

Computational Models of Miscommunication Phenomena

Matthew Purver

Cognitive Science Research Group,
School of Electronic Engineering and Computer Science,
Queen Mary University of London

Julian Hough

Dialogue Systems Group,
Faculty of Linguistics and Literature,
Bielefeld University

Christine Howes

Centre for Linguistic Theory and Studies in Probability (CLASP),
Department of Philosophy, Linguistics and Theory of Science,
University of Gothenburg

Keywords:

miscommunication, dialogue, repair, disfluency, incrementality, parallelism, sparsity

Corresponding Author: Matthew Purver, m.purver@qmul.ac.uk

School of Electronic Engineering and Computer Science, Queen Mary University of
London, Mile End Road, London E1 4NS, United Kingdom

Abstract

Miscommunication phenomena such as repair in dialogue are important indicators of the quality of communication. Automatic detection is therefore a key step towards tools which can characterize communication quality, and thus help in applications from call centre management to mental health monitoring. However, most existing computational linguistic approaches to these phenomena are unsuitable for general use in this way, and particularly for analysing human-human dialogue: although models of other-repair are common in human-computer dialogue systems, they tend to focus on specific phenomena (e.g. repair initiation by systems), missing the range of repair and repair initiation forms used by humans; and while self-repair models for speech recognition and understanding are advanced, they tend to focus on removal of “disfluent” material important for full understanding of the discourse contribution, and/or rely on domain-specific knowledge. We explain the requirements for more satisfactory models, including incrementality of processing and robustness to sparsity. We then describe models for self- and other-repair detection which meet these requirements (for the former, an adaptation of an existing repair model; for the latter, an adaptation of standard techniques), and investigate how they perform on datasets from a range of dialogue genres and domains, with promising results.

Computational Models of Miscommunication Phenomena

Repair Phenomena

One of the primary strategies by which interaction participants achieve and maintain shared understanding is *repair*: a set of strategies for highlighting and/or resolving instances of miscommunication or potential miscommunication. Not only are repair phenomena pervasive in conversation, and highly systematic, but their presence can reveal much about quality of communication, interaction and the participants themselves. A speaker can repair their own utterance, to adjust or clarify their talk (“*self-repair*”); this can be performed as the utterance is produced (example (1)), or later in a subsequent utterance (2). (In examples throughout, we show the *antecedent* (the material to be repaired) underlined, and the repair itself in **bold**.) These self-initiated examples reflect how hard speakers work on a turn-by-turn level to produce and fine-tune talk that is understandable to their specific conversational partner:

(1)¹ Deb: Kin you wait till we get home? We’ll be home in five minutes.
 Anne: Ev//en less th’n that.
 Naomi: But c’d we **c’d** I stay u:p?

(2)² L: I read a very interesting story today,
 M: uhm, what’s that.
 L: **w’ll not today, maybe yesterday, aw who knows when, huh,**
 it’s called Dragon Stew.

However, a speaker can also repair another’s utterance (3), or signal misunderstanding in order to elicit repair from the original speaker (4). These other-initiated examples (which we jointly term *other-repair* here) reflect how much effort speakers make to clarify understanding and address misunderstanding, in order to

¹(Schegloff, Jefferson, & Sacks, 1977) example (17)

²(Schegloff et al., 1977) example (22)

19 reach shared understanding:

- (3)³ Anon 3: Last year I was fifteen for the third time round.
 Grace: Yeah. <laugh> Fifteen for the first time round.
 Anon 3: **Third**.
 Grace: Third time round.
 Anon 3: Third time round.

- (4)⁴ Sarah: Leon, Leon, sorry she's taken.
 Leon: **Who?**
 Sarah: Cath Long, she's spoken for.

22 Self-repairs are conventionally regarded as symptomatic of problems with
 23 communication on the part of the speaker, caused by self-monitoring or production
 24 issues (Bard, Lickley, & Aylett, 2001; Levelt, 1983). However, they are often associated
 25 with more interactive aspects of dialogue – many occur as we tailor our talk for specific
 26 addressees, or as a direct result of feedback from our interlocutors (Goodwin, 1979).
 27 There is also evidence that they do not just indicate miscommunication, but contribute
 28 to improving the effectiveness of interaction. For example, the presence of self-repairs
 29 can aid referential success (Brennan & Schober, 2001), affect grammaticality judgements
 30 (Ferreira, Lau, & Bailey, 2004) while leaving repaired material available for processing,
 31 and increase the frequency of backchannel responses by which listeners indicate their
 32 continued attention and understanding (Healey, Lavelle, Howes, Battersby, & McCabe,
 33 2013). Other-repair too, despite the conventional view that it indicates negative aspects
 34 of miscommunication, has been shown to play a key role in semantic coordination (Mills
 35 & Healey, 2006), with evidence that increased levels of other-repair can improve task
 36 performance and speed up convergence on ways of referring (Mills, 2013).

37 Repair occurs across languages: cross-linguistic studies have shown that other
 38 initiation of repair is a standard function of questions, although the frequency of this
 39 can vary (see Stivers & Enfield, 2010, and others in that volume), and that many
 40 languages share the same repair mechanisms (Dingemanse et al., 2015) and even the
 41 surface form of the basic repair initiator “Huh?” (Dingemanse, Torreira, & Enfield,
 42 2013). Rates of repair vary with a startling variety of factors, though; for example,

³BNC file KPE, sentences 326–331

⁴BNC file KPL, sentences 347–349

43 different domains and dialogue roles (Colman & Healey, 2011), modalities (Oviatt,
44 1995), dialogue moves (Lickley, 2001), gender and age groups (Bortfeld, Leon, Bloom,
45 Schober, & Brennan, 2001). This is particularly well illustrated in the psychiatric
46 domain, where aspects of doctor-patient communication are known to be associated
47 with patient outcomes, in particular patient satisfaction, treatment adherence and
48 health status (Ong, De Haes, Hoos, & Lammes, 1995), and studies specifically
49 investigating repair show associations between repair and factors of clinical significance.
50 Lake, Humphreys, and Cardy (2011) found that participants on the autistic spectrum
51 revised their speech less often than controls, and used fewer filled pauses. For patients
52 with schizophrenia, different rates of repair have been linked to assessments of the
53 therapeutic relationship (McCabe, 2008), to specific types of symptoms such as verbal
54 hallucinations (Leudar, Thomas, & Johnston, 1992), to difficulties with turn-taking that
55 increase patient's social exclusion (Howes, Lavelle, Healey, Hough, & McCabe, 2017)
56 and to likelihood of a patient's future adherence to their treatment (McCabe et al.,
57 2013), with both self- and other-repair playing a significant role.

58 Detecting and characterising both self- and other-repair can therefore give useful
59 information about the interactive processes which drive dialogue, the participants
60 involved, and the quality of their interaction— see e.g. Healey, Colman, and Thirlwell
61 (2005), who propose a repair-based approach to indexing coordination in dialogue.
62 However, approaches of this kind have so far relied on time-consuming manual
63 annotation. The ability to *automatically* quantify repair would therefore open a new
64 opportunity for large scale (potentially real-time) dialogue analysis.

65 Existing computational models, however, have focused either on production
66 mechanisms or on the needs of human-computer dialogue systems, rather than those of
67 analysing human-human conversation. In the human-computer setting, the ability to
68 deal smoothly with human input containing self-repair has long been recognised (see
69 e.g. Johnson & Charniak, 2004), as has the need for systems to perform repair initiation
70 when confidence in automatic speech recognition (ASR) or understanding is low (e.g.
71 San-Segundo, Montero, Ferreiros, Córdoba, & Pardo, 2001); some recent systems can

72 also perform self-repair in their own output when available information changes
 73 (Skantze & Hjalmarsson, 2010). However, current techniques are limited in their
 74 abilities (see next section); and recognition of repair initiation by users has received very
 75 little attention. As examples like (5) (from a real system prototyping session) show,
 76 though, the lack of this ability leads to errors if users try to initiate repair as usual:

77 (5)⁵

Sys: Go East on Miranda Ave
Usr: what is the next turn
Sys: Make sharp right on <u>Arastradero Rd</u>
Usr: what road
Sys: I'm sorry, I do not understand what you mean by: what road. Please try rephrasing that.
Usr: what road do i turn on
Sys: I'm sorry, I do not understand what you mean by: what road do i turn on. Please try rephrasing that.
Usr: what road do i turn on and where is alma street from there
Sys: Did you want to plan a route from there?

78 The ability to automatically detect repairs is therefore important not only to
 79 provide tools for human-human interaction analysis (with potential applications
 80 including medical diagnosis and treatment monitoring), but also to improve
 81 human-computer dialogue systems as user behaviour becomes more natural. Here, we
 82 investigate models for self- and other-repair detection, and test how well they generalise
 83 between domains, with particular interest in the clinical domain.

84 Requirements and Existing Models

85 Types of Repair

86 In the conversation analysis (CA) literature (e.g. Schegloff et al., 1977), repair has
 87 long been a key subject of study, and is characterised in terms of who *initiates* the
 88 (need for) repair (oneself or another), who *completes* the repair (self or other), and in
 89 what *position* the repair is completed. Cases such as example (1) above, in which a
 90 speaker repairs their own utterance in the course of producing it, are thus termed
 91 *position one self-initiated self-repair* (P1SISR); repairing one's own antecedent
 92 utterance following an interlocutor's utterance, as in (2), a *position three self-initiated*

⁵Original data from prototype testing, CHAT project (Weng et al., 2007).

93 *self-repair* (P3SISR). An adjacent repair of another speaker's utterance, as in (3), is a
 94 *position two other-initiated other repair* (P2OIOR), and a clarification request as in (4)
 95 is a *position two next turn repair initiator* (P2NTRI). If the original speaker is then
 96 prompted to repair their problematic antecedent, as in the final utterance in each of (4),
 97 (6)–(9), this constitutes *position three other-initiated self repair* (P3OISR).

98 Colman and Healey (2011) show that by far the most common of these (more
 99 frequent than all other repair types combined), in both general conversation and
 100 task-oriented dialogue, is P1SISR self-repair (which is further subcategorised as
 101 articulation and reformulation), in line with CA's observations on the preference for
 102 self-repair in conversation (Schegloff et al., 1977). P2NTRI other-repair initiation is the
 103 next most common, and much more so than direct repair in that position (P2OIOR);
 104 responses to those in the form of P3OISR come next, with other types much less
 105 frequent. We therefore focus here on the most common forms of self- and other-repair
 106 (P1SISR, P2NTRI), noting also that McCabe et al. (2013) identify these as major
 107 informative factors in their predictive clinical model.

