

A corpus for studying addressing behaviour in multi-party dialogues

Natasa Jovanovic · Rieks op den Akker · Anton Nijholt

Received: ■ / Accepted: ■

© Springer Science+Business Media B.V. 2006

Abstract This paper describes a multi-modal corpus of hand-annotated meeting dialogues that was designed for studying addressing behaviour in face-to-face conversations. The corpus contains annotated dialogue acts, addressees, adjacency pairs and gaze direction. First, we describe the corpus design where we present the meetings collection, annotation scheme and annotation tools. Then, we present the analysis of the reproducibility and stability of the annotation scheme.

Keywords Addressing · Multi-party dialogues · Multimodal corpora · Annotation schemas · Reliability analysis

1 Introduction

Current tendencies in modelling human-computer as well as human-human interactions are moving from a two-party model to a multi-party model. One of the issues that becomes salient in interactions involving more than two parties is addressing. Addressing as an aspect of every form of communication has been extensively studied by conversational analysts and social psychologists (Clark & Carlson, 1992; Goffman, 1981; Goodwin, 1981). Recently, addressing has received considerable attention in interaction modelling in the context of mixed human-human and human-computer interaction (Bakx, van Turnhout, & Terken, 2003; van Turnhout, Terken, Bakx, & Eggen, 2005), human-human-robot interaction (Katzenmaier, Stiefelhagen, & Schultz, 2004), mixed human-agents and multi-agents interaction (Traum, 2004) and multi-party human-human interaction (Jovanovic & op den Akker, 2004; Otsuka, Takemae, Yamato, & Murase, 2005).

Addressing is carried out through various communication channels, such as speech, gestures or gaze. To explore interaction patterns in addressing behaviour

N. Jovanovic · R. op den Akker (✉) · A. Nijholt
Human Media Interaction Group, University of Twente, PO Box 217, Enschede 7500 AE,
The Netherlands
e-mail: infrieks@ewi.utwente.nl

25 and to develop models for automatic addressee prediction, we need a collection of
26 audio and video interaction recordings that contains a set of annotations relevant to
27 addressing. Meetings as complex interplays of interacting participants represent a
28 relevant domain for the research on different aspects of interactions involving more
29 than two participants who employ a variety of channels to communicate with each
30 other.

31 In the context of the meeting research, several corpora have already been
32 developed. Some of the existing meeting corpora, such as the ICSI (Janin et al.,
33 2004) and ISL (Burger & Sloane, 2004) corpora—currently widely used to study
34 linguistic phenomena in natural meetings—are limited to audio data only. The NIST
35 audio–visual meeting corpus (Garofolo, Laprun, Michel, Stanford, & Tabassi, 2004)
36 is designed to support the development of audio and video recognition technologies
37 in the context of meetings. Currently, it provides transcriptions of the meetings to
38 enable the research on automatic speech recognition in meetings. To support the
39 research on higher-level meetings understanding, the VACE (Chen et al., 2006) and
40 AMI (Carletta et al., 2006) multi-modal meeting corpora are currently being pro-
41 duced. The VACE corpus is being developed to support research on multimodal
42 cues, such as speech, gaze, gestures and postures, for understanding meetings. The
43 AMI data collection is being developed to enhance research in various areas related
44 to the development of meeting browsing technologies and remote meeting assistants,
45 including speech recognition, computer vision, discourse and dialogue modelling,
46 content abstraction, human–human and human–computer interaction modelling. It
47 contains a range of annotations including, among others, speech transcription, dia-
48 logue acts, topic segmentation, focus of attention, head and hand communicative
49 gestures, and summaries.

50 In this paper, we describe a multi-modal corpus of hand-annotated meeting dia-
51 logues, designed for studying addressing behaviour in face-to-face conversations.
52 The meetings were recorded in the IDIAP meeting room in the research program of
53 the European M4¹ and AMI² projects. The recordings are available through the
54 MultiModal Media File Server.³ Currently, the corpus contains hand-annotated
55 dialogue acts, adjacency pairs, addressees and gaze directions of meeting partici-
56 pants. A set of the corpus' annotations of the M4 meetings is available as a part of
57 the M4 meeting collection.⁴

58 Apart from the corpus description which includes the description of the meeting
59 data, annotation scheme, annotation tools and the corpus format, this paper reports
60 the reliability of the overall annotation scheme as well as a detailed analysis of
61 detected sources of unreliability.

62 2 Meeting data

63 The corpus consists of 12 meetings recorded at the IDIAP smart meeting room
64 (Moore, 2002). The room is equipped with fully synchronized multi-channel audio

¹ The M4 (MultiModal Meeting Manager) project: <http://www.m4project.org>

² The AMI (Augmented Multi-party Interaction) project: <http://www.amiproject.org>

³ MMM File Server <http://www.mmm.idiap.ch>

⁴ <http://www.mmm.idiap.ch/M4-Corpus/annotations/NXTbasedAnnotation/>

65 and video recording devices (see Fig. 1). Of the 12 meetings, 10 were recorded
 66 within the scope of the M4 project. These meetings are scripted in terms of type
 67 and schedule of group actions that participants perform in meetings such as pre-
 68 sentation, discussion or note taking, but content is natural, spontaneous and
 69 unconstrained. Spontaneous behaviour of participants in these meetings allows us
 70 to examine observable patterns of addressing behaviour in small group discussions.
 71 More natural, scenario-based, meetings have been recorded in the scope of the
 72 AMI project. One of the AMI pilot meetings recorded at the IDIAP meeting room
 73 is included in our corpus. The meeting involves a group focused on the design of a
 74 TV remote control. The last meeting in our corpus is one of a series of meetings
 75 recorded at IDIAP for the exploration of argumentative structures in meeting
 76 dialogues.

77 Research on small group discussions presented in (Carletta, Anderson, &
 78 Garrod, 2002) has shown that there is a noticeable difference in the interaction
 79 patterns between large and small groups. A small group discussion involving up to
 80 seven participants resembles two-way conversations that occur between all pairs of
 81 participants and every participant can initiate conversation. A large group dis-
 82 cussion is more like a series of conversations between a group leader and various
 83 individuals with the rest of participants present but silent. In the M4 and AMI data
 84 collection each meeting consists of 4 participants. Hence, the meetings in our
 85 corpus satisfy the interaction patterns of small group discussions. There are 23
 86 participants in the corpus. The total amount of recorded data is approximately
 87 75 min.

