1    **Computational approaches to the pragmatics problem**

2    **Abstract**

3    Unlike many aspects of human language, pragmatics involves a systematic many-to-many

4    mapping between form and meaning.  This renders the computational problems of encoding

5    and decoding meaning especially challenging, both for humans in normal conversation and

6    for artificial dialogue systems that need to understand their users' input.  A particularly

7    striking example of this difficulty is the recognition of speech act or dialogue act types.  In

8    this review, we discuss why this is a problem, and why its solution is potentially relevant both

9    for our understanding of human interaction and for the implementation of artificial systems.

10    We examine some of the theoretical and practical attempts that have been made to overcome

11    this problem, and consider how the field might develop in the near future.

12    **Introduction**

13    What constitutes human communication?  One possible answer is to claim that it requires a

14    sender and a recipient, and that information is encoded by the sender, transmitted, and

15    decoded by the recipient.  This concept of communication was famously formalised by

16    Shannon (1948).  However, Grice (1957) argued that communication between people was

17    also characterised by the process of intention recognition.  Specifically, he identified the

18    notion of "non-natural meaning", in which sense a speaker "means" something if, firstly, they

19    intend to induce a belief in the hearer as a consequence of that utterance, and secondly, they

20    intend for this to happen as a result of the hearer recognising the intention (conveyed by the

21    utterance) to bring about this belief.  For instance, a speaker who says "Please sit down"

22    intends for the hearer to sit down, and for this to occur because the hearer recognises that this

23    is what the speaker wants to convey by these words.  From this perspective, as Levinson

24    (1983: 15) puts it, "communication involves the notions of intention and agency".

25    Grice's view of inter-personal communication has been enormously influential in linguistic

26    pragmatics and related fields.  A striking point of contrast with the Shannon model, as Grice

27    himself immediately noted (1957: 387), is that the intentional view of communication admits

28    the possibility of indeterminacy.  On the Gricean view, it is possible for the same signal to

29    correspond to different intentions, in which case it is necessary to appeal to context in order

30    to understand what the speaker actually intends on this particular occasion.  Shannon,

31    conversely, adopts a model in which encoding and decoding of a signal are one-to-one

32    mapping processes, and in which context and the mental state of the sender are irrelevant to

33    the recipient's understanding of the message.

34    It seems undeniable that human communication does indeed have the systematic ambiguity

35    that Grice posits, whether this is a consequence of the polysemy of words or the multi-

36    functional nature of various actions: Grice's own examples are the word 'pump' and the

37    action of putting one's hand in a pocket.  So clearly some elaboration of the Shannon model

38    is called for.  And intuitively, it seems credible that the goal of the hearer is to understand the

39    intention of the speaker, as Grice argues.  However, given that many different intentions may

40    be realised by the same signal, the task of recovering the speaker's intention given a signal is

41    logically intractable (Levinson 1995: 231) – there is not enough information in the signal to

42    tell the hearer, precisely and unambiguously, what the intention was.  In order for the

43    Gricean, intentional analysis of communication to be tenable, we therefore need to be able to

44    explain how hearers are so often successful in solving this 'pragmatics problem', and

45    understanding what intention underlies the speaker's choice of utterance. Given the

46    ramifications of this model for our understanding of human interaction, foundational

47    questions about the validity of the model are of substantial theoretical importance.

48    In this paper, we focus on a particular subcase of the pragmatics problem that has attracted

49    widespread interest from philosophers of language and builders of computational systems

50    alike: namely, the way in which we identify dialogue act types.  The following section

51    discusses why this is an important issue for both human-human interactions and for artificial

52    spoken dialogue systems.  We then outline some of the most productive linguistic and

53    computational attempts to address this issue.  We conclude by considering how these

54    methods might usefully be synthesised into a coherent interdisciplinary approach to dialogue

55    act type recognition.

56    **Dialogue act recognition in interaction**

57    As pointed out by Austin (1962), our use of language does not just consist of asserting

58    propositions.  More broadly, we perform "speech acts".  That is to say, we "do things with

59    words" – we use utterances to achieve particular effects.  We may request an action,

60    acknowledge a request, ask for information, and so on.  From this perspective, we can see

61    language as a tool that we can use in order to accomplish things that we would not be able to

62    accomplish by other forms of physical action.  We can also analyse individual instances of

63    language use as social actions that are performed in order to elicit specific responses, which

64    might involve obtaining information or causing interlocutors to act upon the physical world

65    in particular ways.