108 Even these categories, however, can take a variety of surface forms. P2NTRIs (or
 109 *clarification requests* (CRs), see e.g. (Ginzburg & Cooper, 2004)) can appear not only
 110 as *wh*-words as in (4), but short fragments (6), longer reprises or echoes (but not
 111 necessarily verbatim) (7), and more explicit or conventional indicators (8)–(9) (Purver,
 112 Ginzburg, & Healey, 2003):

- | | | | |
|-----|------------------|--|--|
| 113 | (6) ⁶ | | Lara: <u>There's only two people</u> in the class.
Matthew: Two people?
Unknown: For cookery, yeah. |
| 114 | (7) ⁷ | | Anon 5: <u>Oh he's started this other job</u>
Margaret: Oh he's started it?
Anon 5: Well, he he <pause> he works like the clappers he does! |

⁶BNC file KPP, sentences 352–354

⁷BNC file KST, sentences 455–457

115

Cassie: You did get off with him?

116 (8)⁸Catherine: Twice, but it was totally non-existent kissing soCassie: **What do you mean?**

Catherine: I was sort of falling asleep.

117 (9)⁹Anon 2: Gone to the cinema tonight or summat.Kitty: **Eh?**

Anon 2: Gone to the cinema

118 **Manual Analysis and Annotation**

119 [Healey et al. \(2005\)](#) present a protocol for coding repair in interaction which
 120 identifies the different CA types of repair described above. Reliability of the protocol
 121 was shown to be encouraging — in an exercise re-coding a corpus of examples from the
 122 CA literature, 75% were assigned the same category as in the original — although
 123 detection agreement rates were not reported. Many more general annotation schemes
 124 for dialogue acts or utterance functions include repair initiation as a category (e.g.
 125 [Jurafsky, Shriberg, & Biasca, 1997](#); [Stivers & Enfield, 2010](#)). Some use more
 126 fine-grained categorisations: P2NTRI repair initiators have been subcategorised
 127 according to various aspects of syntactic form, semantic structure and pragmatic level
 128 of intention (see e.g. [Purver et al., 2003](#); [Rodríguez & Schlangen, 2004](#)). All such efforts
 129 we are aware of treat complete utterances or speaker turns as the candidate units for
 130 annotation: other-repair is by its nature a between-speaker phenomenon, and therefore
 131 naturally bounded by speaker changes.

132 Self-repair, on the other hand, can begin and end within a single speaker turn, so
 133 P1SISRs are often characterised using a word-level structural schema ([Shriberg, 1994](#)):

134 (10)

$$\underbrace{\text{John and Bill}}_{\text{original utterance}} \quad \underbrace{[\text{like} + \{\text{uh}\}]}_{\text{reparandum}} \quad \underbrace{\text{love}}_{\text{interregnum}} \quad \underbrace{]}_{\text{repair}} \quad \underbrace{\text{Mary}}_{\text{continuation}}$$

135 This structure affords three principal subtypes of self-repairs: *repetitions*,
 136 *substitutions* and *deletions*. Repetitions (‘articulations’ in CA terms) have identical
 137 reparandum and repair phases; substitutions have a repair phase that differs from its

⁸BNC file KP4, sentences 521–524

⁹BNC file KPK, sentences 580–582

138 reparandum phase lexically but is clearly substitutive of it; and deletions have no
139 obvious repair phase that is substitutive of their reparandum, with utterance-initial
140 deletions often termed *restarts* (both substitutions and deletions are ‘reformulations’ in
141 CA). Despite the information such an approach provides, inter-annotator agreement is
142 often low, and the consideration of gradient boundaries between categories may be more
143 useful in some cases (Hough & Purver, 2013). Presence of a repair (or repair initiator)
144 alone is agreed upon more often than structure or specific category.

145 Formal linguistic analyses of some repair mechanisms have been given, with some
146 offering a unified treatment of self- and other-repair (e.g. Ginzburg, Fernández, &
147 Schlangen, 2007); the differences in their form have so far kept annotation and
148 computational approaches separate, though, and we maintain that distinction here.

149 **Requirements for Models of Repair**

150 These repair phenomena illustrate how dialogue participants manage and resolve
151 (potential) misunderstandings as they arise, through and within interaction. For any
152 computational model that hopes to capture them, whether in order to analyse
153 human-human conversation, or produce a human-like dialogue system, this imposes
154 several fairly challenging requirements; and few existing computational models meet
155 these requirements with any degree of generality.

156 ***Parallelism with context.*** While both self- and other-repair can take many
157 forms (1)–(9), all involve a reference to the antecedent material in context; ascribing a
158 semantic interpretation must therefore require a model of this context (see e.g. Purver
159 et al., 2003). Even if detection, rather than full interpretation, is the focus, many forms
160 (e.g. the very common reprise NTRI forms in (4), (6)) can only be interpreted by
161 detecting this reference via some form of similarity or parallelism with the antecedent;
162 while many self-repair models are based on this, most other-repair models are not. This
163 must go beyond simple lexical or syntactic repetition: some cases exploit similarities
164 which are semantic (11), phonological (12) or even orthographic (13), and might be

165 understood by one participant but not intended by the other (13):

166 (11)¹⁰ Dr: Are you suspicious are you suspicious of people
 P: Suspicious
 Dr: Paranoid
 P: **Jealous**
 Dr: Jealous yeah

167 (12)¹¹ Dr: Paroxetine
 P: **Fluoxetine**
 Dr: Ah Fluoxetine

168 (13)¹² Ustr: how long
 Wiz: dave's house is six minutes away
 Ustr: **was that one six or six zero minutes**
 Wiz: six minutes away

169 **Incrementality.** Repair phenomena are inherently incremental: both self- and
 170 other-repair often occur mid-utterance with little regard for conventional notions of
 171 grammatical constituency or completeness (Howes, Purver, Healey, Mills, &
 172 Gregoromichelaki, 2011) – see (14). Detection models must be able to operate over
 173 incomplete utterances; in the case of human-computer dialogue systems, reacting
 174 suitably as soon as is appropriate.

175 (14)¹³ A: And er they X-rayed me, and took a urine sample, took a blood
 sample. Er, the doctor
 B: **Chorlton?**
 A: Chorlton, mhm, he examined me, erm, he, he said now they were
 on about a slide <unclear> on my heart.

176 A model for other-repair detection can rely on speaker changes to indicate
 177 potential repair points, but must be able to handle incomplete context and antecedent
 178 material. A self-repair detection model, however, must operate incrementally at a
 179 finer-grained level, considering individual words and even partial words.

180 **Monotonicity.** Another key requirement that stems from the incrementality of
 181 language processing is that the reparandum must be kept available for future processing.
 182 Psycholinguistic evidence shows that people do not discard repaired material (Brennan

¹⁰Doctor-patient interaction data, (McCabe et al., 2013).

¹¹Doctor-patient interaction data, (McCabe et al., 2013).

¹²Original data from prototype Wizard-of-Oz testing, CHAT project (Weng et al., 2007).

¹³BNC file KPY, sentences 1005–1008

183 & Schober, 2001; Ferreira et al., 2004), and a model of context cannot therefore remove
184 or overwrite antecedents, which can be anaphorically referred to (15), or crucial in the
185 final interpretation of the utterance (16) (see Hough & Purver, 2012).

186 (15)¹⁴ | Nancy: Um The interview was, **it** was alright

187 (16)¹⁵ | A: Peter went swimming with Susan, or rather **surfing**, yesterday

188 **Robustness to sparsity.** Repair phenomena can be sparse. This is
189 particularly clear for other-repair: P2NTRIs typically make up only 3-6% of utterances
190 (3-4% (Purver et al., 2003), 5.8% (Rodríguez & Schlangen, 2004), 5.1% (Rieser &
191 Moore, 2005)). However, in some domains, rates can be much lower: in the clinical
192 dialogue domain of interest here, rates of P2NTRIs in patient speech can be as low as
193 0.8% (McCabe et al., 2013). Self-repair is, on the face of it, much more common, with
194 16-24% of utterances in general conversation containing a P1SISR (Hough, 2015);
195 however, the proportion of *words* which begin a P1SISR is low (3.7-5.3%, Hough, 2015;
196 Hough & Purver, 2013). As P1SISR is a within-utterance phenomenon, in which any
197 word could potentially begin a repair, the sparsity problem is therefore still very real.

198 Computational Models

199 Despite progress in psycholinguistic modelling of production problems, most
200 notably by Levelt (1983, 1989), most practical computational self-repair models have
201 been designed for use in ASR and dialogue systems; while detection accuracy can be
202 high, most take an approach of ‘cleaning’ speech of disfluent elements. This means they
203 generally remove reparanda (antecedents), operate non-incrementally, and rely on
204 relatively domain-specific dependency parsing rather than more general parallelism (e.g.
205 Honnibal & Johnson, 2014; Rasooli & Tetreault, 2014) – thus failing to meet our
206 requirements above. Some recent systems are incremental, and use more general
207 statistical language model information (Zwarts, Johnson, & Dale, 2010), but still focus
208 on removing antecedent material, not meeting our monotonicity requirement. They also

¹⁴From H. H. Clark (1996, p266)

¹⁵From Hough and Purver (2012)

209 generally use cleaned-up data with cut-off words removed. In contrast, the model we
 210 use below (STIR, [Hough & Purver, 2014](#)) meets all our incremental, domain-general,
 211 context-maintaining requirements and here we adapt it to handle cut-off words.

212 Computational models of other-repair initiation have generally focused on
 213 production, allowing systems to clarify errorful ASR input. Naturalness is typically
 214 limited (see (17), from the Let’s Go! system, [Raux, Langner, Black, & Eskenazi, 2005](#)),
 215 although recent developments permit more natural, targeted NTRIs where uncertainty
 216 can be localised (18) ([Stoyanchev, Liu, & Hirschberg, 2014](#)):

(17)¹⁶ U: When’s the next bus to Wood Street?
 S: **Sorry, I didn’t understand that. Please repeat.**
 ...
 U: When’s the next bus to Wood Street?
 S: **Going to WOOD STREET. Did I get that right?**
 U: Yes.

(18)¹⁷ U: Do you have anything other than these [XXX] plans
 S: **Which plans? / Anything other than what?**

219 On the interpretation side, attention has been given to user correction (see e.g.
 220 [Kitaoka, Kakutani, & Nakagawa, 2005](#); [Lemon & Gruenstein, 2004](#); [Litman, Hirschberg,](#)
 221 [& Swerts, 2006](#)). When users notice system errors, they produce P2OIOR repairs, often
 222 using characteristic syntactic and prosodic forms (e.g. repetition with hyperarticulation)
 223 which then cause further misrecognition problems. Detection of corrections can
 224 therefore aid error recovery, and accuracies can be good ([Kitaoka et al. \(2005\)](#) report
 225 c.90% F-scores, although [Litman et al. \(2006\)](#) only 72% on different data, and [Lopes et](#)
 226 [al. \(2015\)](#) similar levels on a specific sub-task, repetition detection).