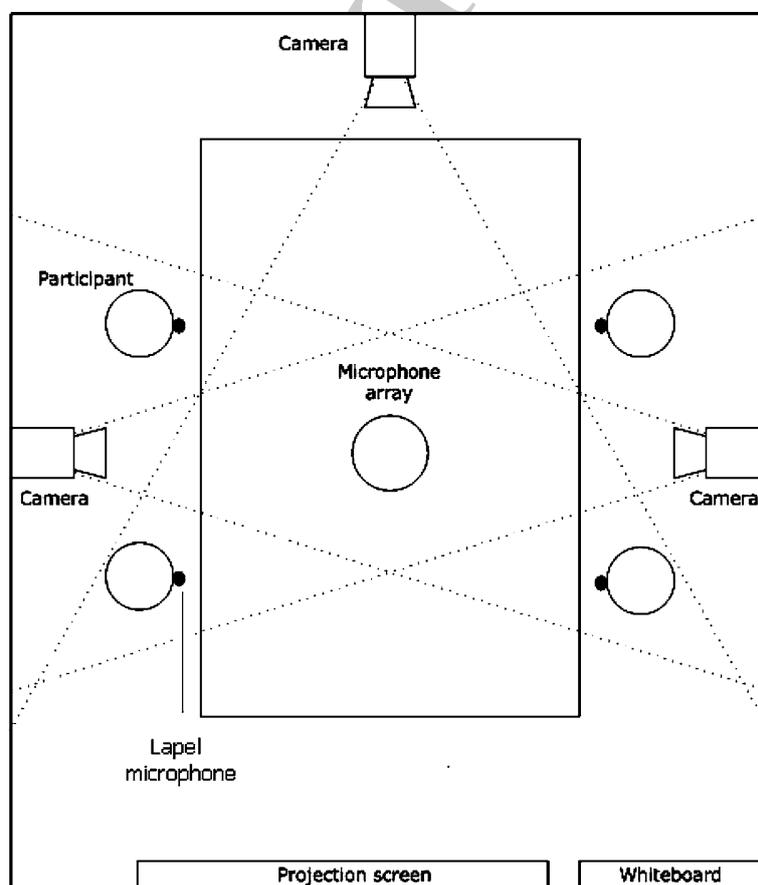


Fig. 1 The configuration of the IDIAP meeting room (M4 data collection)

88 **3 Annotation scheme**

89 In two-person dialogues, it is usually obvious to the non-speaking participant who is
 90 the one being addressed by the current speaker. In a multi-party case, the speaker
 91 has not only the responsibility to make his speech understandable for the listeners,
 92 but also to make clear to whom he is addressing his speech.

93 Analysis of the mechanisms that people use in identifying their addressees leads
 94 to a model of a conversation that describes the features that play a role in these
 95 mechanisms. Our annotation scheme is based on the model presented in (Jovanovic
 96 & op den Akker, 2004). The features described in the model are of three types:
 97 verbal, nonverbal and contextual. For example, utterances that contain the proper
 98 name of a conversational participant may be addressed to that participant. Also
 99 speaker gaze behaviour may be a feature that gives a hint to the intended
 100 addressee. The history of the conversation is important as well, since most of the
 101 utterances that are related to the previous discourse are addressed to one of the
 102 recent speakers.

103 Although the model contains a rich set of features that are relevant for observers
 104 to identify the participants the speaker is talking to, due to time constraints, the
 105 meetings were annotated with a subset of the selected properties. In addition to
 106 addressee annotation, the corpus currently contains annotations of dialogue acts,
 107 adjacency pairs and gaze direction. We also considered coding of deictic hand ges-
 108 tures as they can be used as a means of addressing. However, it was found that
 109 deictic hand gestures occur very rarely in the data.

110 **3.1 Dialogue acts**

111 Annotation of dialogue acts involves two types of activities: marking of dialogue acts
 112 segment boundaries and marking of dialogue acts themselves.

113 Utterances within speech transcripts, also known as prosodic utterances, were
 114 segmented in advance using prosody, pause and syntactical information. In our
 115 scheme, a dialogue act segment may contain a part of a prosodic utterance, a whole
 116 prosodic utterance, or several contiguous prosodic utterances of the same speaker.

117 Our dialogue act tag set is based on the MRDA (Meeting Recorder Dialogue
 118 Act) set (Dhillon, Bhagat, Carvey, & Shriberg, 2004). The MRDA tag set represents
 119 a modification of the SWDB-DAMSL tag set (Jurafsky, Shriberg, & Biasca, 1997)
 120 for an application to multi-party meeting dialogues. Each functional utterance in
 121 MRDA is marked with a label, made up of one or more tags from the set. The
 122 analysis of the MRDA tag set presented in (Clark & Popescu-Belis, 2004) shows that
 123 the number of possible labels reaches several millions. For that reason, the usage of
 124 the complete set may lead to a low quality of manual annotations.

125 Unlike MRDA, each utterance in our dialogue act annotation scheme is marked
 126 as *Unlabelled* or with exactly one tag from the tag set that represents the most
 127 specific utterance function. For addressee identification, it is less important whether
 128 an utterance is a suggestion in the form of a question or in the form of a statement.
 129 More important is that the speaker suggests to the addressee to perform an action,
 130 informing all other participants about that suggestion.

131 Our dialogue act tag set is created by grouping some of the MRDA tags into 17
 132 categories that are divided into seven groups, as follows:

- 133 • Statements
- 134 • *Statement* [MRDA: Statement]. The Statement tag marks utterances which are
135 objective and factual statements as well as utterances which are opinions and
136 other subjective statements.
- 137 • Acknowledgements and Backchannels
- 138 • *Acknowledgement* [MRDA: Acknowledgement, Backchannel]. The Acknowl-
139 edgement tag is a common tag used for acknowledgements and backchannels.
140 Acknowledgements are utterances in which a speaker acknowledges a previous
141 speaker's utterances or a significant portion of a previous speaker's utterance.
142 They are neither positive nor negative. Backchannels have a function to show
143 that a listener is paying attention. They are made in the background by a speaker
144 who does not have the floor.
- 145 • *Assessment/Appreciation* [MRDA:Assessment/Appreciation]. The Assessment/
146 Appreciation tag marks utterances that are acknowledgements directed to an-
147 other speaker's previous utterance with slightly more emotional involvement.
148 They can be positive, such as "that's great", "wow!", or negative, such as "not
149 good enough", "that's impossible".
- 150 • Questions
- 151 • *Information Request* [MRDA: Wh-Question, Y/N Question, OR-Question, Or
152 Clause After Y/N Question]. The Information Request tag marks questions that
153 require specific answers. Examples include "what kind of preprocessing are you
154 using?" or "but do you often cook at night?"
- 155 • *Open-ended Question* [MRDA: Open-ended Question]. The Open-ended
156 Question tag marks questions that do not require a specific answer; they are
157 rather asked in a broad sense (e.g., "What about you?" or "anything else?").
- 158 • *Rhetorical Question* [MRDA: Rhetorical Question]. The Rhetorical Question tag
159 marks questions that are used for rhetorical effects. No answer is expected to
160 those questions. Examples include "who knows?" or "who would have thought
161 that it was possible?"
- 162 • Responses
- 163 • *Positive Response* [MRDA: (Partial) Accept, Affirmative Answer]. The Positive
164 Response tag marks utterances that exhibit a (partial) agreement to or a (partial)
165 acceptance of or an affirmative answer to a previous speaker's proposal, state-
166 ment or question.
- 167 • *Negative Response* [MRDA: (Partial) Reject, Dispreferred Answer, Negative
168 Answer]. The Negative Response tag marks utterances which exhibit a (partial)
169 disagreement to or a (partial) reject of or an explicit or implicit negative answer
170 to a previous speaker's proposal, statement or question.
- 171 • *Uncertain Response* [MRDA: Maybe, No Knowledge]. The Uncertain Response
172 tag marks utterances which express a lack of a speaker's knowledge regarding
173 some subject or that a speaker's utterance is probable, yet not definite (e.g.,
174 "maybe", "I am not sure").
- 175 • Action Motivators