66    The usefulness of linguistic acts in enabling specific social accomplishments cannot easily be

67    treated in terms of truth conditions: it doesn't generally make sense to describe a request as

68    "true" or "false", for instance.  Austin introduced the notion of "illocutionary act" to describe

69    this kind of function, a notion which was later elaborated by Searle (1975).  Although this

70    research tradition is referred to as speech act theory, here we will use the term "dialogue act"

71    rather than "speech act" to emphasise that the relevant actions may be achieved by other

72    means than through speech (for instance, gesture, eye-gaze, and so on).  There is little

73    consensus as to what constitutes an appropriate typology of dialogue acts, but we might

74  distinguish dialogue act types by appeal to a notion like "what kind of response is

75  appropriate".

76  In order for the speaker's dialogue act to be effective, it is generally necessary (under the

77  Gricean assumptions discussed above) for the hearer correctly to identify it, as without doing

78  so, it is impossible for the hearer to respond in such a way as to satisfy the speaker's goals.

79  However, as has long been observed, this is not a straightforward matter. Consider for

80  example the potential dialogue act of 'asking a question'. Nearly all human languages

81  possess the interrogative sentence-type, which is usually distinguished from the declarative

82  by some complex of morphosyntactic and intonational factors. It is tempting to assume that

83  the task of recognising the dialogue act 'asking a question' is reducible to that of recognising

84  an interrogative sentence. But this is simply not true: a formally declarative sentence may

85  perform a questioning function ("You'll let me know"), and a formally interrogative sentence

86  may function as a request ("Could you close the window?") Indeed, interrogative forms can

87  easily be ambiguous between various dialogue act types depending on context ("Can you

88  come?" could be a question, a request or an invitation). Moreover, the notion of 'asking a

89  question' might not even constitute a single coherent dialogue act type: it might include such

90  distinct dialogue acts as 'asking a polar question', 'asking a wh- question', 'asking a check

91  question', and so on. If these need to be distinguished, that clearly cannot rely on appeal to

92  the sentence-type alone, which is typically the same (interrogative) in all cases.

93  The recognition of dialogue act types can thus be seen as a specific case of intention

94  recognition, and one that succumbs to the pragmatics problem: given that several different

95  intentions may be expressed by the same form, how can the hearer locate the right one? And

96  just as we ask this question for human interactors, so we can ask it for artificial systems, and

97  in particular spoken dialogue systems – that is, systems that are designed to converse with

98  humans. To get computers to understand one another, we can program them to communicate

99    unambiguously: but the ultimate goal for a spoken dialogue system is to be able to

100   accommodate all the ambiguity and uncertainty of normal human discourse.  (In practice,

101   humans tend to adjust their choice of words to match the abilities of artificial systems (see

102   Branigan et al. 2011), but ideally this would not be necessary.)  Moreover, the system must

103   understand what the speaker is actually trying to achieve, rather than merely formalising the

104   content of the speaker's utterance in some way.  This kind of understanding also proves

105   useful in enabling the system correctly to identify individual words that would otherwise not

106   have been correctly parsed (Stolcke et al. 2000, Taylor et al. 2000).  In order to allow systems

107   of this kind to approach human performance levels, it would be helpful to have a fuller and

108   clearer account of how humans actually recognise dialogue act types.

109   A growing body of evidence underscores the impressive nature of human performance in this

110   particular domain.  Our own experience suggests that competent language users are able

111   correctly to identify the intended dialogue act in the vast majority of cases, as shown by the

112   appropriateness of their responses.  For instance, a hearer asked "Could you pass the salt?"

113   will usually do so, unless they deliberately choose to misinterpret the speaker's intention and

114   merely say "Yes".  In cases such as this, the formal ambiguity of the utterance is not

115   necessarily noticed by the dialogue participants, unless it is pointed out by a response that is

116   inappropriate to the speaker's actual intention.

117   The success of dialogic communication speaks to the accuracy of the conclusions arrived at

118   by hearers about the speakers' intentions.  Experimental work suggests that hearers are not

119   only accurate but also remarkably fast in identifying the speaker's intention in ongoing

120   utterances.  Relevant evidence here comes from turn-taking.  De Ruiter, Mitterer and Enfield

121   (2006) demonstrated that, in spontaneous Dutch conversation, almost half of the new

122   conversational turns started within 250ms (either way) of the end of the current turn.  Stivers

123   et al. (2009) generalised this result to a typologically mixed sample of 10 languages: for each

124     language, the mean duration of the gap between turns was less than half a second, the

125     "fastest" being Japanese with a mean gap of just 7ms.  This supports the observation by

126     Levinson (1995: 237) that a half-second delay in responding can (in English) be interpreted

127     as conveying some pragmatic effect (in that case, the impossibility of the hearer responding

128     'yes' to a question).