227 Recent approaches use general learning frameworks to induce these functionalities
 228 from data (see e.g. [Young, Gašić, Thomson, & Williams, 2013](#)), but do this by learning
 229 the optimal action for systems to take in a given context; this does not therefore
 230 directly generalise to detecting clarification by human users. While strategies for
 231 responding to user NTRIs could certainly be learned in principle, we are not aware of

¹⁶Let’s Go! system examples ([Stoyanchev & Stent, 2012](#)).

¹⁷From ([Stoyanchev et al., 2014](#)); [XXX] represents a missing or unrecognised word.

232 current implementations; and these would not be suited to third-party analysis, being
233 dependent on system interaction in the dialogue.

234 Dialogue act tagging tools, on the other hand, are designed for third-party
235 analysis; however, they tend to be optimised for general overall accuracy, leading to
236 relatively poor results for sparser classes, including repair and repair initiation. Much
237 work does not attempt to classify these sparse classes (e.g. [Stolcke et al., 2000](#)); where
238 results are given, accuracies are poor. [Surendran and Levow \(2006\)](#) report 43% F-scores
239 on their P2NTRI category (`check`, 8% of turns in their dataset) and only 19% for
240 P3OISRs (`clarify`, 4% of turns); [Schlangen \(2005\)](#) reports 30-40% F-scores on similar
241 classes. [Fernández, Ginzburg, and Lappin \(2007\)](#) report good accuracies but only for a
242 restricted sub-type of P2NTRIs (elliptical noun phrase fragments).

243 Below, we outline and test our own approach to general detection of repair and
244 repair initiation, suitable for human-human as well as human-computer data and
245 compatible with the requirements outlined above.

246 Materials

247 Corpora

248 **Switchboard (SWBD).** Our first corpus is one commonly used for testing
249 computational self-repair models and dialogue act taggers. The Switchboard corpus
250 ([Godfrey, Holliman, & McDaniel, 1992](#)) consists of approximately 2400 dyadic telephone
251 conversations between American participants unfamiliar with one another, on topics
252 assigned from a pre-determined list. For other-repair, we use the Dialogue Act version
253 of Switchboard ([Jurafsky et al., 1997](#)), with 1155 dialogues totalling over 120,000
254 utterances and nearly 1.5m words. For self-repair, we use the disfluency-tagged portion
255 of Switchboard ([Meteer, Taylor, MacIntyre, & Iyer, 1995](#)), with 650 conversations of
256 duration 1.5-10 minutes (average around 6.5 minutes), with a standard division into
257 train, heldout and test sections (see [Hough & Purver, 2014](#); [Johnson & Charniak, 2004](#)).

258 **British National Corpus (BNC-CH, BNC-PGH).** We also investigate how
259 well our methods generalise to more open-domain and multi-party conversation. The

260 BNC-CH corpus (Colman & Healey, 2011) is a subset of the demographic portion
261 (transcribed spontaneous natural conversations made by members of the public) of the
262 British National Corpus (BNC, Burnard, 2000). It contains 31 dialogues annotated for
263 self- and other-repair, with 1933 utterances, 14,034 words produced by 41 people. The
264 BNC-PGH corpus (Purver et al., 2003) is a different, larger subset (c.150,000 words)
265 containing sections from 56 dialogues including specific contexts (e.g. doctor-patient
266 conversations) as well as demographic data, and annotated only for other-repair
267 initiation (in their terminology, clarification requests).

268 **Psychiatric Consultation Corpus (PCC).** To test applicability to a clinical
269 domain, we use a corpus from a study investigating clinical encounters in psychosis
270 (McCabe et al., 2013): transcripts from 51 outpatient consultations of patients with
271 schizophrenia and their psychiatrist, including 51 different patients, and 17
272 psychiatrists. Consultation length varies from only 709 words (c.5 minutes) to 8526
273 (nearly an hour), with mean length 3500 words.

274 **Map Task Corpus (MAPTASK).** To further investigate robustness to
275 change in dialogue style, genre and domain, we also use the HCRC Map Task Corpus
276 (Anderson et al., 1991), with 128 two-person dialogues containing 18,964 turns with
277 c.150,000 words. These conversations concern a very specific task: guiding an
278 interlocutor around a map whose features may not appear identical to the two parties.

279 Annotation

280 SWBD's disfluency annotations include filled pauses, discourse markers, and edit
281 terms, all with standardised spelling (e.g. consistent 'uh' and 'uh-huh' orthography).
282 P1SISRs are bracketed with the structure in (10), with reparandum, interregnum and
283 repair phases marked. Restart repairs (utterance-initial deletions) are coded as two
284 separate units and not in fact annotated as repairs. In the dialogue act corpus,
285 P2NTRIs are tagged as *signal-non-understanding* (br); Jurafsky et al. (1997) report
286 overall inter-annotator agreement of 80% kappa, although figures specifically for this
287 tag are not given.

288 For the BNC-CH and PCC, each transcript is hand-annotated for both self- and
289 other-repair using Healey et al. (2005)’s protocol discussed above. Colman and Healey
290 (2011) and McCabe et al. (2013) report inter-annotator agreement of c.75% kappa.
291 BNC-PGH is annotated only for other-repair initiation P2NTRIs (Purver et al. (2003)
292 report 75%-95% kappa); MAPTASK similarly provides information on P2NTRIs (via
293 check tags) but not self-repair.

294 SWBD, BNC and MAPTASK provide gold-standard part-of-speech (POS) tags;
295 we tagged the PCC using the Stanford POS tagger (Toutanova, Klein, Manning, &
296 Singer, 2003). This is trained on written text; application to spoken dialogue has shown
297 c.10% error rates (Mieskes & Strube, 2006). Here, however, we are not concerned with
298 POS labels *per se*, but in the parallelism between POS sequences - as errors are likely to
299 be fairly consistent (dependent on transcription spelling or spoken dialogue
300 idiosyncracies) we take this as sufficient for our purposes.

301 Detecting Other-Repair

302 In order to detect NTRIs, we define a set of turn-level features that could be
303 extracted from transcripts automatically, and that encode either specific NTRI
304 characteristics (e.g. presence of clarificational words like “pardon”) or more general
305 parallelism features between the turn to be classified and the previous turn by other
306 and same speaker. (This assumes antecedents of clarification are in the immediately
307 preceding turn; Purver et al. (2003) found this to cover 85% of cases). We then train a
308 standard supervised discriminative classifier using these features to detect NTRI turns.

309 This approach meets all our requirements. The notion of incrementality here is at
310 the level of speaker turns: we therefore use only information from current and previous
311 turns so that a classification decision can be made immediately (although subsequent
312 turns can certainly contain useful information). Parallelism with context was captured
313 by designing suitable features: lexical parallelism via simple word string matching;
314 syntactic parallelism by matching part-of-speech tags; and semantic parallelism via
315 neural network models of word similarity. Sparsity varies considerably between datasets:

316 while 11% of MAPTASK utterances are NTRIs, this drops to 4% in the two BNC
317 datasets, 1% in PCC, and only 0.2% in SWBD. To deal with it, we trained the classifier
318 with a weighted cost function, weighting errors on true positive examples more than
319 those on negative ones. Full details of feature calculation and classifier implementation
320 are in the Supplementary Material; the full set of features is shown here in Table 1.

321 [Table 1 about here]

322 Results

323 We test this approach on each of our datasets using 10-fold cross-validation (see
324 Supplementary Material for full details); results are shown in Table 2. Performance is
325 shown against two baselines: always predicting the NTRI class, and using a one-rule
326 classifier with the most helpful single feature. We show performance using our general
327 NTRI and parallelism features (“high-level” features in Table 2), and using all observed
328 unigrams (unique single words, “all” in Table 2). This latter approach illustrates the
329 performance achievable with specific lexical information, but is likely to be highly
330 dataset- and domain-dependent and susceptible to over-fitting, so we treat it as an
331 indicative “ceiling” rather than a suggested robust approach. We also show the
332 performance achieved by [Howes, Purver, McCabe, Healey, and Lavelle \(2012\)](#) on
333 patient-only NTRIs within the PCC dataset, for comparison.

334 Our primary evaluation metrics are F-score (the harmonic mean of precision and
335 recall) for the class of interest (NTRIs), and the area under the precision-recall curve
336 (AUC-PRC): as our weighted classifiers can be adjusted to trade precision against
337 recall, this AUC metric is more informative than F-score alone; and the F-score we
338 show is for the point where precision and recall are balanced. We also show the more
339 familiar receiver-operator curve area (AUC-ROC), although it is less suitable for
340 unbalanced data, as it underestimates the effect of poor performance on the sparser
341 (and here, more interesting) class (see [Saito & Rehmsmeier, 2015](#)).

342 [Table 2 about here]

343 Performance varies with the nature of the dataset: with the open-domain BNC,

344 performances are fairly good with F-scores of 52-55% (AUC-PRC 0.51-0.52); in the
345 more domain-specific clinical PCC, F-scores drop below 50%; and in MAPTASK even
346 further to 38%. (Note that baseline F-scores with such unbalanced data are low, with
347 AUC-PRC scores all below 0.22). Encouragingly, the approach seems fairly robust to
348 sparsity itself, with reasonable performance in both the PCC and more open-domain
349 (but telephone-based) SWBD, where NTRIs make up only 1.3% and 0.2% of utterances
350 respectively. (The lowest performance is in the least sparse data (MAPTASK), in fact).

351 In most datasets, the general high-level features transfer well across domains, with
352 performance similar to the specific unigram features; the exception is MAPTASK and
353 (to a lesser degree) SWBD, suggesting the presence of more domain-specific and/or
354 variable repair mechanisms in those settings. We investigate the most predictive
355 features (by selecting based on information gain); details and feature lists are in the
356 Supplementary Material (Table 6). The most informative are usually the simpler
357 features (interrogative features such as wh-words and question marks; repair keywords;
358 utterance length). Semantic parallelism features (word vector-based similarities) then
359 feature strongly, mixed with the lexical and POS repetition features. However,
360 removing these semantic parallelism features makes little difference to performance:
361 while AUC-PRC tends to drop, indicating less robust performance, the drop is small
362 (1-2%), and the point F-scores do not change; this suggests that the vector-based
363 features capture little information beyond the simpler symbolic ones. Best features for
364 the worst-performing dataset (MAPTASK) are noticeably different, again suggesting
365 different repair mechanisms, with backchannel keywords and repetition seeming to play
366 a stronger role, and wh-words not being useful.

367 **Error analysis.** To investigate the limitations and common sources of error, we
368 trained and tested a version on the same full dataset (BNC-PGH), thus giving an upper
369 bound to performance using this feature set. Performance improved only slightly
370 ($F=0.54$, vs 0.52 using cross-validation), showing that significant limitations exist, and
371 qualitative manual inspection of the errors revealed some common sources of these.
372 NTRI cue words, wh-words, short questions (cued by transcribed question marks) and

373 repetition are all strong indicators, leading to many true positives (19), but are the
 374 main cause of false positives (20), (21) (shown *bold italic*):

375 (19)¹⁸ | Unknown: As most of the main towns in Suffolk have reviews every two years
 are you contemplating having er those, that sort of interview of erm
 public hearing.
 376 | Guy: **Er what every two years sorry?**
 Unknown: They have traffic management erm reviews every two years.