- 176 • *Influencing-listeners-action* [MRDA: Command, Suggestion]. The Influencing-
 177 listeners-action tag marks utterances which influence the listeners' communicative
 178 or non-communicative future actions such as commands, suggestions, proposals,
 179 advices.
- 180 • *Committing-speaker-action* [MRDA: Command, Suggestion]. The Committing-
 181 speaker-action tag marks utterances which indicate that a speaker has committed
 182 himself, in varying degrees of strength, to some future course of action. The
 183 speaker can explicitly commit himself that he will execute an action at some point
 184 in the future, such as "I will prepare a presentation for the next meeting", or he
 185 can suggest that he will do so if listeners accept it, such as "I can say something
 186 about that".
- 187 • Checks
- 188 • *Follow Me* [MRDA: Follow Me]. The Follow Me tag marks utterances by which
 189 a speaker wants to ensure that what he is saying has been understood by lis-
 190 tener(s) (e.g., "do you understand?", "okay?", "this is clear?").
- 191 • *Repetition Request* [MRDA: Repetition Request]. The Repetition Request tag
 192 marks utterances in which a speaker wants another speaker to repeat all or a part
 193 of a previous utterance. This is mostly the case when a speaker could not hear or
 194 could not interpret what another speaker has said and wants to hear it again.
- 195 • *Understanding Check* [MRDA: Understanding Check]. The Understanding
 196 Check tag marks utterances in which a speaker wants to make sure whether he
 197 understands what a previous speaker said or whether he understands some sort of
 198 information. Examples include "you said that machine learning techniques are
 199 applicable?", "so this part is new, right?".
- 200 • Politeness Mechanisms
- 201 • *Thanks* [MRDA: Thanks]. The Thanks tag marks utterances in which a speaker
 202 thanks another speaker(s).
- 203 • *Apology* [MRDA: Apology]. The Apology tag marks utterances in which a
 204 speaker apologizes for something he did (e.g., coughing, interrupting another
 205 speaker) or he plans to do (e.g., to leave meeting earlier, to make a phone call
 206 during the meeting).
- 207 • *Other polite* [MRDA: Welcome, Downplayer, Sympathy]. The Other polite tag
 208 marks all other acts of politeness that do not contribute to the overall discussion
 209 but rather have a social impact such as, "you're welcome", "I'm kidding", "good
 210 luck", "success", "you are so nice".

211 The MRDA scheme also allows the annotation of turn-taking (e.g., floor grabber)
 212 and turn-maintaining (e.g., floor holder) mechanisms. The turn managing dimension
 213 of utterances' functions is excluded from our scheme. Utterances that function only
 214 as turn taking, turn giving or turn holding signals are marked as *Unlabelled*. Turn-
 215 taking and addressing as two aspects of conversational interactions are related, but
 216 we were specifically interested in studying how addressing actually works, i.e., how
 217 people address each other, in order to build predictive models for addressee iden-
 218 tification. The scheme also excludes (1) a set of MRDA tags that are related to
 219 restating information such as repetitions and corrections, (2) a set of MRDA tags
 220 that are related to rhetorical roles such as explanations or elaborations and (3) a set

221 of MRDA tags that provide further descriptions of utterance functions such as self
222 talk, third party talk, jokes, meeting agendas or topic change.

223 3.2 Adjacency pairs

224 Adjacency pairs (APs) are minimal dialogic units that consist of paired utterances
225 such as question–answer or statement–agreement. The paired utterances are pro-
226 duced by different speakers. Utterances in an adjacency pair are ordered with the
227 first part (A-part, the initiative) and the second part (B-part, the response). In multi-
228 party conversations, adjacency pairs do not impose a strict adjacency requirement,
229 since a speaker has more opportunities to insert utterances between two elements of
230 an adjacency pair. For example, a suggestion can be followed by agreements or
231 disagreements from multiple speakers.

232 In our scheme, adjacency pairs are labelled at a separate level from dialogue acts.
233 Labelling of adjacency pairs consists of marking dialogue acts that occur as their A-
234 part and B-part. If a dialogue act is an A-part with several B-parts, for each of these
235 B-parts, a new adjacency pair is created. Furthermore, each dialogue act is marked
236 as a B-part of at most one and as an A-part of an arbitrary number of adjacency
237 pairs. Although it is theoretically possible that a B-part is related to several A-parts,
238 for example, an utterance may answer two questions, the analysis of the data showed
239 that these cases hardly occur.

240 3.3 Addressees

241 In a group discussion, many of the speaker's utterances are addressed to the group as
242 a whole. However, the speaker may show by verbal or non-verbal behaviour that he
243 intends to affect one selected participant or a subgroup of participants in particular,
244 that he expects that participant or that subgroup to react on what he says. In this
245 case, the selected participant or the subgroup is the addressee of the dialogue act
246 performed by the speaker.

247 Given that each meeting in the corpus consists of four participants, the addressee
248 tag set contains the following values:

- 249 • a single participant: P_x
- 250 • a subgroup of participants: P_x, P_y
- 251 • the whole audience: P_x, P_y, P_z
- 252 • *Unknown* $x, y, z \in \{0, 1, 2, 3\}$; P_x denotes speaker at the channel x . The *Unknown* tag is
253 used when the annotator cannot determine to whom the dialogue act is addressed.
254 It is also used for utterances marked as *Unlabelled* and for backchannels.

255 3.4 Gaze direction

256 Annotation of gaze direction involves two types of activities: labelling the changes in
257 the gazed targets and labelling the gazed targets themselves.

258 For addressee identification, the only targets of interest are meeting participants.
259 Therefore, the tag set contains tags that are linked to each participant (P_x) where
260 $x \in \{0, 1, 2, 3\}$ and the *NoTarget* tag that is used when the speaker does not look at any
261 of the participants. The set can be further refined by adding some objects of interests
262 in the meeting room such as whiteboard, projector screen or notebook.

263 Since the meeting room was not equipped with close-up cameras the gaze
 264 information was induced based on the side and central cameras (see Fig. 1). This was
 265 the main reasons for not imposing the requirement for a high precision in labelling
 266 changes in the gazed targets.

267 4 Annotation tools

268 The corpus was created using two annotation tools developed at the University of
 269 Twente: the DACoder (Dialogue Act Coder) and the CSL (Continuous Signal
 270 Labelling) tools (Reidsma, Hofs, & Jovanovic, 2005). The DACoder supports
 271 annotation of dialogue acts, addressees and any kind of relations between dialogue
 272 acts such as adjacency pairs or rhetorical relations. The CSL tool supports labelling
 273 of time-aligned annotation layers directly related to the signal files. Any annotation
 274 layer that consists of simple labelling of non-overlapping segments of the time line
 275 can be coded using this tool (e.g., gaze directions or postures).