129     Recent work on dialogue act recognition (Gisladottir et al. 2012) demonstrates directly that

130     hearers are able accurately to identify dialogue acts off-line.  Hence, given the content of a

131     speaker's turn (and awareness of the contrast), it should not be a problem for the hearer to

132     identify the speaker's dialogue act type.  However, it seems profoundly implausible that this

133     could happen in the gaps between turns documented by Stivers et al. (2009).  In the first

134     place, many of the languages they test exhibit frequent overlap in turn transitions, which

135     indicates that hearers cannot be waiting for the speaker's turn to be complete before they start

136     planning their own conversational response.  In the second place, research on utterance

137     planning (for instance, Brown-Schmidt and Tanenhaus 2006) appears to indicate that even a

138     latency of 500ms would not be enough for the hearer even to formulate a response *ab initio*.

139     Given that the responses are usually faster than this, usually pertinent, and usually conform to

140     the dialogic strictures laid down by the speaker (for instance, a question will be met with an

141     answer), this strongly suggests that the hearer must often be aware of the nature of the

142     speaker's dialogue act before it is complete.

143     In a similar vein, we might interpret the nature of back-channel responses (Yngve 1970) as

144     evidence that the hearer can identify aspects of the speaker's communicative intention

145     incrementally and on-line.  Back-channel responses are utterances by the hearer that are not

146     attempts to initiate a turn.  Schegloff (1982) refers to a subset of these as "continuers", on the

147     basis that they serve to assure the speaker of the hearer's attention and indicate that the turn

148     can continue.  Various utterances can fulfil this function, among them "uh-huh" and "yeah".

149 However, it appears likely that the appropriate choice of back-channel response depends to a

150 certain extent upon the dialogue act being performed by the speaker – for instance, "yeah"

151 would not be an appropriate back-channel if the speaker is formulating a request, unless the

152 hearer intends to comply (cf. Schegloff 1993: 107). If this intuition is correct, it further

153 suggests that hearers may be able to access information about the speaker's dialogue act type

154 from relatively early in the utterance.

155 In sum, there appears to be quite convincing evidence that human dialogue participants are

156 able to draw rich inferences about dialogue act types from very early on in a dialogue turn.

157 In the following section, we examine some approaches to explaining how this process might

158 take place.

159 **Approaches to dialogue act recognition**

160 A linguistic approach to dialogue act recognition was offered by Gazdar (1981), who

161 formulated the Literal Meaning Hypothesis. According to this account, every utterance

162 possesses some kind of illocutionary force that is built into its surface form. Declaratives are

163 used to make statements, interrogatives to question, imperatives to order or request, and verbs

164 such as "promise", "deny" and so on (performatives, in Austin's terms) are used to

165 accomplish whichever function their verb specifies. However, as discussed earlier, utterances

166 are frequently used to accomplish other discourse functions than their surface form would

167 suggest, and the same utterance may be used for multiple functions. So at the very least we

168 need to supplement the Literal Meaning Hypothesis with some mechanism that enables

169 hearers to calculate the alternative non-literal or "indirect" meanings that may arise.

170 One possibility is to appeal to traditional pragmatic notions of cooperativity and, in

171 particular, relevance. Gordon and Lakoff (1971) suggest that reanalysis occurs when the

172 hearer realises that the surface meaning of the utterance is inappropriate given the context.

173     For instance, a speaker asking "Could you pass the salt?" typically knows that the hearer is

174     able to do so, and the hearer can infer from this that the purpose of the utterance is not to

175     enquire as to their salt-passing capabilities. For the utterance not to be a waste of effort,

176     therefore, there must be some other purpose to it.  Searle (1975) tells a slightly different

177     story: on his account, the 'natural' answer to the question "Could you pass the salt?" (namely:

178     yes, the hearer could do so) must be relevant to the speaker.  A possible reason for this is that

179     the speaker wants the salt; and the hearer, being cooperative, should therefore pass the salt to

180     the speaker, without an explicit request being necessary.