377 (20)¹⁹ | e bust: If it's no I think what we agreed Glynis if it was going to be a stone
 it could go in the wall where it could be seen from outside.
 g herbert: **Oh right yes sorry I beg your pardon .**
 378 | e bust: But if we were deciding on a brass plaque or something

379 (21)²⁰ | Neal: <pause> Same thing as as I mentioned before. It all fell out. Bags.
 Unknown: <laugh>.
 Unknown: <unclear>?
 380 | Neal: Yes, certainly. She'd got all the clothes she'd ever had.

381 Omission of question marks in transcription can therefore also cause false negatives.
 382 Other false negatives give more interesting insight about what our features fail to
 383 capture. In some cases, the key parallelism is not captured by simple sequence and
 384 vector-similarity approaches (22); even more challenging are examples with no explicit
 385 parallel elements, e.g. P2NTRIs which offer elaborating material (23) or possible
 386 continuations (24) (in what Purver et al. (2003) call *gap filler CRs*).

387 (22)²¹ | Anon 1: Four.
 Malcolm: Yep.
 Anon 1: Six. Nine.
 Malcolm: <tut> **How many ?**
 Anon 1: <unclear> <pause> Nine.
 Malcolm: Nine.

¹⁸BNC file KN3, sentences 299–301.

¹⁹BNC file KM8, sentences 599–601.

²⁰BNC file KNC, sentences 1075–1080.

²¹BNC file KND, sentences 567–573.

388

(23)²² e bust: Have have you found out any more the cost Harry of this?
 h rickett: **Yeah for a stone that is ?**
 e bust: Yes.

(24)²³ e bust: Ruby <unclear> she'll have she'll have some children though because I mean they're somewhere down in ...
 d kemp: <unclear>
 e bust: they're somewhere down in Gillingham down in ...
 d kemp: **Kent**
 e bust: Yeah they're down in Kent.

391

Detecting Self-Repair

392 For self-repair detection we use STIR ('STrongly Incremental Repair detection')
 393 (Hough & Purver, 2014).²⁴ STIR takes a local, incremental approach, detecting the
 394 structure in (10) and isolated edit terms (such as 'uh', 'um' and 'you know'), assigning
 395 appropriate structural labels – see Figure 1. While sparsity is handled similarly to our
 396 other-repair experiments, we now generalise the approach to parallelism: instead of
 397 using specific syntactic or semantic knowledge from POS taggers or word vectors, STIR
 398 uses a range of information-theoretic measures to capture parallelism in a more general
 399 fashion. The notion of incrementality is also different, as a fully word-by-word approach
 400 is required (as discussed above).

401 [Figure 1 here]

402 Rather than detecting the repair structure in its left-to-right string order,
 403 detection consists of 4 time-steps as words are encountered: STIR first detects edit
 404 terms (possibly interregna) at step T1; then repair onsets rp_{start} at T2; if one is found,
 405 it searches backwards to find the reparandum start rm_{start} at T3; then finally finds the
 406 repair end rp_{end} at T4. Step T1 relies mainly on lexical probabilities; T2 exploits
 407 features of divergence from 'fluent' language; T3 uses fluency of utterances without the
 408 hypothesised reparanda, and parallelism between repair and reparandum; and T4 the
 409 similarity between distributions after reparandum and repair end points (indicated by

²²BNC file KM8, sentences 534–536.

²³BNC file KM8, sentences 741–744; ellipsis ... added to show putative 'antecedent'.

²⁴Available from <http://bitbucket.org/julianhough/stir>.

410 the dotted edge between S3 and S4 in Figure 1). Each stage implements these insights
411 via multiple related features in a statistical classifier, and the four stages are connected
412 together in a pipeline (Figure 2). The output is a graph-like structure (Figure 1). STIR
413 has previously been applied to SWBD; here, we investigate its transfer to our other
414 datasets, and the nature of its errors, while updating it to handle cut-off words.

415 [Figure 2 about here]

416 Classifiers and features

417 Each individual classifier has its own error function, allowing trade-off of
418 immediate accuracy, run-time and stability, and balance in the face of sparsity. Each
419 classifier also uses its own specific combination of features, but all derived from basic
420 information-theoretic measures from n-gram language models (LMs). N-gram LMs are
421 easily applied incrementally, require no commitment to any particular grammar
422 formalism, and can be extended to model levels other than the purely lexical, e.g.
423 grammaticality judgements (A. Clark, Giorgolo, & Lappin, 2013). We train our LMs on
424 the standard Switchboard training data, following Johnson and Charniak (2004) by
425 cleaning the data of all edit terms and reparanda, to approximate a ‘fluent’ LM. We
426 train two such models, one for words and one for POS tags;²⁵ this allows us to derive
427 features giving syntactic as well as lexical information, both by using POS tags directly
428 and via A. Clark et al. (2013)’s Weighted Mean Log (*WML*) measures which factor out
429 lexical probability to approximate syntactic plausibility. From these basic LMs we then
430 derive features that characterise (dis)fluency, via probability and *surprisal* for observed
431 words; uncertainty in a context, via the *entropy* of possible continuations, and increases
432 and reductions therein; and similarity or parallelism between contexts, via the
433 Kullback-Leibler (KL) divergence between distributions. We handle partial words
434 within the LM scoring itself, assigning penalties when partial words are encountered.
435 Full details of feature calculation and classifier implementation are given in the
436 Supplementary Material; we give a brief overview here.

²⁵Below, measures from the word LM are indicated by the superscript *lex* and the POS LM by *pos*.
When referring to the same measure from both LMs, these are suppressed.

437 **Edit term detection.** The first classifier uses the word surprisal s^{lex} from a
 438 specific edit word bigram LM (edit words will have high probability and therefore lower
 439 s^{lex}), and the trigram surprisal s and syntactic fluency WML from the standard fluent
 440 LMs described above (the intuition here being that general fluency will seem lower for
 441 trigrams containing an edit term). This also helps interregnum recognition, due to the
 442 inclusion of interregnum vocabulary within edit term vocabulary (Hough & Purver,
 443 2013), and provides a useful feature for repair detection in subsequent steps (Hough &
 444 Purver, 2014; Lease, Johnson, & Charniak, 2006).

445 **Repair start detection.** The second step to detect rp_{start} is arguably the most
 446 crucial component: the greater its accuracy, the better the input for downstream
 447 components and the lesser the overhead of filtering false positives required. This
 448 classifier uses a combination of simple alignment features (e.g. whether a word is
 449 identical to a predecessor), and a series of features describing local changes in LM
 450 fluency. Figure 3 shows the main intuition: that repair onsets correspond to troughs in
 451 lexical and syntactic probability measures (in Figure 3, WML^{lex}).

452 [Figure 3 about here]

453 **Reparandum start detection.** We now detect rm_{start} positions given a
 454 hypothesised rp_{start} , using two main intuitions. First, we use the noisy channel intuition
 455 of Johnson and Charniak (2004) that removing the reparandum (from rm_{start} to rp_{start})
 456 increases fluency of the utterance (captured via WML features), while removing
 457 non-reparandum words decreases it. Second, we can measure parallelism between rp_{start}
 458 and rm_{start} , via the KL divergence between their LM distributions.

459 **Repair end detection and structure classification.** Finally, detection of
 460 rp_{end} and the final structure of the repair exploits the notion of parallelism. This can be
 461 measured as divergence between the conditional probability distributions θ^{lex} at the
 462 reparandum-final word rm_{end} and the repair-final word rp_{end} : for repetition repairs, KL
 463 divergence will trivially be 0; for substitutions, it will be higher; for deletes, even higher.
 464 It can also be captured via *ReparandumRepairDifference*, the difference in probability
 465 between an utterance cleaned of the reparandum and the utterance with its repair

466 phase substituting its reparandum. In the running example from Figure 1, this would
 467 be as in equation (1).

$$\text{ReparandumRepairDifference}(\text{"John [likes + loves]"}) = \\
 WML^{lex}(\text{"John loves"}) - WML^{lex}(\text{"John likes"}) \quad (1)$$

468 Results

469 Hough and Purver (2014) show state-of-the-art performance for incremental
 470 self-repair detection (77.9% accuracy at detecting reparandum words in Switchboard
 471 test data); they removed cut-off words which on average occur every 118 words (0.84%
 472 of all words) in the Switchboard heldout data. Here we test with cut-off words included,
 473 a realistic approach for transcripts and incremental ASR output, and potentially
 474 providing further cues about repair onset. By way of comparison, we also test the
 475 performance of Hough and Schlangen (2015)’s Recurrent Neural Network (RNN)-based
 476 disfluency detector.²⁶ In all cases, we derive LM features from the SWBD training set
 477 using 10-fold cross-validation (full details in the Supplementary Material); we then train
 478 and test classifiers using a standard training/test split for each corpus.

479 We report accuracy of repair onset detection on a per-utterance level, as that is
 480 the most relevant measure for dialogue-level analysis; we also report the overall
 481 Spearman’s rank correlation of the repair rate (per utterance) between the gold
 482 standard transcripts and STIR’s output. These allow comparison with the PCC and
 483 BNC-CH annotations, which use a different annotation schema from Switchboard (see
 484 above), and (for BNC-CH) do not mark repair onset point. For Switchboard, we also
 485 report the standard per-word reparandum detection result (*F_{rm}*), in line with previous
 486 work— see Table 3. This per-word evaluation tells us about ability to identify the precise
 487 location of repairs, important for dialogue system development; but the per-utterance
 488 figures also give us a useful, if less precise, metric for practical applications such as the
 489 analysis of patient-doctor dialogues.

²⁶Code available from https://github.com/dsg-bielefeld/deep_disfluency

490 On Switchboard, accuracy of reparandum word detection reaches 78.1% on the
491 test set, and per-utterance detection accuracy is 85.0%. The correlation for repair rates
492 is very high and significant (*Spearman's rank*=0.956). This marginally improves over
493 [Hough and Purver \(2014\)](#)'s results with partial words removed; and training and testing
494 on the SWBD data with partial words removed in our experimental setup reduces
495 accuracy even more, to 76.8%. This shows the potential utility (rather than hindrance)
496 of using partial words for disfluency detection if adapted appropriately. The RNN
497 model, which is not adapted for partial words, shows the opposite pattern, dropping
498 from 66.8% to 63.8% when *introducing* partial words – see Table 4.