276 The tools were built using NXT (NITE XML Toolkit) (Carletta et al., 2003).
 277 NXT uses a stand-off XML data storage format which consists of several inter-
 278 related xml-files. The structure and location of the files are represented in a
 279 “metadata” file. The NXT stand-off XML format enables the capture and efficient
 280 manipulation of complex hierarchical structures across different modalities. Fur-
 281 thermore, it supports an easy extension of the corpus with new annotation layers
 282 without influencing exiting annotations. For exploitation of annotated data, NXT
 283 provides the NXT Search tool for the execution of the queries expressed in the NXT
 284 Query Language (NQL).⁵

285 5 Distributional statistics

286 In this section, we provide distributional statistics for the dialogue act tags and
 287 addressee tags in the corpus. The corpus contains 1457 dialogue act segments out of
 288 which 131 segments (8.99%) are labelled as *Unlabelled*. Table 1 shows the distri-
 289 bution of DA tags after discarding those segments that are marked as *Unlabelled*.

290 The distribution of the addressee tags over those segments that are marked with a
 291 DA label is presented in Table 2. All subgroup addressee tags (P_x, P_y) are grouped
 292 into the SUB category and all tags that denote the whole audience (P_x, P_y, P_z) are
 293 grouped into the ALLP category.

294 6 Reliability

295 In order to obtain valid research results, data on which they are based must be
 296 reliable. We have performed two reliability tests proposed by Krippendorff (1980):
 297 stability (intra-annotator reliability) and reproducibility (inter-annotator reliability).
 298 Stability is the degree to which an annotator’s judgments remain unchanged over
 299 time. It is measured by giving the same annotator a set of data to annotate twice, at
 300 different times. Reproducibility is the degree to which different annotators can

⁵ NXT Query Language: <http://www.ims.uni-stuttgart.de/projekte/nite/>

Table 1 Distribution of DA tags

Statement	44.19%	Influencing-listeners-action	3.32%
Information request	9.58%	Committing-speaker-action	2.64%
Open-ended question	2.11%	Follow me	0.15%
Rhetorical question	0.60%	Repetition request	0.53%
Acknowledgement	15.61%	Understanding check	1.58%
Assessment/Appreciation	2.19%	Thanks	0.38%
Positive response	11.99%	Apology	0.08%
Negative response	3.24%	Other polite	0.30%
Uncertain response	1.51%		

Table 2 Distribution of addressee tags

Addressee	
P0	12.97%
P1	14.63%
P2	17.50%
P3	16.59%
ALLP	34.16%
SUB	1.73%
Unknown	2.41%

301 produce the same annotation. It is measured by giving several annotators the same
302 data to annotate independently, following the same coding instructions.

303 Reliability is a function of agreement achieved among annotators. In the dialogue
304 and discourse processing community, the Kappa agreement coefficient (κ) has been
305 adopted as a standard (Carletta, 1996; Cohen, 1960). In recent years, there have been
306 some discussions about the usage of Kappa as an appropriate reliability metric.
307 Krippendorff's Alpha (α) has been proposed as a more adequate metric for assessing
308 reliability of subjective codings (Krippendorff, 1980, 2004).

309 To estimate reliability of dialogue act, addressee and gaze annotation, we applied
310 both agreement coefficients. The obtained Kappa and Alpha values were identical.
311 Therefore, in the following sections we report only Kappa values. In contrast to
312 dialogue act and addressee annotation, adjacency pairs annotation cannot be con-
313 sidered as a simple labelling of annotation units with categories. Therefore, we
314 developed our own approach that represents annotated APs in a form of categorical
315 labelling and measures agreement on APs annotation using Alpha.

316 For the evaluation of Alpha and Kappa values, we used Krippendorff's scale that
317 has been adopted as standard in the discourse and dialogue processing community
318 (Krippendorff, 1980). According to that scale, any variable with an agreement
319 coefficient below .67 is disregarded as unreliable, between .67 and .8 allows drawing
320 tentative conclusions and above .80 allows drawing definite conclusions.

321 6.1 Detecting sources of unreliability

322 Detecting causes of disagreement may be of great use to obtain reliable data or to
323 improve data reliability. A source of unreliability can be a coding unit, a category, a
324 subset of categories or an annotator (Krippendorff, 1980). Even if a category is well
325 defined, annotators may still have different interpretations of the category. Fur-
326 thermore, annotators may show a correlated disagreement. For example, annotator
327 A_1 uses category C_1 whenever annotator A_2 uses category C_2 .

328 To identify which categories are sources of unreliability we measured single-category
 329 reliability (Krippendorff, 1980). Single-category reliability assesses the extent to which
 330 one category is confused with all other categories in a set. It is estimated by grouping
 331 the remaining categories into one category and measuring the agreement among
 332 annotators regarding the assignment of units to these two categories. A low agreement
 333 can be the result of an ambiguous definition of the category or of the coders' inability to
 334 interpret the meaning of the category.

335 7 Inter-annotator reliability

336 In this section we present the analysis of inter-annotator reliability of the annotation
 337 scheme applied on the M4 meeting data.

338 Six trained annotators were involved in the corpus creation. They were divided
 339 into two groups: the DA (Dialogue Act) group and the VL (Video Labelling) group.
 340 The DA group, involving 4 annotators, annotated dialogue acts, addressees and
 341 adjacency pairs. The VL group, involving 2 annotators, annotated gaze direction.
 342 The corpus was divided into two sets of meetings. The DA group was divided into 2
 343 subgroups of 2 annotators: the B&E group and the M&R group. Each of these
 344 subgroups annotated exactly one set of meeting data. Each annotator in the VL
 345 group annotated one set of meeting data. Additionally, two meetings were annotated
 346 by both annotators in the VL group in order to test reliability of gaze annotation. In
 347 summary, each meeting in the corpus was annotated with dialogue acts, addressees
 348 and adjacency pairs by exactly two annotators, and with participants' gaze directions
 349 by at most two annotators.

350 7.1 Reliability of dialogue acts annotation

351 We first measured agreements among annotators on how they segmented dialogues
 352 into dialogue act segments. Then, we tested reliability of dialogue act classification
 353 on those segments for which annotators agreed on their boundaries.

354 7.1.1 Segmentation reliability

355 In the discourse and dialogue community, several approaches have been proposed
 356 for assessing segmentation reliability using various metrics: percent agreement
 357 (Carletta et al., 1997; Shriberg, Dhillon, Bhagat, Ang, & Carvey, 2004), precision
 358 and recall (Passonneau & Litman, 1997), and κ (Carletta et al., 1997; Hirschberg &
 359 Nakatani, 1996).

360 Since there is no standardized technique to estimate segmentation agreement, we
 361 developed our own approach based on percent agreement. We defined four types of
 362 segmentation agreement:

- 363 • *Perfect agreement (PA)*—Annotators completely agree on the segment
 364 boundaries.
- 365 • *Contiguous segments of the same type (ST)*—A segment of one annotator is
 366 divided into several segments of the same type by the other annotator. Segments
 367 are of the same type if they are marked with the same dialogue act tag and the

368 same addressee tag. An additional constraint is that segments are not labelled as
369 parts of APs.

