provokes (L03 "Oh"). In other words, in this case, a

509    transgressive value is ascribed to the robot's turn L02, and P is attributed the initiative

510    to move in two steps towards the end of the interaction.

511        However, an alternative to this analysis is possible, and can be derived from the

512    affective computing annotation itself. Specifically, in the annotation, there is a comment

513    associated with SED under the "cause" section: *Robot says goodbye too soon*. The first

514    comment we can make is that such a "cause" is consistent with categorizing a segment

515    as a "decrease in engagement". However, we can go even further. This comment is an

516    analysis of an entirely different level than, for example, the Cue1 section with "head".

517    This is a categorial and sequential analysis that has a significant influence on the

518  subsequent interpretation of P's behaviour. The fact that the annotator comments on R's

519  turn, "another time" as "Robot says goodbye too soon", and because this is not a

520  "goodbye" in the strict sense, shows that he considers this turn to be a behaviour

521  projecting a closing, such as a "goodbye": this is precisely the work accomplished by a

522  pre-closing, or a *junction*[7] in general (Button, 1991). If this is the case, P's "oh" may be

523  an indication of disappointment due not to technical incompetence revealed by the

524  emergence of an overlap (with a thematically incongruent turn), but to the participant

525  analyzing Pepper's turn as a pre-closing. Hence, P's turn, and in particular the unit "je

526  ne sais pas qué: qu'est-ce qué tou (1,1s) dire\", could be analyzed as an alignment

527  (constructed through topicalizing a non-understanding) with the end of the interaction,

528  initiated by the robot: something along the lines of "I guess we don't understand each

529  other". Following this analysis, the disengagement is therefore not initiated by P but

530  rather, from P's viewpoint, by R. Contrary to a disengagement, it is rather an *affiliation*

531  (Stivers, 2008) of the participant towards the disengagement of the robot.

532      These two comparative analyses show to what extent the interpretation of a

533  behaviour - even that of a robot - is not as univocal as it may seem, as soon as it is

534  examined in its interactional framework. Moreover, regardless of the interpretation

535  prioritized, a central and well-described sequence in conversation analysis literature

536  emerges as an essential phenomenon for analyzing disengagement: pre-closing.

537      Cues, such as verbal (the exclamation "oh") and non-verbal cues (posture, gaze)

538  and conversation analysis' interpretation of them also facilitate our understanding of

---

[7] The *junction* refers to a conversational pivot action that switches from the current topic to the closing. It is done in the current topic. There are several forms of action that put interaction on a closing track, e.g. projecting future activities, announcement of departure, or formulating summaries.

539    engagement breakdowns and can subsequently be integrated into the annotation scheme

540    or automatically annotated in order to be quantified. If they are sufficiently frequent and

541    representative, they can be added to the features extracted from videos for the machine

542    learning of the system to detect engagement breakdowns.
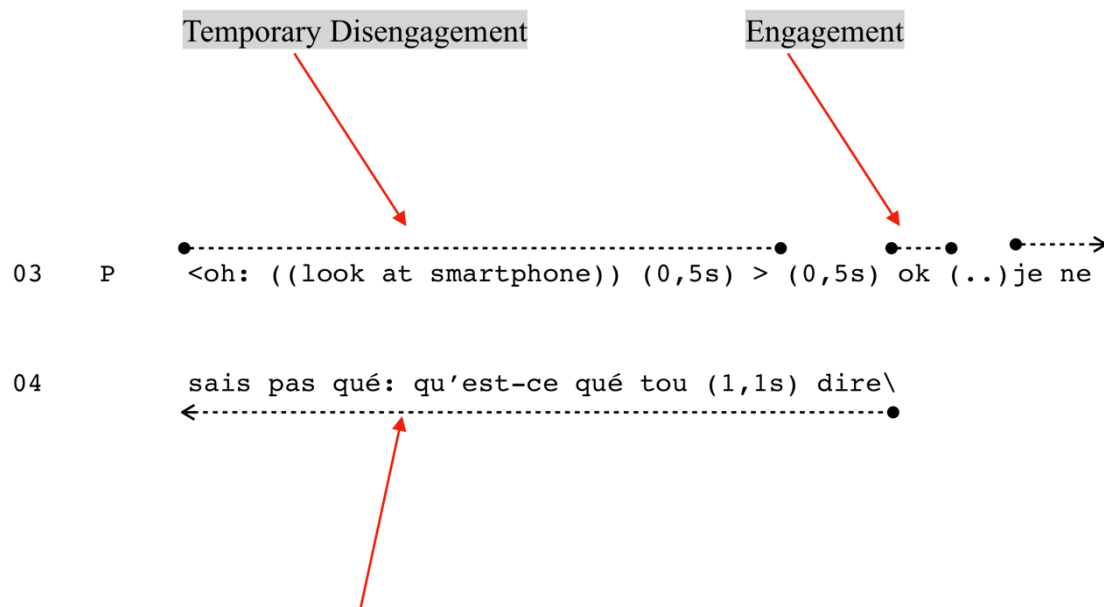
543    **c. Segmentation**