181     Can we, however, reconcile this kind of account with the data on turn-taking discussed

182     above?  Timing presents a serious problem.  Both versions of the pragmatic account take as

183     their starting point the realisation that the literal meaning of the utterance is in some way

184     inadequate given the conversational context, and has to be enriched.  However, if the

185     reasoning in the previous section is correct, this process has to begin before the utterance is

186     complete.  The problem is, how can the hearer determine that the literal meaning of the

187     utterance is inadequate before knowing what the utterance is?  A sentence beginning "Could

188     you…", or even "Could you pass…", could certainly be a genuine question that was not a

189     request ("Could you pass for 21?").  More generally, we might observe that almost any

190     sentence beginning "Could you…" might conceivably be used either as a question or as a

191     request, and for many such cases, it is easy to imagine contexts in which either use might be

192     intended ("Could you teach a course in psycholinguistics?")  In order to know that "Could

193     you pass the salt?" cannot (normally) be intended as a question about the hearer's

194     capabilities, the hearer must identify the meaning of the sentence and realise that the speaker

195     knows the answer to the question that is ostensibly being posed.  This is completely

196     reasonable *post hoc*, but as an account of online reasoning it doesn't appear to give the hearer

197     enough time to formulate their response.

198 One conceivable way of rescuing this account is to propose that the hearer in fact guesses

199 how the sentence will end, and reasons on the basis of that guess, thus being able to draw the

200 inferences discussed above before the end of the speaker's turn. After all, Sacks, Schegloff

201 and Jefferson (1974) proposed that hearers anticipate the end of speakers' turns in order to

202 achieve smooth transitions; and Magyari and De Ruiter (2012) provide evidence that the

203 accuracy of this anticipation is correlated with the rapidity of turn transition. However, as an

204 account of dialogue act type recognition, this explanation is in danger of becoming circular: a

205 hearer may well guess that the sentence "Could you pass…" concludes with the words "the

206 salt", but this continuation only makes sense if the utterance is a request, whereas by

207 hypothesis the hearer currently takes the utterance to be a question. To put it another way:

208 intuitively, we might expect the words "the salt" because we guess that the speaker wants the

209 salt passed to them. But how did we guess that the speaker wanted something passed to

210 them? Presumably because "Could you pass…" tends to signal that this is the case,

211 notwithstanding that it is formally part of an interrogative sentence-form.

212 An alternative approach, foreshadowed by Levinson (1983), is to dispense with the Literal

213 Meaning Hypothesis, and instead treat the identification of dialogue act type as a puzzle to be

214 solved by any means available. That is not to propose that the hearer ignores the sentence-

215 type: that might be a valuable clue to the dialogue act type. However, according to Levinson,

216 most speech acts are indirect, in the sense that they do not correspond to the surface form of

217 the sentence. Fortunately, there are many other forms of information that might be helpful to

218 the hearer. Within the speech signal itself, other indications of the likely dialogue act type

219 are present. These include the prosody, as discussed by Bolinger (1964) and extensively

220 explored by Shriberg et al. (1998) among many others. It is also likely that specific lexical

221 choices are strongly associated with particular dialogue acts. For instance, "I want you to…"

222 strongly suggests that the current sentence has the character of a request, even though the

223     sentence-type is purely declarative. Even more generally, the use of "please" seems typically

224     to mark a request whether it is appended to a declarative ("The door should be closed,

225     please"), imperative ("Close the door, please") or interrogative ("Could you close the door,

226     please?") sentence-type.

227     At a higher level, there are considerations deriving from the structure of dialogue, as studied

228     within the research tradition of conversation analysis: for instance, the idea of adjacency pairs

229     (Schegloff and Sacks 1973). If the preceding dialogue turn was a question, the current turn is

230     likely to be an answer, even if its form suggests otherwise. If the previous turn was an offer,

231     the current turn is likely to involve accepting or declining that offer. Thus, when we

232     encounter the first turn of an adjacency pair, we might (with some degree of confidence)

233     expect that the second turn of that pair will follow. Adjacency pairs can also have non-

234     linguistic constituents, as argued by Schegloff (1968). Clark (2004) originates the notion of

235     'projective pair' to cover cases where a non-linguistic communicative act such as a gesture

236     serves to trigger a particular kind of communicative act in response. He later argues (Clark

237     2012) that we can identify wordless exchanges that are analysable as question-answer

238     sequences. At a still higher level of discourse organisation, an awareness of the overarching

239     purpose of the dialogue and of the participants' roles in it might help a hearer disambiguate

240     dialogue act types. In a restaurant, for instance, if a customer states the names of dishes, this

241     is likely to be a request; if a waiter does so, it is more likely to be an offer (or effectively a

242     multiple-choice question).

243     Computational implementations of dialogue act recognition have predominantly adopted this

244     kind of permissive, inclusive approach, in which all available forms of information are used

245     to make the relevant decisions. This cue-based approach essentially dispenses with the

246     assumption of literal meaning elaborated by the kind of stepwise inference discussed earlier,

247     although that approach has also been explored computationally (from Perrault and Allen 1980

248     to Allen et al. 2007). The role of the cue-based model is simply to identify which dialogue

249     act is instantiated by a given utterance, appealing as necessary to lexical, syntactic, prosodic

250     and conversational-structural factors, among others.