- 370 • *Unlabelled-DA (UDA)*—A segment of one annotator is divided into two seg-
371 ments by the other annotator where one of those segments is marked as *Unla-*
372 *belled* and the other one with a dialogue act tag.
- 373 • *Conjunction-Floor(CF)*—Two adjacent segments differ only in a conjunction or a
374 floor mechanism at the end of the first segment. The following example shows the
375 segmentation agreement of this type:

- 376 (1) I can do that—but I need your help
- 377 (2) I can do that but—I need your help

378

379 The approach takes one annotator's segmentation as a reference (*R*) and com-
380 pares it with the other annotator's segmentation (*C*) segment by segment. As a
381 result, it gives a new segmentation (*C'*) that represents the modification of (*C*) to
382 match the reference segmentation (*R*) according to identified types of agreement. In
383 addition to measuring segmentation agreement, the modified segmentation (*C'*) is
384 used for assessing reliability of dialogue act classification, addressee classification
385 and adjacency pairs annotation. Table 3 shows overall segmentation results for each
386 annotation group.

387 Most of the segmentation disagreements are of the following three types. First,
388 while one annotator labelled a segment with the *Acknowledgement* tag, the other
389 one included the segment in the dialogue act that follows. Second, while one
390 annotator marked a segment with one of the response tags, the other annotator split
391 the segment into a response and a statement that has a supportive function such as
392 explanation, elaboration or clarification. Third, while one annotator split a segment
393 into two or more segments labelled with the same dialogue act tag but different
394 addressee tags, the other annotator marked it as one segment.

395 7.1.2 Reliability of dialogue act classification

396 Reliability of dialogue act classification is measured over those dialogue act seg-
397 ments for which both annotators agreed on their boundaries. Since the number of
398 agreed segments for each R–C pair is different, we calculated reliability of dialogue
399 act classification for each pair. The results are shown in Table 4. According to
400 Krippendorff's scale annotators in each DA group reached an acceptable level of
401 agreement that allows drawing tentative conclusions from the data.

402 We applied a single-category reliability test for each dialogue act tag to assess the
403 extent to which one dialogue tag was confused with the other tags in the set. Table 5

Table 3 Segmentation agreement (R–C pair: reference annotator (R)–comparison annotator (C))

R–C	Agreement types				Agree	Total	Agree %
	PA	ST	UDA	CFM			
B–E	326	22	16	2	366	406	90.15
E–B	326	32	17	2	377	411	91.73
M–R	317	29	10	2	358	419	85.44
R–M	317	33	15	2	367	426	86.14



Table 4 Inter-annotator agreement on DA classification

Group	R-C pair	N	κ
M&R	M-R	358	0.70
	R-M	367	0.70
B&E	B-E	366	0.75
	E-B	377	0.77

Table 5 Single-category reliability for DA tags (Kappa values)

Category	B-E	M-R
Statement	0.82	0.72
Acknowledgement	0.87	0.75
Assessment/Appreciation	0.32	0.39
Information request	0.70	0.84
Open-ended question	0.74	0.84
Repetition request	1.00	1.00
Rhetorical question	0.00	0.66
Influencing-listeners-action	0.58	0.70
Committing-speaker-action	0.86	0.74
Positive response	0.70	0.52
Uncertain response	0.80	0.50
Negative response	0.67	0.61
Understanding check	0.32	-0.01
Other polite	0.00	-
Thanks	1.00	1.00
Follow me	-	- 0.003

404 shows the results of performing the Kappa tests for only one R-C pair in each DA
405 group.

406 Annotators in the B&E group used different ranges of categories: the *Other polite*
407 and *Rhetorical Question* categories, which occur rarely in the data, were employed
408 only by annotator B. For that reason, Kappa values for these categories are zero.
409 Negative Kappa values for *Understanding Check* and *Follow me* categories indicate
410 that annotator agreement is below chance: in all cases where one annotator identifies
411 one of these two categories, the other annotator does not. The results show an
412 unacceptably low agreement on *Assessment/Appreciation* and *Understanding Check*
413 categories in both groups. The *Assessment/Appreciation* category was mainly con-
414 fused with *Positive Response* and *Statement* categories. The *Understanding Check*
415 category was mostly confused with *Information Request* and *Statement* categories.
416 Annotators in the M&R group reached a lower agreement on the responses tags
417 than annotators in the B&E group. The responses tags were mostly confused with
418 the *Statement* tag. Additionally, annotators in the M&R group had a little more
419 difficulty distinguishing *Positive Response* from *Assessment/Appreciation* and
420 *Acknowledgement*. The low Kappa value for the *Influencing-listener-actions* category
421 in the B&R group is a result of the confusion with the *Statement* category.

422 7.2 Reliability of addressee annotation

423 As for dialogue act classification, reliability of addressee annotation is measured
424 over those dialogue act segments for which both annotators agreed on their
425 boundaries.

Table 6 Inter-annotator agreement on addressee annotation

Group	R-C pair	N	κ
M&R	M-R	358	0.68
	R-M	367	0.70
B&E	B-E	366	0.79
	E-B	377	0.81

426 The Kappa values for addressee annotation are shown in Table 6.

427 The results show that annotators in the B&E group reached good agreement on
428 addressee annotation, whereas annotators in the M&R group reached an acceptable
429 level of agreement.

430 Annotators mainly disagreed on whether an individual or a group had been ad-
431 dressed. When annotators agreed that an individual had been addressed, they agreed
432 in almost all cases which individual it had been. There were only a few instances in
433 the data labelled with categories that represent subgroup addressing. In both DA
434 groups, annotators failed to agree on those categories. Annotators had problems
435 distinguishing subgroup addressing from addressing the group as a whole.

436 We measured single-category reliability for those addressee tags that represent
437 individual and group addressing. Single-category reliability is measured using the
438 Kappa test for one R-C pair in each group. Addressee values that consist of three
439 participants such as p_0, p_1, p_3 or p_1, p_2, p_3 were grouped into one category that repre-
440 sents the whole audience (ALLP). Annotators in the B&E group reached a good
441 agreement ($\kappa \geq 0.80$; $N = 369$) on all categories representing a single participant.

442 Agreement on ALLP was $\kappa = 0.77$. Annotators in the M&R group reached a
443 lower agreement on each category than annotators in the B&E group. They had a
444 little more difficulty distinguishing ALLP ($\kappa = 0.63$; $N = 366$) as well as p_3 ($\kappa = 0.59$;
445 $N = 366$) from a remaining set of categories. For all other categories representing a
446 single participant Kappa was $0.71 \leq \kappa < 0.80$.

447 7.3 Reliability of adjacency pairs annotation

448 According to our scheme for annotation of adjacency pairs, each dialogue act can be
449 marked as a B-part of at most one and as an A-part of an arbitrary number of
450 adjacency pairs. The sets of adjacency pairs produced by two annotators may differ
451 in several ways. First, the annotators may disagree on dialogue acts that are marked
452 as A-parts of adjacency pairs. Second, they may assign a different number of B-parts
453 as well as different B-parts themselves to the same A-part.