544    The allocation of categories by the annotation schemes used in social robotics

545    discussed above is generally based on a prior stage of segmenting audio/video feeds that

546    allows one to delimit annotation segments. In this case, we present a conversation

547    analysis viewpoint of the segments and phenomena thus delimited.

548    Annotating consists of two main stages. In the first stage the segments

549    associated with the above-mentioned categories are identified according to the

550    following steps:

551    1. Detecting the occurrences of one of these categories of phenomena (TD, EBD,

552        etc., see 2.a).

553    2. Segmenting those phenomena over time by defining their time boundaries

554        (segments that appear in the form of rectangles in Figures 2 and 5). Attributing

555        the category identified to these time segments. These annotation time segments

556        are called annotation units.

557    We believe to be crucial, first, discussing the segmentation of the interaction

558    flow and, second, the categorization of the segments as well as the analysis level that

559    they underpin. If we once again consider the comparative analysis of affective

560    computing annotation vs. transcription presented above, a second conclusion that

561    emerges is that annotation segments are "too macro". Concretely, P's turn in L03-04

562    can be successively analyzed as an indication of disengagement, with the "oh + looks at

563    cell phone", and a re-engagement initiated by "ok". Specifically (Figure 5):

Temporary Disengagement                    Engagement

03    P    <oh: ((look at smartphone)) (0,5s) > (0,5s) ok (..)je ne

04          sais pas qué: qu'est-ce qué tou (1,1s) dire\

Early sign of future engagement BreakDown

564

565    *Figure 5*. Assemblage diagram with CA transcription (from Transcript 2) and affective

566    annotation (from Figure 4)

567

568          Such a breakdown suggests a methodological viewpoint already mentioned

569    above: transcription attempts to provide analysis with the means of rendering

570    (accounting for) participants' ways of doing. This is particularly true for segmentation.

571    Even though analysis requires one to extract and isolate interaction segments, this task

572    can run the risk of losing its "natural" explanatory basis by becoming de-

573    contextualized—with researchers in that case running behind something that was

574    nonetheless already there. For if we, the second-hand observers, manage to follow a

575    conversation or to understand a fragment of an interaction, it is because an initial

576    analysis undertaking *in praesentia* took place.

577       By reconsidering the example in Figure 5, we note that the work accomplished

578    by the participant through what he says and does can be described in three stages. Three

579    stages for a turn; three *turn-constructional units* (Sacks, Schegloff, & Jefferson, 1974;

580    Ford & Thompson, 1996) which each accomplish something different, as we describe

581    above (Paragraph 2.a). A turn is a participation unit. It can be composed of multiple

582    sub-units which are often delimited in terms of syntax, intonation, semantics, or

583    pragmatics—and even gesturally or rhythmically. A turn is a contribution to the

584    progressiveness of a situated interaction, and an initial signifying segmentation

585    undertaking is conducted by the interactants themselves in the here and now of their

586    social activities.

587       From a machine learning viewpoint, the question of the analysis unit is

588    fundamental: what time frames of analysis should be used to extract social signals, and

589    which units should be considered a frame for decision-making for the phenomenon to

590    predict? Conversation analysis contributes to understand how people themselves break

591    down the stream of language, and can help provide indications around the choice of

592    analysis unit or the decision frame.