251     It would perhaps be fair to say that cue-based implementations are primarily focused on

252     improving the performance of systems, rather than necessarily providing insights into the

253     process of dialogue act recognition *per se*. However, the models are linguistically informed,

254     in important respects. They are trained on labelled corpora, from which they can learn the

255     strengths of association between specific signals and specific dialogue acts. The choice of

256     signals may, and typically does, reflect empirically-determined findings as to which aspects

257     of the utterance are likely to constitute informative cues. Identifying potentially useful

258     signals is a non-trivial problem in domains such as prosody, where it is unclear precisely

259     what properties of the acoustic pattern have informational value (see for example Rangarajan

260     Sridhar, Bangalore and Narayanan 2009).

261     Although traditional linguistics and computational modelling approaches find common cause

262     when it comes to identifying signals, the customary meaning of 'dialogue act' varies

263     significantly between the two traditions. As Thomson (2010: 10) puts it, "In the traditional

264     definitions of both speech and dialogue acts, the semantic information is completely

265     separated from the act". That is to say, the utterance "Could you pass the salt?" is an instance

266     of a dialogue act type like REQUEST rather than one like REQUEST-SALT. From a linguistic

267     point of view, the motivation for this is fairly clear: the notion of dialogue act type captures

268     the idea that there are commonalities between all forms of REQUEST, regardless of what is

269     being requested. However, from a dialogue systems standpoint, this is not necessarily an

270     advantage. If the goal of the system is to fulfil the user's request, then merely identifying the

271     utterance as 'some kind of request' is not helpful: it does not enable the system to formulate a

272     response, as this response will depend upon what is being requested. Unless the system has

273     an abstract understanding of how to fulfil generic requests, the 'type' level of dialogue acts is

274     not useful here.

275     Moreover, by dispensing with the 'type' level, it may be possible for a system to identify

276     dialogue acts more efficiently than a human could.  Consider the case of a robot receptionist

277     (as implemented, for example, by Paek and Horvitz 2000).  Suppose that John Smith is an

278     employee at the company and that the robot is programmed with only one action that relates

279     to John Smith, namely putting a call through to him.  Confronted with the input "Could you

280     call John Smith?", the robot can use the words "John Smith" as a cue to the action it should

281     take, and thus use the name as evidence that it should put a call through.  A more capable

282     robot, just like a human, would be disadvantaged here, because if it could take various

283     different actions with respect to John Smith, recognising the name would not suffice to

284     identify which one should be performed.  Of course, the simple robot may misidentify

285     dialogue acts that are outside its knowledge base ("My name is John Smith"), but it has no

286     problem using lexical cues to choose among its limited repertoire of abilities.

287     The question arises of whether the traditional notion of dialogue act type is at all helpful for

288     implementations of spoken dialogue systems.  Traum (1999) considers this point, coming to

289     the conclusion that dialogue act types may not be strictly necessary but are potentially useful

290     as an intermediate step in communication planning.  The practice of identifying dialogue acts

291     at a finer level of granularity (REQUEST-SALT, CALL-JOHN-SMITH) certainly has implications

292     for the scalability of dialogue systems, as the number of distinct dialogue acts increases

293     drastically as the coverage of the system expands to multiple conversational domains

294     (whereas, by hypothesis, the number of dialogue act types is relatively small even for the

295     whole of human interaction).  This becomes especially pertinent when we consider

296     statistically-driven dialogue systems of the kind surveyed by Young et al. (2013).  These

297     models use the approach named POMDP (partially observable Markov decision processes)

298    and treat dialogue as a Markov process, in which transitions between dialogue states are

299    modelled probabilistically. Even within a small domain, it is impractical to track dialogue

300    state fully in such a model; for a general spoken dialogue system, the resulting state space

301    would be intractably large (Young et al. 2010: 152).

302    In particular, a domain-general system that identified highly specific dialogue acts would

303    necessarily have to incorporate thousands of distinct dialogue acts.  Consider the receptionist

304    scenario: a person entering the building might request the receptionist to make a call to any

305    individual in the building, using the form of words "Could you call X?"  A system that treats

306    every such request completely separately, depending on the identity of X, could not make

307    useful generalisations across this set of requests.  For instance, if the name of X is mumbled

308    or unfamiliar, it will not be able to respond "Sorry, who?" unless it identifies the utterance as

309    a request: it could only announce its inability to respond to the request as a whole, which

310    might prompt futile reformulations ("I would like to talk to X").  That is, although such a

311    system might be very efficient at learning the mappings between specific strings and specific

312    tasks, it will struggle to generalise these mappings in any remotely human-like way.