499 We also test on the out-of-domain PCC and BNC-CH datasets. With PCC,
500 per-utterance detection performance is very encouraging even with no optimization
501 (62.0%), and correlation of repair rates to the gold standard is also high (*Sp. R*=0.805).
502 For BNC-CH, per-utterance results are far worse (41.7%) — we attribute this to the
503 annotation protocol, which lacked the exact identification of reparandum and repair
504 phases used in the other two corpora — however, the correlation of repair rates is still
505 moderately strong (*Sp. R*=0.583, $p<0.001$). Table 3 shows that using POS LM features
506 helps detection performance in each corpus, particularly boosting correlation score for
507 our most challenging dataset, PCC (0.583 vs. 0.530); this suggests that syntactic-level
508 information can help detect repair structures.

509 **Error analysis.** The detailed Switchboard annotation format permits a
510 quantitative analysis of the error distribution, and comparison between STIR and the
511 comparable RNN model. Table 5 (a) shows the F-score with different combined
512 reparandum and interregnum lengths, where correct detection is counted if both repair
513 onset and reparandum onset are predicted correctly. All three systems show reduced
514 performance as length increases. However, reduction is less for STIR; its explicit
515 backwards search mechanism alleviates the problem of long-distance dependency, while
516 the RNN relies on internally learned memory structure and struggles further than 5
517 words back from the repair onset. Table 5 (b) shows performance for different repair
518 types. Repetitions are the easiest, followed by substitutions, then deletes; but STIR

519 performs far better on substitutions and deletions than the RNN. Both of these rarer
 520 types rely on more complex notions of parallelism and fluency, rather than the presence
 521 of verbatim repeats.

522 [Tables in Table 5 about here]

523 A qualitative survey of the errors when changing domain shows that many are due
 524 to the transcription and annotation protocols (as discussed by Howes, Hough, Purver, &
 525 McCabe, 2014), not merely poor system performance. As shown in examples (25)-(27)²⁷
 526 from the PCC, false positives occur when STIR tags embedded repairs as multiple
 527 instances, but the annotator views this as part of one longer repair (25). False negatives
 528 include confusion between editing phrases and repairs (26), a distinction in SWBD but
 529 not in Healey et al. (2005)’s annotation protocol; and missing repairs entirely (27), as
 530 utterance-initial deletions are not marked in SWBD but treated as separate utterances.

531 (25)²⁸ (a) D: ... and if you tell me that **that**_[RP_{START}] that the depressions kicks
 in ...
 (b) D: ... and if you tell me that **that**_[rp_{start}] **that**_[rp_{start}] the depressions
 kicks in ...

532 (26)²⁹ (a) D: and so **I**_[RP_{START}] mean otherwise I’m not too concerned about
 your mental health...
 (b) D: and so **I**_[ed] **mean**_[ed] otherwise I’m not too concerned about your
 mental health...

533 (27)³⁰ (a) P: I don’t **I’m**_[RP_{START}] not like hearing voices...
 (b) P: I don’t I’m not like hearing voices...

534 Discussion and Conclusions

535 Our experiments show that detection of both self-repair and other-repair initiation
 536 is possible with reasonable accuracy. For the self-repair case, by-utterance F-scores can
 537 reach 85% when trained on in-domain data, and up to 62% even when transferring a
 538 model to other (here, face-to-face clinical) data. For the much sparser other-repair case,
 539 F-scores can reach 60%, but depend on the nature of the data; while robust to sparsity

²⁷Hand annotation tags are shown in (a) in each case with STIR’s annotations shown in (b).

²⁸(Howes et al., 2014) example (10)

²⁹(Howes et al., 2014) example (11)

³⁰(Howes et al., 2014) example (12)

540 itself in Switchboard where NTRIs are particularly sparse (0.2% of turns), some
 541 domains cause bigger drops, although in the sparse clinical data F-scores still reach
 542 46%. These results are encouraging as they use general models which exploit features of
 543 repair-indicating vocabulary and parallelism, hence giving robustness across datasets
 544 and being applicable to the general case of third-party dialogue analysis.

545 Examination of the effect of features suggests that the key to good performance is
 546 capturing parallelism, reflecting the nature of repair as a resource for querying and
 547 reformulating material. However, this seems hard to achieve using general models of
 548 word meaning (as in our other-repair classifier): using general lexical matching and
 549 suitably trained information-theoretic models of word distributions, as STIR does for
 550 self-repair, seems more successful, and more robust across domains and phenomena
 551 than more directly lexically driven approaches (here, the comparison RNN). A possible
 552 direction for future research would be to investigate whether similar methods could help
 553 with the challenging cases of implicit parallelism seen with other-repair.

554 The effect of changing domains and genres suggests that some domains show
 555 different repair phenomena and mechanisms. Inspection of the task-driven Map Task
 556 data shows that the challenging other-repair types are more common (e.g. offering
 557 elaboration and reformulation), as is long-range clarification, where participants check
 558 their understanding of whole sequences of instructions (rare in the other datasets).
 559 Many of the domain-related effects, though, are associated with differences in
 560 transcription and annotation standards, as discussed above for self-repair. This is also a
 561 factor with other-repair data; for example, the Map Task annotations tag some forms of
 562 NTRI question as belonging instead to an ‘other question’ category (28).

563 (28)³¹ G: until you you get over the top of the slate mountain
 F: **over the top of the**
 G: slate mountain
 F: don't have a slate mountain

564 However, in many cases these differences in annotation approach stem from
 565 genuine ambiguity or multifunctionality. We have seen cases of self-repair where

³¹Map Task corpus, dialogue q1ec2, utterances 59-62.

566 alternate analyses are possible (25)-(27), cases of other-repair which perform repair
567 initiation simultaneously with offering possible repair (23)-(24), and many forms (e.g.
568 repeated fragments) can also perform acknowledgement or answer questions.
569 Recognising and handling this ambiguity is of course crucial for dialogue systems,
570 although resolving it is not always possible or desirable — hence the success of
571 probabilistic models which maintain uncertainty (Young et al., 2013) — and this
572 suggests that repair identification should be approached and evaluated in a probabilistic
573 fashion, not a categorical one. This also points to the limitations of using transcripts as
574 our source material. For human annotators, one of the signals of an NTRI is whether
575 the *following* turn contains a *position 3 other initiated self repair* – i.e. whether the
576 other dialogue participant has interpreted the preceding turn as requesting repair; our
577 incremental approach means we cannot benefit from this information. Of course,
578 participants in dialogue must decide whether to treat a turn as initiating repair as and
579 when they encounter it – so this cannot be how humans identify these whilst engaged in
580 dialogue. Evidence suggests that in real dialogue, feedback (positive or negative) is cued
581 by or accompanied by gaze (Hjalmarsson & Oertel, 2012), intonation (Gravano &
582 Hirschberg, 2009) or gesture (Healey et al., 2013; Healey, Plant, Howes, & Lavelle,
583 2015), suggesting that we may improve our performance if we include these features.

584 Despite these limitations, these models go a long way toward fulfilling our
585 desiderata: they operate *incrementally* (utterance-by-utterance for P2NTRIs,
586 word-by-word for P1SISRs) and *monotonically* (STIR leaves reparandum material
587 available for later processing); they use general measures of *parallelism with context*;
588 and they are relatively robust to the *sparsity* of NTRIs and rarer and longer self-repairs.
589 Such models therefore have potential not only to help make human-computer dialogue
590 systems more human-like, via more robust, incremental self-repair and other-repair
591 detection; but also to improve our ability to analyse and evaluate the quality of
592 communication in settings like clinical psychiatry.

593

Acknowledgements

594 Purver was partially supported by EPSRC (grant EP/J010383/1) and by the
595 ConCreTe project, which acknowledges the financial support of the Future and
596 Emerging Technologies (FET) programme within the Seventh Framework Programme
597 for Research of the European Commission, under FET grant number 611733. Hough
598 was supported by the Cluster of Excellence Cognitive Interaction Technology ‘CITEC’
599 (EXC 277) at Bielefeld University, funded by the German Research Foundation (DFG),
600 and the DFG-funded DUEL project (grant SCHL 845/5-1). Howes was supported by
601 two grants from the Swedish Research Council (VR); 2016-0116 – Incremental
602 Reasoning in Dialogue (IncReD) and 2014-39 for the establishment of the Centre for
603 Linguistic Theory and Studies in Probability (CLASP) at the University of
604 Gothenburg. We thank David Schlangen for extensive discussions on the topic.

References

605

- 606 Anderson, A., Bader, M., Bard, E., Boyle, E., Doherty, G., Garrod, S., . . . Weinert, R.
607 (1991). The HCRC map task data. *Language and Speech*, 34(4), 351–366.
- 608 Bard, E. G., Lickley, R. J., & Aylett, M. P. (2001). Is disfluency just difficulty? In *Isca*
609 *tutorial and research workshop (itrw) on disfluency in spontaneous speech*.
- 610 Bortfeld, H., Leon, S. D., Bloom, J. E., Schober, M. F., & Brennan, S. E. (2001).
611 Disfluency rates in conversation: Effects of age, relationship, topic, role, and
612 gender. *Language and speech*, 44(2), 123–147.
- 613 Brennan, S., & Schober, M. (2001). How listeners compensate for disfluencies in
614 spontaneous speech. *Journal of Memory and Language*, 44(2), 274–296.
- 615 Burnard, L. (2000). *Reference guide for the british national corpus (world edition)*.
616 Oxford University Computing Services. Retrieved from
617 <http://www.natcorp.ox.ac.uk/docs/userManual/>
- 618 Clark, A., Giorgolo, G., & Lappin, S. (2013, August). Statistical representation of
619 grammaticality judgements: the limits of n-gram models. In *Proceedings of the*
620 *fourth annual workshop on cognitive modeling and computational linguistics*
621 *(cmcl)* (pp. 28–36). Sofia, Bulgaria: Association for Computational Linguistics.
622 Retrieved from <http://www.aclweb.org/anthology/W13-2604>
- 623 Clark, H. H. (1996). *Using language*. Cambridge University Press.
- 624 Colman, M., & Healey, P. G. T. (2011). The distribution of repair in dialogue. In
625 L. Carlson, C. Hoelscher, & T. F. Shipley (Eds.), *Proceedings of the 33rd annual*
626 *meeting of the cognitive science society* (pp. 1563–1568). Boston, MA: Cognitive
627 Science Society.
- 628 Dingemanse, M., Roberts, S. G., Baranova, J., Blythe, J., Drew, P., Floyd, S., . . .
629 others (2015). Universal principles in the repair of communication problems. *PloS*
630 *one*, 10(9), e0136100. Retrieved from [http://journals.plos.org/plosone/](http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0136100)
631 [article?id=10.1371/journal.pone.0136100](http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0136100)
- 632 Dingemanse, M., Torreira, F., & Enfield, N. J. (2013). Is “Huh?” a universal word?
633 Conversational infrastructure and the convergent evolution of linguistic items.