593    **3. Conclusion and Prospects**

594       In this article we have presented an affective computing annotation protocol

595    dedicated to a project to detect disengagement in human-robot interactions. In the frame

596    of an interdisciplinary collaboration between machine learning and conversation

597    analysis within social robotics, this critical and constructive viewpoint can be applied to

598    the analysis of:

599 • Context and relevancy: combining a multidimensional viewpoint with that of the

600     interacting participants, affords a perspective that encompasses the situated

601     nature of social behaviours.

602 • Issues surrounding annotation and the segmenting of interaction flows: as stated

603     in Paragraph 2a., breaking down and categorizing the behaviours appearing in

604     interactions are not meaningless practices: they relate to forms of representation

605     with respect to the relationship between the verbal and the non-verbal, as well as

606     regarding the level of exogenous ('etic') production of meaning.

607 • The question of sequentiality as a new "explanatory feature": this is a dimension

608     that is all too often neglected but can nonetheless constitute a fundamental

609     resource in the endogenous production of meaning, in the same way as the

610     syntactic, semantic, melodic, and pragmatic levels. The identification of

611     phenomena such as pre-closings and junctions as cues of disengagement not

612     only concerns a set of typical actions ("I gotta go", "talk to you later", etc.), but

613     also a space of interaction that highlights the significant relationship between

614     "what has just happened" and "what could happen next".

615     An initial line of collaboration between the two disciplines is based on the idea

616 that interaction is more fluid through robotic behaviours that tend towards a form of

617 ordinariness (that is, practices that are recognizable as being able to appear in an

618 ordinary, daily, and routine social interaction (Sacks, 1984)). In other words,

619 collaboration in creating scenarios can consist in providing designers with the viewpoint

620 of a competent participant of everyday life who has developed a reflexive perspective

621 with respect to his or her own (ordinary) practices, which is then refined through

622   detailed observations of diverse social interactions. The following interaction cues serve

623   as examples:

- the acoustic forms of a robot's utterances which project an action in a sequential
  process;

- the construction of turns as pragmatic units that are not necessarily primarily
  based on syntactic or semantic considerations; and

- the orientation of actions from a sequential viewpoint, that is, in a logic of turn-
  taking system and establishing interactional episodes or activities.

630   As another prospect for collaboration, attention can be drawn to the explanatory

631   cues of an annotation category. In the scheme presented above (Section 1), a number of

632   cues ranging from the human-robot distance to acoustic features and spatial orientation

633   are highlighted. Moreover, analyzing the interaction excerpt between the robot Pepper

634   and a human (cf. Paragraph 2.a) reveals a fundamental feature in explaining

635   disengagement in social interaction in general. Specifically, beyond the pre-closing

636   phenomenon as such, the sequential dimension appears to be central and the machine

637   learning chosen must be able to integrate this sequentiality. To analyze how a robot's

638   turn (such as "another time") is treated by the participant ("je ne sais pas qu'est-ce qué

639   tou dire"), semantic, acoustic and pragmatic contents are not enough: the sequential

640   positioning highly contributes on the grasp of what takes place, what follows, and what

641   the participant makes accountable. It thus becomes necessary to establish a way of

642   categorizing / codifying this sequential dimension each time that the annotator, aware of

643   this dimension, observes its consequential nature. This awareness means considering a

644   certain teaching process intended for annotators, and questioning how far this can be

645   taken. The first step would be to sensitize annotators on pre-closings and conclusive

646     junctions as familiar phenomena they do experience in their ordinary life even if they've

647     never 'conceptualized' it – a step we've been just started to test. Affective computing

648     annotation work is already moving in this direction, consisting of the canonical

649     sequential format in interaction, and namely adjacent pairs (Langlet & Clavel, 2014),

650     such as the question-answer sequence. Rather than leaving this field completely open to

651     the annotator, the idea here is to enrich the explanatory cues of a macro-category (TD,

652     EBD, BD, SED) by implementing the "Cause" section of the best-demarcated sub-cues:

653     first comes the question of sequentiality, but we can also consider the pragmatic or even

654     topical dimension. In this sense, the precision regarding the annotation segment

655     addressed in the Cause section becomes crucial. These sub-cues can either be used as

656     subcategories to be predicted by machine learning, or for the design of the input features

657     of machine learning in order to improve disengagement prediction models.