313    Similarly, if it is possible to make generalisations about dialogue act sequences (e.g.

314    question-answer, apology-acceptance, check-confirmation, and so on), these generalisations

315    will not be as evident when the coarse-grained dialogue act types are broken down into fine-

316    grained ones.[i]  If each particular kind of apology must be separately associated with a kind of

317    acceptance, a large volume of data may be required for the pattern to be learnt by the system

318    across all pertinent occasions.

319    However, this observation, like Traum's (1999) discussion, relates primarily to the operation

320    of relatively complex dialogue agents with sophisticated 'mental' states.  For simpler

321    systems, dialogue act type recognition in the traditional sense is clearly less useful: in the

322    limiting case, if a system does nothing but (attempt to) satisfy requests, coding a module to

323    identify every input as a REQUEST is clearly not going to add anything to the system's

324    efficacy. What the system needs to do is to identify what is being requested: only then can it

325    initiate the appropriate response behaviour. Unless the system has a generic handling

326    procedure for requests, it cannot benefit from the inclusion of this additional level of analysis.

327    By contrast, systems that actually attempt to emulate human behaviour have the potential to

328    benefit from including a dialogue act level. A recent example of such a system the virtual

329    agent implemented by DeVault, Sagae and Traum (2011), designed to help soldiers practice

330    negotiation skills. The agent uses a natural language understanding module to convert the

331    content of the human user's utterance into a semantic frame representation. One of the

332    attributes within this semantic frame is 'speech act type', so the artificial agent could be said

333    to be calculating and exploiting information about the human speaker's purpose. Moreover,

334    the agent can be configured to guess the content of the semantic frame based on partial

335    utterances, thus effectively engaging in incremental identification of dialogue act type.

336    The catch, however, is that semantic frames are treated as atomic within DeVault et al.'s

337    model, even though they are decomposable in principle. That is, their model postulates a

338    finite set of semantic frames and aims to identify, based on the user's utterance, which one is

339    currently being instantiated by the speaker. Each semantic frame happens to have an attribute

340    that is called 'speech act type', but this specific attribute is not exploited in any way:

341    responses are selected based upon the entire semantic frame that is identified. There is, in

342    effect, no commonality between semantic frames that contain the same speech act type. The

343    decision to treat semantic frames as atomic reflects a deliberate simplification, justified on the

344    basis that it does not impair performance on the constrained domain in which the model

345    operates. However, for the model to be scalable, some form of non-atomic approach would

346    be necessary, which might involve the exploitation of dialogue act types in a more traditional

347    way.

**Towards an interdisciplinary perspective on dialogue act recognition**

As the above discussion indicates, insights from theoretical linguistics have already been brought to bear productively upon the implementation of artificial spoken dialogue systems. However, our psycholinguistic questions about the process of dialogue act recognition and behaviours such as turn-taking are not directly addressed by this practical computational work. Most of the computational work has so far taken place in highly constrained domains, while we are interested in the full sweep of human communicative interaction. Moreover, computational approaches have predominantly attempted to achieve effective behaviour by any means necessary, but this may involve means that are not available to, or not exploited by, human interactors. For instance, computational models do not have the working memory limitations of humans, and can in principle use probabilistic cues that are outside of humans' knowledge (for instance, because they involve relations over too long a distance, or patterns that humans are not disposed to spot). They do not have the experiential limitations of humans: they can be trained on larger corpora than a human would ever experience. And they typically do not operate under the same time pressure as humans, assuming that they can initiate responses faster than humans can program their own motor functions.

Nevertheless, the application of these methods already gives us a useful insight into what might work, and which theoretical ideas add value in a practical context. For instance, Young et al. (2010) use a bigram model of dialogue act type, which is informed by the work of Schegloff and Sacks (1973) on adjacency pairs, to help identify the user's response to their artificial agent's questions. DeVault, Sagae and Traum (2011) use a rich array of lexical cues from the input string to support the semantic classification of the user's utterances. As discussed earlier, this latter model can also be made to operate incrementally, while the bigram approach of Young et al. also informs us about the likely nature of the current dialogue act before it is complete. It would seem quite conceivable to take these

373  mechanisms, and others like them, equip them with a notion of dialogue act type, and use

374  them to classify utterances in natural human-human interactions.