- 634 *PloS one*, 8(11), e78273. Retrieved from [http://journals.plos.org/plosone/](http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0078273)
635 [article?id=10.1371/journal.pone.0078273](http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0078273)
- 636 Fernández, R., Ginzburg, J., & Lappin, S. (2007). Classifying ellipsis in dialogue: A
637 machine learning approach. *Computational Linguistics*, 33(3), 397–427.
- 638 Ferreira, F., Lau, E. F., & Bailey, K. G. D. (2004). Disfluencies, language
639 comprehension, and tree adjoining grammars. *Cognitive Science*, 28(5), 721–749.
- 640 Ginzburg, J., & Cooper, R. (2004). Clarification, ellipsis, and the nature of contextual
641 updates in dialogue. *Linguistics and Philosophy*, 27(3), 297–365.
- 642 Ginzburg, J., Fernández, R., & Schlangen, D. (2007). Unifying self- and other-repair. In
643 R. Artstein & L. Vieu (Eds.), *Proceedings of the 11th workshop on the semantics*
644 *and pragmatics of dialogue (DECALOG)*. Rovereto, Italy: SemDial.
- 645 Godfrey, J. J., Holliman, E., & McDaniel, J. (1992). SWITCHBOARD: Telephone
646 speech corpus for research and development. In *Proceedings of IEEE ICASSP-92*
647 (pp. 517–520). San Francisco, CA: IEEE.
- 648 Goodwin, C. (1979). The interactive construction of a sentence in natural conversation.
649 In G. Psathas (Ed.), *Everyday language: Studies in ethnomethodology* (pp.
650 97–121). New York: Irvington Publishers.
- 651 Gravano, A., & Hirschberg, J. (2009). Backchannel-inviting cues in task-oriented
652 dialogue. In M. Uther, R. Moore, & S. Cox (Eds.), *Interspeech* (pp. 1019–1022).
653 Brighton, UK: ISCA.
- 654 Healey, P. G. T., Colman, M., & Thirlwell, M. (2005). Analysing multi-modal
655 communication: Repair-based measures of human communicative co-ordination.
656 In J. van Kuppevelt, L. Dybkjaer, & N. Bernsen (Eds.), *Natural, intelligent and*
657 *effective interaction in multimodal dialogue systems* (Vol. 30, pp. 113–129).
658 Dordrecht: Kluwer.
- 659 Healey, P. G. T., Lavelle, M., Howes, C., Battersby, S., & McCabe, R. (2013, July).
660 How listeners respond to speaker's troubles. In M. Knauff, M. Pauen, N. Sebanz,
661 & I. Wachsmuth (Eds.), *Proceedings of the 35th annual conference of the cognitive*
662 *science society* (pp. 2506–2511). Berlin: Cognitive Science Society.

- 663 Healey, P. G. T., Plant, N., Howes, C., & Lavelle, M. (2015). When words fail:
664 Collaborative gestures during clarification dialogues. In S. Andrist, D. Bohus,
665 E. Horvitz, B. Mutlu, & D. Schlangen (Eds.), *2015 aai spring symposium series:
666 Turn-taking and coordination in human-machine interaction* (pp. 23–29).
667 Stanford, CA: AAAI Press.
- 668 Hjalmarsson, A., & Oertel, C. (2012). Gaze direction as a back-channel inviting cue in
669 dialogue. In J. Edlund, I. de Kok, R. Poppe, & D. Traum (Eds.), *Iva 2012
670 workshop on realtime conversational virtual agents* (Vol. 9). Santa Cruz, CA.
- 671 Honnibal, M., & Johnson, M. (2014). Joint incremental disfluency detection and
672 dependency parsing. *Transactions of the Association of Computational Linguistics
673 (TACL)*, 2, 131-142.
- 674 Hough, J. (2015). *Modelling incremental self-repair processing in dialogue* (Unpublished
675 doctoral dissertation). Queen Mary University of London.
- 676 Hough, J., & Purver, M. (2012, September). Processing self-repairs in an incremental
677 type-theoretic dialogue system. In S. Brown-Schmidt, J. Ginzburg, & S. Larsson
678 (Eds.), *Proceedings of the 16th SemDial workshop on the semantics and
679 pragmatics of dialogue (SeineDial)* (pp. 136–144). Paris, France: SemDial.
- 680 Hough, J., & Purver, M. (2013, December). Modelling expectation in the self-repair
681 processing of annotated, un-, listeners. In R. Fernández & A. Isard (Eds.),
682 *Proceedings of the 17th SemDial workshop on the semantics and pragmatics of
683 dialogue (DialDam)* (pp. 92–101). Amsterdam: SemDial.
- 684 Hough, J., & Purver, M. (2014). Strongly incremental repair detection. In A. Moschitti,
685 B. Pang, & W. Daelemans (Eds.), *Proceedings of the 2014 conference on empirical
686 methods in natural language processing (emnlp)*. Doha, Qatar: Association for
687 Computational Linguistics.
- 688 Hough, J., & Schlangen, D. (2015). Recurrent neural networks for incremental
689 disfluency detection. *Interspeech 2015*.
- 690 Howes, C., Hough, J., Purver, M., & McCabe, R. (2014, September). Helping, i mean
691 assessing psychiatric communication: An application of incremental self-repair

- 692 detection. In V. Rieser & P. Muller (Eds.), *Proceedings of the 18th SemDial*
693 *workshop on the semantics and pragmatics of dialogue (DialWatt)* (pp. 80–89).
694 Edinburgh: SemDial.
- 695 Howes, C., Lavelle, M., Healey, P. G. T., Hough, J., & McCabe, R. (2017). Disfluencies
696 in dialogues with patients with schizophrenia. In G. Gunzelmann, A. Howes,
697 T. Tenbrink, & E. Davelaar (Eds.), *Proceedings of the 39th annual meeting of the*
698 *cognitive science society*. London, UK: Cognitive Science Society.
- 699 Howes, C., Purver, M., Healey, P. G. T., Mills, G. J., & Gregoromichelaki, E. (2011).
700 On incrementality in dialogue: Evidence from compound contributions. *Dialogue*
701 *& Discourse*, 2(1), 279–311. Retrieved from
702 <http://dad.uni-bielefeld.de/index.php/dad/article/view/362> doi:
703 10.5087/dad.2011.111
- 704 Howes, C., Purver, M., McCabe, R., Healey, P. G. T., & Lavelle, M. (2012). Helping
705 the medicine go down: Repair and adherence in patient-clinician dialogues. In
706 S. Brown-Schmidt, J. Ginzburg, & S. Larsson (Eds.), *Proceedings of the 16th*
707 *SemDial workshop on the semantics and pragmatics of dialogue (SeineDial)* (pp.
708 155–156). Paris: SemDial.
- 709 Johnson, M., & Charniak, E. (2004). A tag-based noisy channel model of speech
710 repairs. In *Proceedings of the 42nd annual meeting on association for*
711 *computational linguistics*. Stroudsburg, PA, USA: Association for Computational
712 Linguistics. Retrieved from <http://dx.doi.org/10.3115/1218955.1218960>
713 doi: <http://dx.doi.org/10.3115/1218955.1218960>
- 714 Jurafsky, D., Shriberg, E., & Biasca, D. (1997). *Switchboard subd-damsl*
715 *shallow-discourse-function annotation coders manual, draft 13* (Tech. Rep. No.
716 97-02). University of Colorado, Boulder. Institute of Cognitive Science.
- 717 Kitaoka, N., Kakutani, N., & Nakagawa, S. (2005). Detection and recognition of
718 correction utterances on misrecognition of spoken dialog system. *Systems and*
719 *Computers in Japan*, 36(11), 24–33. Retrieved from
720 <http://dx.doi.org/10.1002/scj.20341> doi: 10.1002/scj.20341

- 721 Lake, J. K., Humphreys, K. R., & Cardy, S. (2011). Listener vs. speaker-oriented
722 aspects of speech: Studying the disfluencies of individuals with autism spectrum
723 disorders. *Psychonomic bulletin & review*, *18*(1), 135–140.
- 724 Lease, M., Johnson, M., & Charniak, E. (2006). Recognizing disfluencies in
725 conversational speech. *Audio, Speech, and Language Processing, IEEE
726 Transactions on*, *14*(5), 1566–1573.
- 727 Lemon, O., & Gruenstein, A. (2004). Multithreaded context for robust conversational
728 interfaces: Context-sensitive speech recognition and interpretation of corrective
729 fragments. *ACM Transactions on Computer-Human Interaction*, *11*(3), 1–27.
- 730 Leudar, I., Thomas, P., & Johnston, M. (1992). Self-repair in dialogues of
731 schizophrenics: Effects of hallucinations and negative symptoms. *Brain and
732 Language*, *43*(3), 487 - 511.
- 733 Levelt, W. (1983). Monitoring and self-repair in speech. *Cognition*, *14*(1), 41–104.
- 734 Levelt, W. (1989). *Speaking: From intention to articulation*. MIT Press.
- 735 Lickley, R. J. (2001). Dialogue moves and disfluency rates. In *Isca tutorial and research
736 workshop (itrw) on disfluency in spontaneous speech* (pp. 93–96). Edinburgh:
737 ISCA.
- 738 Litman, D., Hirschberg, J., & Swerts, M. (2006). Characterizing and predicting
739 corrections in spoken dialogue systems. *Computational Linguistics*, *32*(3),
740 417–438.
- 741 Lopes, J., Salvi, G., Skantze, G., Abad, A., Gustafson, J., Batista, F., . . . Trancoso, I.
742 (2015, September). Detecting repetitions in spoken dialogue systems using
743 phonetic distances. In S. Möller, H. Ney, B. Möbius, E. Nöth, & S. Steid (Eds.),
744 *Proceedings of interspeech* (pp. 1805–1809). Dresden: ISCA.
- 745 McCabe, R. (2008). *Doctor-patient communication in the treatment of schizophrenia: Is
746 it related to treatment outcome?* (Tech. Rep.). Final report on G0401323 to
747 Medical Research Council.
- 748 McCabe, R., Healey, P. G. T., Priebe, S., Lavelle, M., Dodwell, D., Laugharne, R., . . .
749 Bremner, S. (2013). Shared understanding in psychiatrist-patient communication:

- 750 Association with treatment adherence in schizophrenia. *Patient Education and*
751 *Counselling*.
- 752 Meteer, M., Taylor, A., MacIntyre, R., & Iyer, R. (1995). *Disfluency annotation*
753 *stylebook for the Switchboard Corpus* (Tech. Rep.). Department of Computer and
754 Information Science, University of Pennsylvania. Retrieved from
755 <ftp://ftp.cis.upenn.edu/pub/treebank/swbd/doc/DFL-book.ps>
- 756 Mieskes, M., & Strube, M. (2006). Part-of-speech tagging of transcribed speech. In
757 N. Calzolari et al. (Eds.), *Proceedings of Irec* (pp. 935–938). Genoa: LREC.
- 758 Mikolov, T., Yih, W.-t., & Zweig, G. (2013). Linguistic regularities in continuous space
759 word representations. In *Proceedings of naacl-hlt* (pp. 746–751).
- 760 Mills, G. (2013, December). Establishing a communication system: Miscommunication
761 drives abstraction. In R. Fernández & A. Isard (Eds.), *Proceedings of the 17th*
762 *SemDial workshop on the semantics and pragmatics of dialogue (DialDam)*.
763 Amsterdam: SemDial.
- 764 Mills, G., & Healey, P. G. T. (2006, September). Clarifying spatial descriptions: Local
765 and global effects on semantic co-ordination. In *Proceedings of the 10th workshop*
766 *on the semantics and pragmatics of dialogue (SEMDIAL)*. Potsdam, Germany.
- 767 Ong, L., De Haes, J., Hoos, A., & Lammes, F. (1995). Doctor-patient communication:
768 a review of the literature. *Social science & medicine*, 40(7), 903–918.
- 769 Oviatt, S. (1995). Predicting spoken disfluencies during human–computer interaction.
770 *Computer Speech & Language*, 9(1), 19–35.
- 771 Purver, M., Ginzburg, J., & Healey, P. G. T. (2003). On the means for clarification in
772 dialogue. In R. Smith & J. van Kuppevelt (Eds.), *Current and new directions in*
773 *discourse & dialogue* (pp. 235–255). Kluwer Academic Publishers.
- 774 Rasooli, M. S., & Tetreault, J. (2014). Non-monotonic parsing of fluent umm I mean
775 disfluent sentences. In S. Wintner, S. Goldwater, & S. Riezler (Eds.), *Eacl 2014*
776 (pp. 48–53). Gothenburg: Association for Computational Linguistics.
- 777 Raux, A., Langner, B., Black, A., & Eskenazi, M. (2005). Let’s go public! taking a
778 spoken dialog system to the real world. In I. Trancoso (Ed.), *Proceedings of*

- 779 *interspeech 2005 (eurospeech)*. Lisbon, Portugal: ISCA.
- 780 Rieser, V., & Moore, J. (2005, June). Implications for generating clarification requests
781 in task-oriented dialogues. In K. Knight, H. T. Ng, & K. Oflazer (Eds.),
782 *Proceedings of the 43rd annual meeting of the ACL* (pp. 239–246). Ann Arbor.
- 783 Rodríguez, K., & Schlangen, D. (2004, July). Form, intonation and function of
784 clarification requests in German task-oriented spoken dialogues. In J. Ginzburg &
785 E. Vallduví (Eds.), *Proceedings of the 8th workshop on the semantics and*
786 *pragmatics of dialogue (SEMDIAL)* (pp. 101–108). Barcelona, Spain: SemDial.
- 787 Saito, T., & Rehmsmeier, M. (2015). The precision-recall plot is more informative than
788 the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLoS*
789 *ONE*, 10(3), e0118432. Retrieved from
790 <https://doi.org/10.1371/journal.pone.0118432>
- 791 San-Segundo, R., Montero, J. M., Ferreiros, J., Córdoba, R., & Pardo, J. M. (2001,
792 September). Designing confirmation mechanisms and error recover techniques in a
793 railway information system for Spanish. In D. Traum, J. Hirschberg, R. Smith, &
794 J. van Kuppevelt (Eds.), *Proceedings of the 2nd SIGdial workshop on discourse*
795 *and dialogue* (pp. 136–139). Aalborg, Denmark: Association for Computational
796 Linguistics.
- 797 Schegloff, E., Jefferson, G., & Sacks, H. (1977). The preference for self-correction in the
798 organization of repair in conversation. *Language*, 53(2), 361–382.
- 799 Schlangen, D. (2005, June). Towards finding and fixing fragments: Using machine
800 learning to identify non-sentential utterances and their antecedents in multi-party
801 dialogue. In K. Knight, H. T. Ng, & K. Oflazer (Eds.), *Proceedings of the 43rd*
802 *annual meeting of the association for computational linguistics (ACL)* (pp.
803 247–254). Ann Arbor, MI.
- 804 Shriberg, E. (1994). *Preliminaries to a theory of speech disfluencies* (Unpublished
805 doctoral dissertation). University of California, Berkeley.
- 806 Skantze, G., & Hjalmarsson, A. (2010, September). Towards incremental speech
807 generation in dialogue systems. In *Proceedings of the sigdial 2010 conference* (pp.

- 808 1–8). Tokyo, Japan: Association for Computational Linguistics. Retrieved from
809 <http://www.sigdial.org/workshops/workshop11/proc/pdf/SIGDIAL01.pdf>
- 810 Stivers, T., & Enfield, N. J. (2010). A coding scheme for question–response sequences
811 in conversation. *Journal of Pragmatics*, 42(10), 2620–2626.
- 812 Stolcke, A., Ries, K., Coccaro, N., Shriberg, E., Bates, R., Jurafsky, D., . . . Meteer, M.
813 (2000). Dialogue act modeling for automatic tagging and recognition of
814 conversational speech. *Computational Linguistics*, 26(3), 339–373.
- 815 Stoyanchev, S., Liu, A., & Hirschberg, J. (2014, April). Towards natural clarification
816 questions in dialogue systems. In *Aisb symposium on questions, discourse and*
817 *dialogue: 20 years after Making it Explicit*.
- 818 Stoyanchev, S., & Stent, A. (2012). Concept type prediction and responsive adaptation
819 in a dialogue system. *Dialogue & Discourse*, 3(1), 1–31.
- 820 Surendran, D., & Levow, G.-A. (2006). Dialog act tagging with support vector
821 machines and hidden Markov models. In R. M. Stern (Ed.), *Proceedings of*
822 *interspeech*. Pittsburgh, PA: ISCA.
- 823 Toutanova, K., Klein, D., Manning, C., & Singer, Y. (2003). Feature-rich part-of-speech
824 tagging with a cyclic dependency network. In M. Hearst & M. Ostendorf (Eds.),
825 *Proceedings of HLT-NAACL* (pp. 252–259). Edmonton: Association for
826 Computational Linguistics.
- 827 Turian, J., Ratinov, L.-A., & Bengio, Y. (2010, July). Word representations: A simple
828 and general method for semi-supervised learning. In *Proceedings of the 48th*
829 *annual meeting of the association for computational linguistics* (pp. 384–394).
830 Uppsala, Sweden: Association for Computational Linguistics. Retrieved from
831 <http://www.aclweb.org/anthology/P10-1040>
- 832 Weng, F., Yan, B., Feng, Z., Ratiu, F., Raya, M., Lathrop, B., . . . Peters, S. (2007,
833 September). CHAT to your destination. In *Proceedings of the 8th SIGdial*
834 *workshop on discourse and dialogue* (p. 79-86). Antwerp, Belgium. Retrieved from
835 [http://godel.stanford.edu/twiki/pub/Public/SemlabPublications/](http://godel.stanford.edu/twiki/pub/Public/SemlabPublications/weng-et-al07sigdial.pdf)
836 [weng-et-al07sigdial.pdf](http://godel.stanford.edu/twiki/pub/Public/SemlabPublications/weng-et-al07sigdial.pdf)

- 837 Young, S., Gašić, M., Thomson, B., & Williams, J. D. (2013, May). POMDP-based
838 statistical spoken dialog systems: A review. *Proceedings of the IEEE*, 101(5),
839 1160-1179. doi: 10.1109/JPROC.2012.2225812
- 840 Zwarts, S., Johnson, M., & Dale, R. (2010). Detecting speech repairs incrementally
841 using a noisy channel approach. In *Proceedings of the 23rd international*
842 *conference on computational linguistics* (pp. 1371–1378). Stroudsburg, PA, USA:
843 Association for Computational Linguistics. Retrieved from
844 <http://portal.acm.org/citation.cfm?id=1873781.1873935>

PRE-PUBLICATION DRAFT

Appendix

845

846 **Materials for Replication**

847 The PCC corpus is confidential due to its sensitive nature; all other data and
848 experiment processing scripts are publicly available. The scripts for the other-repair
849 experiments can be accessed via the Open Science Framework at
850 <http://osf.io/w4dmz>; scripts and pre-processed data for the self-repair experiments
851 can be accessed via the git repository <http://bitbucket.org/julianhough/stir>. The
852 original datasets can be obtained as follows:

- 853 • **SWBD**: The original corpus is available from
854 http://www.stanford.edu/~jurafsky/swb1_dialogact_annot.tar.gz; we also
855 used the associated Python package available at
856 <http://compprag.christopherpotts.net/swda.html>.
- 857 • **BNC**: The original corpus is available from <http://purl.ox.ac.uk/ota/2554>.
858 The BNC-PGH and BNC-CH annotations are included with our experiment
859 scripts on the OSF.
- 860 • **MAPTASK**: The original corpus is available from
861 <http://groups.inf.ed.ac.uk/maptask/>; we used the V2.1 NXT format
862 annotations.

Supplementary Material

Experimental Details: Other-Repair

Our turn-level NTRI classifier uses two categories of features. One is designed to capture characteristic surface properties of NTRIs: presence of wh-words, specific clarification keywords (e.g. “pardon”), and behaviour associated with repair such as fillers, pauses and overlaps. The second is designed to capture parallelism. Lexical parallelism is captured via simple word string matching; syntactic parallelism by matching POS tags. For semantic parallelism, we measure word and turn similarity using distributed vector representations from two neural models, (Mikolov, Yih, & Zweig, 2013)’s word2vec with 300 dimensions trained on the Google News corpus, and Turian, Ratinov, and Bengio (2010)’s implementation of (Collobert & Weston, 2008) with 100 dimensions trained on the Reuters RCV1 corpus. The full set of features is given in Table 1.

For the majority of features, we extracted one raw feature (the numeric count or binary indicator, see Table 1) and one proportional feature (the proportion of the turn made up of the feature in question, from 0 to 1). For vector-based similarity features, we extracted four features (minimum, mean and maximum pairwise word cosine similarities, and overall turn cosine similarity summing word vectors within turns following Mitchell and Lapata (2010)). Parallelism features were calculated between the turn being classified, and the preceding turns by the same and other speaker separately. The features ranked in terms of information gain (using Weka’s default implementation with best-first search) are shown in Table 6 – see above for discussion.

We experimented with logistic regression, decision tree and support vector machine (SVM) classifiers, as implemented in the Weka toolkit (Hall et al., 2009). Results given here use SVMs with radial basis function kernels; logistic regression and linear-kernel SVMs performed very similarly for all datasets other than SWBD, while decision trees were usually worse. For each dataset, we used 10-fold cross-validation (using Weka’s built-in stratified cross-validation routines): the dataset is randomised and split into 10 equally-sized parts, and the classifier tested on each 10% part in turn

892 while training on the other 90%.