658         These prospects raise the question of the extent to which the phenomena

659     observed in the data can be sufficiently formalized to be processed for machine learning

660     in order to improve machines' ability in detecting these phenomena. Moreover, to what

661     extent is the tension between the principle of describing the uniqueness of cases -

662     defining the analytical mentality of Conversation Analysis - and the requirement of

663     generalization for the training of automatic models bearable? That is to say, how closely

664     can we model the uniqueness and emerging nature of social interaction?

665

666                                                       References

667     Bartneck, C., Belpaeme, T., Eyssel, F., Kanda, T., Keijsers, M., & Šabanović, S. (2019).

668             *Human-robot interaction: An introduction*. Cambridge: Cambridge University

669             Press.

670     Beach, W. A. (1993). Transitional regularities for 'casual' "Okay" usages. *Journal of*

671             *Pragmatics, 19*, 325-352.

672     Ben-Youssef, A., Clavel, C., Essid, S., Bilac, M., Chamoux, M., & Lim, A. (2017). *UE-*

673             *HRI: a new dataset for the study of user engagement in spontaneous human-*

674             *robot interactions*. Paper presented at the 19th ACM International Conference

675             on Multimodal Interaction, Glasgow, UK.

676      Button, G. (1991). Conversation-in-a-series. In D. Boden & D. H. Zimmerman (Eds.),

677             *Talk and social structure*. *Studies in Ethnomethodology and Conversation*

678             *Analysis*  (pp. 251-277). Cambridge: Polity Press.

679     Campano, S., Clavel, C., & Pélachaud, C. (2015). *"I like this painting too": when an*

680             *ECA shares appreciations to engage users*. Paper presented at the 14th

681             International Conference on Autonomous Agents and Multiagent Systems

682             (AAMAS'15), Istanbul, Turkey.

683      Cassell, J., Torres, O., & Prevost, S. (1999). Turn taking vs. discourse structure: How

684             best to model multimodal conversation. In Y. Wilks (Ed.), *Machine*

685             *Conversations* (pp.143–154). The Hague: Kluwer.

686     Chapman, D. (1992). Computer rules, conversational rules. *Computational Linguistics, 18*(4),

687             531-536.

688     Clavel, C., Vasilescu, I., & Devillers, L. (2011). Fiction support for realistic portrayals

689    of fear-type emotional manifestations. *Computer Speech & Language*, *25(1)*, 63-

690    83.

691    Clavel, C., Cafaro, A., Campano, S., & Pelachaud, C. (2016). Fostering User

692    Engagement in Face-to-Face Human-Agent Interactions: A Survey. In A.

693    Esposito & L. C. Jain (Eds.), *Toward Robotic Socially Believable Behaving*

694    *Systems – Volume II: Modeling Social Signals* (pp. 93-120). Cham: Springer

695    International Publishing.

696    Clavel, C., & Callejas, Z. (2016). Sentiment analysis: from opinion mining to human-

697    agent interaction. *IEEE Transactions on affective computing*, *7(1)*, 74-93.

698    Cowie, R., & Cornelius, R. R. (2003). Describing the emotional states that are

699    expressed in speech. *Speech communication*, *40(1-2)*, 5-32.

700    Dautenhahn, K. (2007). Socially intelligent robots: dimensions of human-robot

701    interaction. *Philosophical transactions of the Royal Society of London. Series B,*

702    *Biological sciences, 362*(1480), 679-704.

703    Dickerson, P., Robins, B., & Dautenhahn, K. (2013). Where the action is: A

704    conversation analytic perspective on interaction between a humanoid robot, a

705    co-present adult and a child with an ASD. *Interaction Studies, 14*(2), 297-316.

706    Duranti, A. (1997). *Linguistic anthropology*. New York: Cambridge University Press.

707    Duranti, A., & Goodwin, C. (1992). Rethinking context: an introduction. In A. Duranti

708    & C. Goodwin (Eds.), *Rethinking Context, language as an interactive*

709    *phenomenon* (pp. 1-42). Cambridge University Press.

710    Fong, T., Nourbakhsh, I., & Dautenhahn, K. (2003). A survey of socially interactive

711    robots. *Robotics and Autonomous Systems, 42(3-4)*, 143-166.

712    Garfinkel, H. (1967). *Studies in ethnomethodology*. Engelwood Cliffs: Prentice-Hall.

713    Goffman, E. (1973). *La mise en scène de la vie quotidienne* (Vol. 2. Les relations en

714         public). Paris: Les Editions de minuit

715    Goffman, E. (1983). The interaction order. *American Sociological Review, 48(1)*, pp: 1-

716         17.