375  Furthermore, if we are interested in learning about how humans treat dialogue acts, we can

376  calibrate such a model against experimentally verified human behaviour. That is, we can

377  eliminate factors that do not appear to influence human performance, just as we can introduce

378  additional factors that are posited to play a role in humans' classification of dialogue act

379  types. And we can similarly adjust the candidate set of dialogue act types, in accordance with

380  competing theoretical proposals. The ultimate goal of such a programme might be to

381  establish a set of dialogue acts that are descriptively adequate as a characterisation of the

382  components of human dialogic interaction, and which are identifiable sufficiently quickly by

383  appeal only to utterance and contextual properties that humans are known to respond to.

384  Working in the opposite direction, it is also conceivable that a fully adequate theory of

385  dialogue acts could be very useful in the development of domain-general spoken dialogue

386  systems. It is, of course, clear that this is not a substitute for a comprehensive system of

387  semantics – a system that reliably gives 'answers' that don't relate to the question will not

388  survive scrutiny – but it may turn out to be a necessary component if dialogue systems are to

389  behave in a credibly human-like fashion (and thus to allow their human users to behave

390  normally with them). It may also transpire that the use of dialogue acts results in systems

391  being more compact and efficient than would otherwise be the case, just as the analysis of

392  dialogue reveals order in what might otherwise be the limitless variety of human-human

393  interaction.

394  **Endnotes**

395  [1] The potential to draw useful generalisations will depend not only on defining dialogue act types at

396  the right level of granularity, but actually choosing an appropriate set of specific dialogue act types

397 with which to populate the model.  For reasons of space we cannot substantively address this issue

398 here.  See Král and Cerisara (2010) for a discussion of some specific candidate 'tag-sets' for dialogue

399 acts.

400 **References**

401 Allen, J. F., Chambers, N., Ferguson, G., Galescu, L., Jung, H., Swift, M., & Taysom, W.

402 (2007). PLOW: A collaborative task learning agent. National Conference on Artificial

403 Intelligence (AAAI), Vancouver, BC.

404 Austin, J. L. (1962). *How to Do Things with Words*. Oxford: Clarendon Press.

405 Bolinger, D. L. (1964). Intonation across languages. In J. P. Greenberg, C. A. Ferguson & E.

406 A. Moravcsik (eds.), *Universals of Human Language Phonology, vol. 2*. Stanford: Stanford

407 University Press. 471-524.

408 Branigan, H. P., Pickering, M. J., Pearson, J., McLean, J. F., and Brown, A. (2011).  The role

409 of beliefs in lexical alignment: evidence from dialogs with humans and computers.

410 *Cognition*, 121: 41-57.

411 Brown-Schmidt, S., & Tanenhaus, M. K. (2006). Watching the eyes when talking about size:

412 an investigation of message formulation and utterance planning. *Journal of Memory and*

413 *Language*, 54: 592-609.

414 Clark, H. H. (2004). Pragmatics of language performance. In L. R. Horn & G. Ward (eds.),

415 *Handbook of Pragmatics*. Oxford: Blackwell. 365-382.

416 Clark, H. H. (2012). Wordless questions, wordless answers. In J. P. de Ruiter (ed.),

417 *Questions: Formal, functional and interactional perspectives*. Cambridge: Cambridge

418 University Press. 81-100.

419     De Ruiter, J. P., Mitterer, H., & Enfield, N. J. (2006). Predicting the end of a speaker's turn: a

420     cognitive cornerstone of conversation. *Language*, 82: 515-535.

421     DeVault, D., Sagae, K., & Traum, D. (2011). Incremental interpretation and prediction of

422     utterance meaning for interactive dialogue. *Dialogue and Discourse* 2: 143-170.

423     Gazdar, G. (1981). Speech act assignment. In A. K. Joshi, B. L. Webber & I. A. Sag (eds.),

424     *Elements of Discourse Understanding*. Cambridge: Cambridge University Press. 64-83.

425     Gisladottir, R. S., Chwilla, D. J., Schriefers, H., & Levinson, S. C. (2012). Speech act

426     recognition in conversation: experimental evidence. In N. Miyake, D. Peebles, & R. P.

427     Cooper (Eds.), Proceedings of the 34th Annual Meeting of the Cognitive Science Society.

428     Austin, TX: Cognitive Science Society. 1596-1601.

429     Gordon, D., & Lakoff, G. (1971). Conversational postulates. *Papers from the Seventh*

430     *Regional Meeting of the Chicago Linguistic Society*, 63-84.

431     Grice, H. P. (1957). Meaning. *Philosophical Review*, 67: 377-388.

432     Král, P., & Cerisara, C. (2010). Dialogue act recognition approaches. *Computing and*

433     *Informatics*, 29: 227-250.