454 Since there seems to be no standard associated metric for agreement on APs
455 annotation in the literature, we developed a new approach that resembles a method
456 for measuring reliability of co-reference annotation proposed in (Passonneau, 2004).
457 The key of the approach is to represent annotated data as a form of categorical
458 labelling in order to apply standard reliability metrics.

459 Adjacency pairs annotation can be seen as assigning to each dialogue act a *context*
460 that represents the relations that the dialogue act has with surrounding dialogue acts.
461 To encode the contexts of dialogue acts, we define a set of classes that contain
462 related dialogue acts. For each A-part, all its B-parts are collected in one class.
463 Therefore, a class is characterized with its A-part and a set of B-parts (b-set):
464 $\langle a, bset(a) \rangle$ where $bset(a) = \{b | (a, b) \in AP\}$. A dialogue act can belong to at most



465 two classes: a class containing the dialogue act as an A-part (A-class) and a class
 466 containing the dialogue act as a B-part (B-class). Thus, the complete context of a
 467 dialogue act is encoded with an AP label (L) that is compounded of its A-class and
 468 B-class ($L=A\text{-class|B-class}$).

469 Given a list of dialogue acts $DA = [da_1, \dots, da_n]$, a class can be represented in two
 470 different ways: with fixed or relative position of the dialogue acts. The former en-
 471 codes each dialog act in the class with the index of the dialog acts in the list. The
 472 latter encodes the dialogue acts in the class with relative positions to the dialogue act
 473 representing the A-part of the class. In this paper, we use the approach with relative
 474 positions because it significantly decreases the number of possible classes. In our
 475 encoding, each class of the labelled dialogue act da_i (A-class and B-class) has the
 476 form $\langle -n, O \rangle$, where n is an offset of the labelled DA da_i from the A-part of the
 477 class and O is a set of offsets of the dialogue acts in the b-set from the A-part of the
 478 class. Note that for the A-class, n is always 0 since the labelled dialogue act is the A-
 479 part of the class. For the B-class, n is always positive because the labelled dialogue
 480 act is in the b-set and the A-part always precedes dialogue acts in the b-set. Thus, $-n$
 481 refers to the dialogue act that is the A-part of the class. In the case where the
 482 labelled dialogue act is not an A-part or a B-part of an adjacency pair, one or both of
 483 the A-class and the B-class can be empty ($\langle 0, \{\} \rangle$).

484 The proposed encoding makes patterns of disagreements between annotators
 485 directly visible. For example, (1) if one annotator marks the dialogue act 43 as an
 486 A-part of two adjacency pairs with B-parts 44 and 45, respectively, and the dialogue
 487 act 45 as an A-part of an adjacency pair with the B-part 47, and (2) the other
 488 annotator marks the dialogue act 44 as an A-part of an adjacency pair with the
 489 B-part 45 and the dialogue act 45 as an A-part of two adjacency pairs with B-parts 46
 490 and 47, respectively, then the dialogue acts will be labelled as presented in Table 7.
 491 Fig. 2 illustrates the relation between the context of the dialogue act 45 and the AP
 492 label that encodes this context.

493 Encoding context in this way enables us to estimate *for each dialogue act* to what
 494 extent annotators agree on relating that dialogue act with surrounding dialogue acts
 495 in several ways: (1) as being an A-part related to a number of B-parts, (2) as being a
 496 B-part related to other B-parts with the same A-part and (3) not being related at all.
 497 It is to be noted that the context can be encoded in different ways as well. For
 498 example, it is possible to label each dialogue act that is marked as an A-part with its
 499 b-set. In this way, the actual disagreement is estimated only over A-parts. As context
 500 labels are not assigned to dialogue acts marked as B-parts, these dialogue acts would
 501 always be considered as agreed.

502 Agreement on APs annotation is measured over those dialogue acts for which
 503 annotators agreed on their boundaries. For computing agreement between annota-
 504 tors we use Krippendorff's α measure. This measure allows the usage of an appro-

Table 7 An example of adjacency pairs annotation (C_1 and C_2 : original AP annotations; $C_1(1)$ and $C_2(1)$: AP labels)

DA	C_1	C_2	$C_1(1)$	$C_2(1)$
43	1a2a		$\langle 0, \{1,2\} \rangle \langle 0, \{\} \rangle$	$\langle 0, \{\} \rangle \langle 0, \{\} \rangle$
44	1b	1a	$\langle 0, \{\} \rangle \langle -1, \{1,2\} \rangle$	$\langle 0, \{1\} \rangle \langle 0, \{\} \rangle$
45	3a2b	2a3a1b	$\langle 0, \{2\} \rangle \langle -2, \{1,2\} \rangle$	$\langle 0, \{1,2\} \rangle \langle -1, \{1\} \rangle$
46		2b	$\langle 0, \{\} \rangle \langle 0, \{\} \rangle$	$\langle 0, \{\} \rangle \langle -1, \{1,2\} \rangle$
47	3b	3b	$\langle 0, \{\} \rangle \langle -2, \{2\} \rangle$	$\langle 0, \{\} \rangle \langle -2, \{1,2\} \rangle$

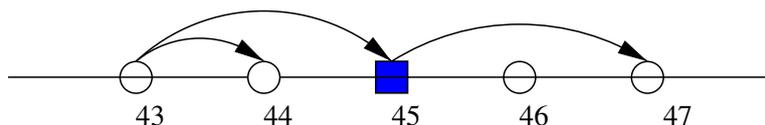


Fig. 2 A graphical representation of the context of dialogue act 45. The label that encodes this context is $\langle 0, \{2\} \rangle \mid \langle -2, \{1, 2\} \rangle$

505 appropriate user defined distance metric on the AP labels. For nominal categories, the
 506 usual α distance metric (δ) is a binary function: $\delta = 0$ if categories are equal,
 507 otherwise $\delta = 1$. We need to use a more refined distance metric, one that is sensitive
 508 for partial agreement of annotators on the context they assign to a dialogue act. The
 509 agreement on the contexts is translated to agreements on the corresponding
 510 A-classes and B-classes. When annotators disagree, their disagreement should be
 511 penalized based on the difference between classes.

512 The intuition is that similarity of two classes with the same A-part depends not
 513 only on the number of elements in the intersection of their b-sets, but also depends
 514 on the size of both sets. Therefore, we define a distance metric δ' that uses the
 515 following similarity measure on sets:⁶

$$\text{sim}(c_1, c_2) = \frac{2|c_1 \cap c_2|}{|c_1| + |c_2|} \quad (1)$$

518 The distance metric (δ') between the corresponding A-classes (or B-classes) of
 519 two AP labels is defined as:

$$\delta'(\langle -n_1, O_1 \rangle, \langle -n_2, O_2 \rangle) = 1, n_1 \neq n_2 \quad (2)$$

$$521 \quad \delta'(\langle -n, O_1 \rangle, \langle -n, O_2 \rangle) = 1 - \text{sim}(O_1, O_2) \quad (3)$$

523 The distance between two AP labels, $L_1 = A_1 \mid B_1$ and $L_2 = A_2 \mid B_2$, is defined as:

$$\delta_\lambda(L_1, L_2) = \lambda \cdot \delta'(A_1, A_2) + (1 - \lambda)\delta'(B_1, B_2), \quad (4)$$

526 where $\lambda \in [0, 1]$ is a factor that determines the relative contribution of the distance
 527 between the corresponding classes the labels consist of.