717    Goodwin, C. (1981). *Conversational organization: interaction between speakers and*

718         *hearers*. New York: Academic Press.

719    Heritage, J. (1991). L'Ethnométhodologie : une approche procédurale de l'action et de

720         la communication. *Réseaux CNET, 50, 89-130*.

721    Jefferson, G. (1978). Sequential aspects of storytelling in conversation. In J. Schenkein

722         (Ed.), *Studies in the organization of conversational interaction* (pp. 219-248).

723         New York: New York Academic Press.

724    Jones, R. A. (2017). What makes a robot "social"? *Social Studies of Science, 47*(4),

725         556-579.

726    Langlet, C., & Clavel, C. (2014). *Modelling user's attitudinal reactions to the agent*

727         *utterances: focus on the verbal content*. Paper presented at the 5th International

728         Workshop on Corpora for Research on Emotion, Sentiment & Social Signals

729         (ES3 2014), Reykjavik, Iceland.

730    Langlet, C., & Clavel, C. (2018). *Detecting User's Likes and Dislikes for a Virtual*

731         *Negotiating Agent*. Paper presented at the 20th ACM International Conference

732         on Multimodal Interaction, Boulder, USA.

733    Levinson, S. C. (1983). *Pragmatics*. Cambridge University Press.

734    Licoppe, C., & Figeac, J. (2014). L'Organisation temporelle des engagements visuels

735         dans des situations de multi-activité équipée en milieu urbain. *Activités, 11*(1).

736    Mondada, L. (2006). Participants' online analysis and multimodal practices: projecting

737      the end of the turn and the closing of the sequence. *Discourse Studies, 8*(1), pp.

738      117-129.

739    Mondada, L. (2008). Documenter l'articulation des ressources multimodales dans le

740      temps : la transcription d'enregistrements vidéos d'interactions. In M. Bilger

741      (Ed.), *Données orales. Les enjeux de la transcription* (Vol. 37), pp. 127-156.

742      Perpignan: Presses Universitaires de Perpignan.

743    Mohri, M., Rostamizadeh, A., & Talwalkar, A. (2012). *Foundations of machine*

744      *learning*. MIT press.

745    Ochs, E. (1979). Transcription as a theory. In B. B. Schieffelin (Ed.), *Developmental*

746      *Pragmatics* (pp. 42-72). Academic Press.

747    Pelachaud, C., & Glas, N. (2015a). *Definitions of Engagement in Human-Agent*

748      *Interaction*. Paper presented at the International Workshop on Engagement in

749      Human Computer Interaction (ENHANCE), Xi'an, China.

750    Pelachaud, C., & Glas, N. (2015b). *Topic transition strategies for an information-giving*

751      *agent*. Paper presented at the the 15th European Workshop on natural Language

752      Generation, Brighton, UK.

753    Pelikan, H. R. M., & Broth, M. (2016). *Why That Nao?: How Humans Adapt to a*

754      *Conventional Humanoid Robot in Taking Turns-at-Talk*. Paper presented at the

755      2016 CHI Conference on Human Factors in Computing Systems, San Jose,

756      California, USA.

757    Pitsch, K., Kuzuoka, H., Suzuki, Y., Süssenbach, L., Luff, P., & Heath, C. (2009). *"The*

758      *first five seconds": Contingent stepwise entry into an interaction as a means to*

759      *secure sustained engagement in Human-Robot-Interaction*. Paper presented at

760          the IEEE International Symposium on Robot and Human Interactive

761          Communication ROMAN 2009, Toyama, Japan.

762   Porcheron, M., Fischer, J. E., Reeves, S., & Sharples, S. (2018). *Voice Interfaces in*

763          *Everyday Life*. Montreal QC, Canada: Association for Computing Machinery.

764   Rollet, N. (2010). *All the things you are*. Activité multimodale, frontières et musiques

765          improvisées en répétition In N. Andrieux-Reix (Ed.), *Frontières. Du linguistique*

766          *au sémiotique* (pp. 279-302). Limoges: Lambert-Lucas.