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

435     Levinson, S. C. (1995). Interactional biases in human thinking. In E. N. Goody (ed.), *Social*

436     *intelligence and interaction*. Cambridge: Cambridge University Press. 221-260.

437     Magyari, L., & De Ruiter, J. P. (2012). Prediction of turn-ends based on anticipation of

438     upcoming words. *Frontiers in Psychology*, 3: 376.

439    Paek, T., & Horvitz, E. (2000). Conversation as action under uncertainty. *Proceedings of the*

440    *Sixteenth Conference on Uncertainty in Artificial Intelligence*. San Francisco: Morgan

441    Kaufmann. 455-464.

442    Perrault, C. R., & Allen, J. F. (1980). A plan-based analysis of indirect speech acts.

443    *Computational Linguistics*, 6: 167-182.

444    Rangarajan Sridhar, V. K., Bangalore, S., & Narayanan, S. (2009). Combining lexical,

445    syntactic and prosodic cues for improved online dialog act tagging. *Computer Speech and*

446    *Language*, 23: 407-422.

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

448    organization of turn-taking for conversation. *Language*, 50: 696-735.

449    Schegloff, E. A. (1968). Sequencing in conversational openings. *American Anthropologist*,

450    70: 1075-1095.

451    Schegloff, E. A. (1982). Discourse as an interactional achievement: some uses of 'uh huh'

452    and other things that come between sentences. In D. Tannen (ed.), *Georgetown University*

453    *Roundtable on Languages and Linguistics*. Washington DC: Georgetown University Press.

454    71-92.

455    Schegloff, E. A. (1993). Reflections on quantification in the study of conversation. *Research*

456    *on Language and Social Interaction*, 26: 99-128.

457    Schegloff, E. A., & Sacks, H. (1973). Opening up closings. *Semiotica* VIII, 4: 289-327.

458    Searle, J. R. (1975). Indirect speech acts. In P. Cole & J. Morgan (eds.), *Syntax and*

459    *Semantics, Vol. 3: Speech Acts*. New York: Academic Press. 59-82.

460 Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical*

461 *Journal*, 27: 379-423.

462 Shriberg, E., Bates, R., Stolcke, A., Taylor, P., Jurafsky, D., Ries, K., Coccaro, N., Martin,

463 R., Meteer, M., & Van Ess-Dykema, C. (1998). Can prosody aid the automatic classification

464 of dialog acts in conversational speech? *Language and Speech*, 41: 439-487.

465 Stivers, T., Enfield, N. J., Brown, P., Englert, C., Hayashi, M., Heinemann, T., Hoymann, G.,

466 Rossano, F., De Ruiter, J. P., Yoon, K. E., & Levinson, S. C. (2009). Universals and cultural

467 variation in turn-taking in conversation. *Proceedings of the National Academy of Sciences of*

468 *the United States of America*, 106: 10587-10592.

469 Stolcke, A., Ries, K., Coccaro, N., Shriberg, E., Bates, R., Jurafsky, D., Taylor, P., Martin,

470 R., Van Ess-Dykema, C., & Meteer, M. (2000). Dialogue act modelling for automatic tagging

471 and recognition of conversational speech. *Computational Linguistics*, 26: 339-373.

472 Taylor, P., King, S., Isard, S., & Wright, H. (2000). Intonation and dialogue context as

473 constraints for speech recognition. *Language and Speech*, 41: 493-512.

474 Thomson, B. (2010). Statistical methods for spoken dialogue management. PhD thesis,

475 University of Cambridge.

476 Traum, D. R. (1999). Speech acts for dialogue agents. In M. Wooldridge & A. Rao (eds.),

477 *Foundations of Rational Agency*. Dordrecht: Kluwer Academic Publishers. 169-201.

478 Yngve, V. (1970). On getting a word in edgewise. In M. A. Campbell (ed.), *Papers from the*

479 *Sixth Regional Meeting, Chicago Linguistics Society*. Chicago: University of Chicago Press.

480 567-578.

481    Young, S., Gasic, M., Keizer, S., Mairesse, F., Schatzmann, J., Thomson, B., & Yu, K.

482    (2010). The Hidden Information State Model: a practical framework for POMDP-based

483    spoken dialogue management. *Computer Speech and Language*, 24: 150-174.

484    Young, S., Gasic, M., Thomson, B., & Williams, J. (2013). POMDP-based statistical spoken

485    dialogue systems: a review.  To appear in *Proceedings of the IEEE.*