528 Applying $\delta_{0.5}$ to the data of exactly one R-C pair in each group gave the following
 529 results: M-R: $\alpha = 0.71$ ($N = 260$), B-E: $\alpha = 0.83$ ($N = 322$). The most frequently
 530 occurring disagreement is when one annotator marks a dialogue act with the empty
 531 label, the other annotator with a non-empty one. If annotators agreed that a dialogue
 532 act is an A-part of an adjacency pair, they mostly agreed, either partially or fully, on
 533 the B-set of this dialogue act. In most cases, the confusion between (1) an AP label
 534 with both A-class and B-class non-empty and (2) an AP label with one of the classes
 535 empty is related to the disagreement on the DA tags assigned by annotators. This
 536 concerns the confusion between (i) *Statement* and *Assessment/Appreciation* tags, (ii)
 537 *Statement* and *Response* tags and (iii) *Understanding Check* and *Information Request*
 538 tags.

⁶ The defined similarity measure is known as *Dice coefficient* (Manning & Schutze, 1999)

539 7.4 Reliability of gaze annotation

540 To evaluate reliability of gaze annotation, we first measured annotators agreement
541 on marking the changes in gazed targets. Then, we measured agreement on labelling
542 of time segments with gazed targets.

543 Marking the changes in gazed targets results in a segmentation of the time-line
544 into non-overlapping, continuous segments that cover the whole input. In other
545 words, the start time of a segment coincides with the end time of the segment that
546 precedes. A segment boundary indicates a change in gazed target.

547 The segmentation agreement is measured over all locations where any of the
548 annotators marked a segment boundary. The number of locations where both
549 annotators agree to some tolerance level is averaged over the total number of
550 locations marked as a boundary. A tolerance level is introduced because the gaze
551 annotation schema does not impose the requirement for a high precision on labelling
552 changes in the gazed targets. It is defined to adjust the difference in whether a
553 change is marked at the moment when the speaker starts changing the gaze direction
554 or at the moment when the new target has been reached. It also adjusts the differ-
555 ence in the reaction of the annotators to the observed changes. Empirical analysis of
556 the data shows that two points of the time-line can be considered equal with a
557 tolerance level of 0.85 s.

558 The agreement on locations where any coder marked a segment boundary is
559 80.40% ($N = 939$). Annotators mostly disagreed on marking the cases when a par-
560 ticipant briefly changes the gaze direction and then looks again at the previous
561 target. Annotators reached very good agreement on gaze labelling ($\kappa = 0.95$) mea-
562 sured over those segments where boundaries were agreed.

563 **8 Intra-annotator reliability**

564 Intra-annotator reliability measures whether the results of a single annotator re-
565 main consistent over time. We assessed intra-annotator reliability of dialogue act
566 and addressee annotation. One meeting from each data subset has been annotated
567 twice by each annotator in the DA group over a period of three months. The
568 results presented in Table 8 show that agreement on dialogue act annotation was
569 good for each annotator indicating intra-annotator consistency in applying the
570 dialogue act scheme. Furthermore, the results show that annotator R had a little
571 more difficulty with addressee annotation than other annotators who reached good
572 agreement.

Table 8 Intra-annotator agreement

Coder	Total	Agree	Segmentation	DA(κ)	ADD(κ)
E	110	104	94.54 %	0.83	0.88
B	107	104	97.20 %	0.89	0.81
M	73	64	87.67 %	0.81	0.87
R	77	72	93.51 %	0.85	0.76

573 **9 Discussions and conclusions**

574 We presented a multi-modal corpus of hand-annotated meeting dialogues that is
 575 designed for studying addressing behaviour in face-to-face conversations involving
 576 four participants. The corpus currently contains dialogue acts, addressees, adjacency
 577 pairs and gaze directions of meeting participants.

578 Annotators involved in the corpus design were able to reproduce the gaze
 579 labelling reliably. The annotations of dialogue acts and addressees were somewhat
 580 less reliable but still acceptable. Since there are only a few instances of subgroup
 581 addressing in the data and annotators failed to agree on them, the corpus cannot be
 582 used for exploring the patterns in addressing behaviour when a subgroup is ad-
 583 dressed. In this paper, we have also presented a new approach for measuring reli-
 584 ability of adjacency pairs annotation. The key of the approach is to represent AP
 585 annotated data as a form of categorical labelling in order to apply standard reliability
 586 metrics.

587 The corpus has already been used for the development of models for automatic
 588 addressee prediction (Jovanovic, op den Akker, & Nijholt, 2006). Apart from
 589 addressing, the corpus can be exploited for studying other interesting aspects of
 590 conversations involving more than two participants. As the NXT stand-off XML
 591 format enables an easy extension of the corpus with new annotation layers without
 592 influencing existing annotations, the corpus can be extended to include, for example,
 593 coding of turn-taking mechanisms which would enable studying this aspect of con-
 594 versational interaction independently as well as in relation to addressing.

595 **Acknowledgements** This work was partly supported by the European Union 6th FWP IST Inte-
 596 grated Project AMI (Augmented Multi-party Interaction, FP6-506811, publication AMI-160). We
 597 would like to thank Dennis Reidsma, Dennis Hof, Lynn Packwood and the annotators who were
 598 involved in the corpus development. We are grateful to Klaus Krippendorff for useful discussions
 599 about reliability metrics.

600 **References**

- 601 Bakx, I., van Turnhout, K., & Terken, J. (2003). Facial orientation during multi-party interaction
 602 with information Kiosks. In *Proceedings of the 9th IFIP TC13 International Conference on*
 603 *Human-Computer Interaction (INTERACT)*, (pp. 701–704). Zurich, Switzerland
 604 Burger, S., & Sloane, Z. (2004). The ISL meeting corpus: Categorical features of communicative
 605 group interactions. In *Proceedings of the NIST ICASSP Meeting Recognition Workshop*. Mon-
 606 treal, Canada
 607 Carletta, J. (1996). Assessing agreement on classification tasks: The kappa statistic. *Computational*
 608 *Linguistics*, 22(2), 249–254
 609 Carletta, J., Anderson, A. H., & Garrod, S. (2002). Seeing eye to eye: An account of grounding and
 610 understanding in work groups. *Cognitive Studies: Bulletin of the Japanese Cognitive Science*
 611 *Society*, 9(1), 1–20
 612 Carletta, J., Ashby, S., Bourban, S. M., Flynn, M. G., Hain, T., Kadlec, J., Karaiskos, V., Kraaij, W.,
 613 Kronenthal, M., Lathoud, G., Lincoln, M., Lisowska, A., McCowan, I., Post, W., Reidsma, D., &
 614 Wellner, P. (2006). The AMI meeting corpus: A pre-announcement. In S. Renals & S. Bengio
 615 (Eds.), *Machine learning for multimodal interaction, Second International Workshop, MLMI*
 616 *2005, Edinburgh, UK, July, 2005, revised selected papers*, vol 3869 of *Lecture Notes in Computer*
 617 *Science* (pp. 28–39). Springer-Verlag. ISBN 3-540-32549-2
 618 Carletta, J., Evert, S., Heid, U., Kilgour, J., Robertson, J., & Voormann, H. (2003). The NITE XML
 619 toolkit: Flexible annotation for multi-modal language data. *Behavior Research Methods,*
 620 *Instruments, and Computers*, 35(3), 353–363