767   Rollet, N. (2013). "D'accord". Approche conversationnelle et multimodale d'une forme

768          située dans les appels au Samu-Centre 15. *L'Information grammaticale, 139*.

769   Rollet, N., Jain, V., Licoppe, C., & Devillers, L. (2017). Towards Interactional

770          Symbiosis: Epistemic Balance and Co-presence in a Quantified Self Experiment.

771          In L. Gamberini, A. Spagnolli, G. Jacucci, B. Blankertz, & J. Freeman (Eds.),

772          *Symbiotic Interaction: 5th International Workshop, Symbiotic 2016, Padua,*

773          *Italy, September 29–30, 2016, Revised Selected Papers* (pp. 143-154). Cham:

774          Springer International Publishing.

775   Sacks, H. (1984). On doing 'being ordinary'. In J. M. Atkinson & J. Heritage (Eds.),

776          *Structures of Social Action* (pp. 413-429). Cambridge: CUP.

777   Sacks, H. (1992). *Lectures on conversation* (Jefferson, G. ed.). Oxford: Blackwell.

778   Sacks, H., & Schegloff, E. (1973). Opening up closings. *Semiotica, 8*, 289-326.

779   Sacks, H., Schegloff, E. A., & Jefferson, G. (1974). A symplest systematics for the

780          organization of turn-taking for conversation. *Language, 50*, 696-731.

781   Sadazuka, K., Kuno, Y., Kawashima, M., & Yamazaki, K. (2007). *Museum Guide*

782          *Robot with Effective Head Gestures*. Paper presented at the International

783          Conference on Control, Automation and Systems, Seoul, Korea.

784   Schegloff, E. (1987). Between micro and macro: contexts and other connections. In B.

785         Giesen, J.C. Alexander, R. Münch, & N.J. Smelser (Eds.), *The Micro-macro link*

786         (pp. 207-234). Los Angeles: University of California Press.

787   Schegloff, E. (2002). Accounts of conduct in interaction. Interruption, overlap, and turn-

788         taking. In T. J.H. (Ed.), *Handbook of sociological theory* (pp. 287-321). New

789         York: Kluwer Academic/Plenum publishers.

790   Schegloff, E. (2007). *Sequence organization in interaction. A Primer in Conversation*

791         *Analysis.* (Vol. 1): Cambridge University Press.

792   Sidner, C. L., & Dzikovska, M. (2002, 16-16 Oct. 2002). *Human-robot interaction:*

793         *engagement between humans and robots for hosting activities.* Paper presented

794         at the 4th IEEE International Conference on Multimodal Interfaces,

795         Pittsburgh, USA.

796   Sidner, C. L., Lee, C., Kidd, C. D., & Rich, C. (2005). Explorations in engagement for

797         humans and robots. *Artificial Intelligence, 166*, 140-164.

798   Stivers, T. (2008). Stance, Alignment, and Affiliation During Storytelling: When

799         Nodding Is a Token of Affiliation. *Research on Language & Social Interaction,*

800         *44*(1), 31-57.

801   Stivers, T. (2015). Coding Social Interaction: A Heretical Approach in Conversation

802         Analysis? *Research on Language and Social Interaction, 48*(1), 1-19.

803   Šabanović, S., & Chang, W.-L. (2016). Socializing robots: constructing robotic sociality

804         in the design and use of the assistive robot PARO. *AI and Society, 31*(4), 537-

805         551.

806     Straub, I. (2016). "It looks like a human!" The interrelation of social presence,

807             interaction and agency ascription: a case study about the effects of an android

808             robot on social agency ascription. *AI and Society, 31*(4), 553-571.

809     Suchman, L. (2007). *Human-Machine Reconfigurations. Plans and situated actions,*

810             *2nd edition*. New York: Cambridge University Press.

811     Yu, Z., Scherer, S., Devault, D., Gratch, J., Stratou, G., Morency, L., & Cassell, J.

812             (2013). *Multimodal Prediction of Psychological Disorder: Learning Verbal and*

813             *Nonverbal Commonality in Adjacency Pairs*. Paper presented at the 17th

814             Workshop Series on the Semantics and Pragmatics of Dialogue, Amsterdam,

815             Netherland.

816     Zimmerman, D. (2006). *How closing matters in emergency telephone calls*. Paper

817             presented at the Annual meeting of the American Sociology Association,

818             Montréal, Canada.