- 621 Carletta, J., Isard, A., Isard, S., Kowtko, J., Doherty-Sneddon, G., & Anderson, A. (1997). The
 622 reliability of a dialogue structure coding scheme. *Computational Linguistics*, 23(1), 13–31
 623 Chen, L., Rose, R. T., Parrill, F., Han, X., Tu, J., Huang, Z., Harper, M., Quek, F., McNeill, D.,
 624 Tuttle, R., & Huang, T. (2006). VACE multimodal meeting corpus. In S. Renals & S. Bengio
 625 (Eds.), *Machine learning for multimodal interaction, Second International Workshop, MLMI*
 626 *2005, Edinburgh, UK, July, 2005, revised selected papers*, vol 3869 of *Lecture Notes in Computer*
 627 *Science* (pp. 40–51). Springer-Verlag. ISBN 3-540-32549-2
 628 Clark, H. H. & Carlson, T. B. (1992). Hearers and speech acts. In H. H. Clark (Ed.), *Arenas of*
 629 *Language Use* (pp. 205–247). University of Chicago Press and CSLI
 630 Clark, A., & Popescu-Belis, A. (2004). Multi-level dialogue act tags. In *Proceedings of the 5th*
 631 *SIGdial Workshop on Discourse and Dialogue* (pp. 163–170). Cambridge, MA
 632 Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological*
 633 *Measurement*, 20, 37–46
 634 Dhillon, R., Bhagat, S., Carvey, H., & Shriberg, E. (2004). Meeting recorder project: Dialogue act
 635 labeling guide. Technical report TR-04-002, International Computer Science Institute (ICSI),
 636 Berkeley, CA, USA
 637 Garofolo, J. S., Laprun, C. D., Michel, M., Stanford, V. M., & Tabassi, E. (2004) The NIST meeting
 638 room pilot corpus. In *Proceedings of the 4th International Conference on Language Resources*
 639 *and Evaluation (LREC)* (pp. 1411–1414). Lisbon, Portugal
 640 Goffman, E. (1981). Footing. In Goffman, E., (Ed.), *Forms of talk* (pp. 124–159). Philadelphia, PA:
 641 University of Pennsylvania Press
 642 Goodwin, C. (1981). *Conversational organization: Interaction between speakers and hearers*. New
 643 York: Academic Press
 644 Hirschberg, J., & Nakatani, C. (1996). A prosodic analysis of discourse segments in direction-giving
 645 monologues. In *In Proceedings of the 34th Annual Meeting of the Association for Computational*
 646 *Linguistics (ACL)* (pp 286–293). Santa Cruz, California
 647 Janin, A., Ang, J., Bhagat, S., Dhillon, R., Edwards, J., Macias-Guarasa, J., Morgan, N., Peskin, B.,
 648 Shriberg, E., Stolcke, A., Wooters, C., & Wrede, B. (2004). The ICSI meeting project: Resources
 649 and research. In *Proceedings of the NIST ICASSP Meeting Recognition Workshop*. Montreal,
 650 Canada
 651 Jovanovic, N., & op den Akker, R. (2004). Towards automatic addressee identification in multi-party
 652 dialogues. In *Proceedings of the 5th SIGdial Workshop on Discourse and Dialogue* (pp. 89–92).
 653 Cambridge, MA
 654 Jovanovic, N., op den Akker, R., & Nijholt, A. (2006). Addressee identification in face-to-face
 655 meetings. In *Proceeding of the 11th Conference of the European Chapter of the Association for*
 656 *Computational Linguistics (EACL)* (pp. 169–176). Trento, Italy
 657 Jurafsky, D., Shriberg, L., & Biasca, D. (1997). Switchboard SWBD-DAMSL Shallow-Discourse-
 658 Function Annotation Coders Manual, Draft 13. Technical report 97–02, University of Colorado,
 659 The Institute of Cognitive Science, Boulder, CO
 660 Katzenmaier, M., Stiefelhagen, R., & Schultz, T. (2004). Identifying the addressee in human–human–
 661 robot interactions based on head pose and speech. In *Proceedings of the International Confer-*
 662 *ence on Multimodal Interfaces (ICMI)* (pp. 144–151). State College, PA
 663 Krippendorff, K. (1980). *Content analysis: An introduction to its methodology*. Sage Publications,
 664 Beverly Hills, CA
 665 Krippendorff, K. (2004). Reliability in content analysis: Some common misconceptions and rec-
 666 ommendations. *Human Communication Research*, 30(3), 411–433
 667 Manning, C., & Schütze, H. (1999). *Foundations of statistical natural language processing*. The MIT
 668 Press, Cambridge, MA
 669 Moore, D. (2002). The IDIAP smart meeting room. Technical report IDIAP-COM-07, IDIAP,
 670 Martigny, Switzerland
 671 Otsuka, K., Takemae, Y., Yamato, J., & Murase, H. (2005). A probabilistic inference of multiparty-
 672 conversation structure based on Markov-switching models of gaze patterns, head directions, and
 673 utterances. In *Proceedings of the International Conference on Multimodal Interfaces (ICMI)* (pp.
 674 191–198). Trento, Italy
 675 Passonneau, R. (2004). Computing reliability for coreference annotation. In *Proceedings of the 4th*
 676 *International Conference on Language Resources and Evaluation (LREC)* (pp. 1503–1506).
 677 Lisbon, Portugal
 678 Passonneau, R. & Litman, D. (1997). Discourse segmentation by human and automated means.
 679 *Computational Linguistics*, 23(1), 103–139

- 680 Reidsma, D., Hofs, D., & Jovanovic, N. (2005). A presentation of a set of new annotation tools based
681 on the NXT API. In *Proceedings of Measuring Behavior* (pp. 512–513). Wageningen, The
682 Netherlands
- 683 Shriberg, E., Dhillon, R., Bhagat, S., Ang, J., & Carvey, H. (2004). The ICSI meeting recorder dialog
684 act (MRDA) corpus. In *Proceedings of the 5th SIGdial Workshop on Discourse and Dialogue*
685 (pp. 97–100). Boston, USA
- 686 Traum, D. (2004). Issues in multi-party dialogues. In F. Dignum (Ed.), *Advances in agent commu-*
687 *nication* (pp. 201–211). Springer-Verlag LNCS
- 688 van Turnhout, K., Terken, J., Bakx, I., & Eggen, B. (2005). Identifying the intended addressee in
689 mixed human–human and human–computer interaction from non-verbal features. In *Proceed-*
690 *ings of the International Conference on Multimodal Interfaces (ICMI)* (pp.175–182). Trento, Italy
691

UNCORRECTED PROOF