To combat sparsity, we weighted the classifiers to give equal precision and recall for the rare class of interest (here, NTRIs), by either adjusting the decision threshold directly (for logistic regression) or by weighting the rarer class more highly in the training error function (for support vector machines). Unweighted versions gave very low recall due to the rarity of NTRIs. For the latter, we used Weka’s built-in class weighting function; the underlying implementations vary for particular classifier types, but share the intuition shown in (2,3). Here, C is a single global cost parameter, I is the set of training instances and ϕ_i the individual error for any given instance under a particular model’s prediction; C_+, C_- are weighted cost parameters for positive and negative classes respectively, and I_+, I_- the sets of associated instances.

$$\text{standard error term: } C \sum_{i \in I} \phi_i \quad (2)$$

$$\text{weighted error term: } C_+ \sum_{i \in I_+} \phi_i + C_- \sum_{i \in I_-} \phi_i \quad (3)$$

893 Both methods allow precision to be traded off against recall, and the best choice in
894 practice will depend on application and aims; we give results for the point where
895 precision and recall are equal, and show the area under the precision-recall curve as a
896 measure which abstracts away from the exact setting.

897 **Experimental Details: Self-Repair**

898 Our word- (and partial-word-)level self-repair classifier uses STIR, with features
899 derived from n-gram language models (LMs). The basic ‘fluent’ LMs are trigram LMs
900 with Kneser-Ney smoothing (Kneser & Ney, 1995), trained on the standard
901 Switchboard training data cleaned of all edit terms and reparanda, giving a total of
902 $\approx 100\text{K}$ utterances, $\approx 600\text{K}$ tokens. We train one LM for words and one for POS tags.³²
903 We call their probabilities p^{lex} and p^{pos} respectively below; if referring to the same
904 calculation for both models we suppress the superscripts.

³²In pre-processing, POS tags in a many-to-one relation to words are concatenated into one token; this had no significant effect on results.

905 From the basic probability values we derive our principal lexical uncertainty
 906 measurement *surprisal* s (equation 4); and, following Clark et al. (2013), the (unigram)
 907 Weighted Mean Log trigram probability (*WML*, eq. 5) – this factors out lexical
 908 frequency to approximate *incremental syntactic probability*.

$$s(w_{i-2} \dots w_i) = -\log_2 p(w_i | w_{i-2}, w_{i-1}) \quad (4)$$

$$909 \quad WML(w_0 \dots w_n) = \frac{\sum_{i=2}^n \log_2 p(w_i | w_{i-2}, w_{i-1})}{-\sum_{j=2}^n \log_2 p(w_j)} \quad (5)$$

$$910 \quad H(w|c) = - \sum_{w \in Vocab} p(w|c) \log_2 p(w|c) \quad \text{where } c = w_{i-2}, w_{i-1} \quad (6)$$

911 As a measure of uncertainty, we can then derive the entropy $H(w|c)$ of possible word
 912 continuations w given a context c , from $p(w_i|c)$ for all words w_i in the vocabulary – see
 913 (6). Calculating distributions over the entire lexicon incrementally is costly, so this is
 914 approximated by calculating directly only for words observed at least once in context c
 915 in training, assuming a uniform distribution over unseen suffixes (see Hough & Purver,
 916 2014). We can then measure increases and reductions in entropy, and similarity between
 917 distributions in two different contexts c_1 and c_2 via the Kullback-Leibler (KL)
 918 divergence (relative entropy, using a similar approximation).

919 **Adaptation for processing partial words.** We adapt our language model
 920 scoring when there is a partial, cut-off word transcribed as the penultimate word in any
 921 n-gram (in our case, the second word of any trigram). This captures the idea of the
 922 fluency dropping after the cut-off word has been processed. We simply assign such
 923 trigrams a minimum probability $\frac{1}{|V|}$ where $|V|$ is the vocabulary size.

924 While this simple method gives good results in practice, we have also developed a
 925 more principled off-line model. We train a simple word completion model $p^{complete}(w|w_i)$
 926 which operates on any annotated partial word prefix w_i to provide a distribution over
 927 possible complete words that it could have started, and thus also the most likely
 928 completion (based on the prefix and unigram co-occurrence). This is combined with the
 929 language model probability p^{lex} within the function p^{fluent} , which for a partial word w_i ,
 930 gives the likelihood of a given word w being its corresponding complete word at the
 931 time of interruption given its two word context is as in (7).

$$p^{fluent}(w \mid w_{i-2}, w_{i-1}, w_i) = \frac{1}{Z} \times p^{lex}(w \mid w_{i-2}, w_{i-1}) \times p^{complete}(w \mid w_i) \quad (7)$$

where Z is a standard normalisation constant to ensure: $\sum_{w \in Vocab} p^{fluent}(w \mid w_{i-2}, w_{i-1}, w_i) = 1$

932 The probability $p^{\hat{fluent}}$ of the most likely complete word guess for w_i is therefore:

$$p^{\hat{fluent}}(w \mid w_{i-2}, w_{i-1}, w_i) = \max_w p^{fluent}(w \mid w_{i-2}, w_{i-1}, w_i) \quad (8)$$

933 The intuition here is that when hearers encounter a partial word, they attempt to
 934 find the fluent word most likely to both complete the partial word and follow the two
 935 preceding words. The probability of a completion ‘remember’ will be higher after “Yes I
 936 remem-” than in a less predictable context e.g. utterance-initial “Re-”.

937 **Classifiers and features.** The 4 individual classifiers in STIR then use
 938 combinations of features derived from these basic measures.

939 **Edit term detection.** We use the word surprisal s^{lex} from a specific edit word
 940 bigram model (expecting low s^{lex} for words likely to be edit terms), and the trigram
 941 surprisal s and syntactic fluency WML from the standard fluent word and POS models
 942 described above. The decision task is to classify at the current position w_i , one, both or
 943 none of words w_i and w_{i-1} as edit terms. We found this simple approach effective and
 944 stable, detecting edit term words with an F-score of 0.938, performing marginally worse
 945 though detecting a broader range of phenomena than [Heeman and Allen \(1999\)](#)’s
 946 discourse marker detector. Some delayed decisions occur in cases where s^{lex} and
 947 WML^{lex} have similar values in both the edit and fluent language models before the end
 948 of the edit, e.g. “I like” \rightarrow “I *{like}* want...”, with classification only achieved at w_{i-1} ;
 949 this could cause some output instability or ‘jitter’.

950 **Repair start detection.** Starting with s , WML , H for word and POS models,
 951 we derive 5 additional information-theoretic features: ΔWML is the difference between
 952 the WML values at w_{i-1} and w_i ; ΔH is the difference in entropy between w_{i-1} and w_i ;
 953 *InformationGain* is the difference between expected entropy at w_{i-1} and observed s at
 954 w_n , a measure that factors out the effect of naturally high entropy contexts;

955 *BestEntropyReduce* is the best reduction in entropy possible by an early rough
 956 hypothesis of reparandum onsets within 3 words; and *BestWMLboost* similarly
 957 speculates on the best improvement of *WML* possible by positing rm_{start} positions up to
 958 3 words back. We also include simple alignment features: binary features which indicate
 959 if the word w_{i-x} is identical to the current word w_i for $x \in \{1, 2, 3\}$. With 6 alignment
 960 features, 16 language model information-theoretic features and a single logical feature
 961 *edit* which indicates the presence of an edit word at position w_{i-1} , rp_{start} detection uses
 962 23 features— see Table 7.

963 Ranking the features by Information Gain using 10-fold cross validation over the
 964 Switchboard heldout data (see Table 7) shows that the language model features are far
 965 more discriminative than the alignment features, with *WML* in both p^{lex} and p^{pos}
 966 models being the most discriminative. Actual lexical or POS *values* (i.e. words and POS
 967 tags) are not used at all in the feature sets, only these information-theoretic measures.

968 ***Reparandum start detection.*** Altogether we use 32 features, and again
 969 information-theoretic ones are most useful. The two best features capture the noisy
 970 channel intuition that removing the reparandum increases fluency: they are
 971 $\Delta WMLboost$ (the drop in *WMLboost* from one backtracked position to the next) for
 972 word and POS models. The third best feature measures parallelism between rp_{start} and
 973 rm_{start} , via the KL divergence between $\theta^{pos}(w \mid rm_{start}, rm_{start-1})$ and
 974 $\theta^{pos}(w \mid rp_{start}, rp_{start-1})$.

975 ***Repair end detection and structure classification.*** Finally, detection of
 976 rp_{end} uses parallelism, measured as KL divergence between the conditional probability
 977 distribution θ^{lex} at the reparandum-final word rm_{end} and the repair-final word rp_{end} .
 978 For repetition repairs, divergence will trivially be 0; for substitutions, it will be higher;
 979 for deletes, even higher. It can also be captured via *ReparandumRepairDifference*, the
 980 difference in probability between an utterance cleaned of the reparandum and the
 981 utterance with its repair phase substituting its reparandum. In the running example
 982 from Figure 1, this would be as in equation (9).

$$\text{ReparandumRepairDifference}(\text{"John [likes + loves]"}) = \\ WML^{lex}(\text{"John loves"}) - WML^{lex}(\text{"John likes"}) \quad (9)$$

983 **Classifier pipeline and training setup.** The classifiers are implemented
 984 using Random Forests (Breiman, 2001), using different error functions for each stage via
 985 MetaCost (Domingos, 1999); in early investigation this outperformed single decision
 986 tree classifiers. The LM-derived features are obtained using a 10-fold cross-validation
 987 method, always using the SBWD training set: for each fold, we train the LMs on 90%
 988 and use them to calculate feature values on the unseen 10% (this avoids over-fitting
 989 probability values). We then use these feature values to train and test the classifiers
 990 using a standard single development/test split for each corpus .

991 In both training and testing, the classifiers are combined in a pipeline as in
 992 Figure 2, where the *ed* classifier only permits non-*ed* words to be passed on to *rp_{start}*
 993 classification. The *rp_{start}* classifier passes positive repair hypotheses to the *rm_{start}*
 994 classifier, which searches backwards up to 7 words which have not been classified as edit
 995 terms *e*. If a *rm_{start}* is classified, the output is passed on for *rp_{end}* classification. Active
 996 repair hypotheses are added to a stack, each consisting of a $\langle rm_{start}, rp_{start}, rp_{end} \rangle$ triple
 997 of word positions; the position of *rp_{end}* may change as more words are consumed. The
 998 *rp_{end}* detector may temporarily *cancel* a hypothesis after two words have been consumed
 999 beyond the repair onset, which does not remove the hypothesis indefinitely but subdues
 1000 its effect in its output before searching for more suitable *rp_{end}* points— this could cause
 1001 output jitter. Repair hypotheses are are popped off the stack when the string is 7 words
 1002 beyond its *rp_{start}* position. Putting limits on the stack’s storage space is a way of
 1003 controlling for processing overhead and complexity. Embedded repairs whose *rm_{start}*
 1004 coincide with another’s *rp_{start}* are easily dealt with as they are added to the stack as
 1005 separate hypotheses.³³ In terms of complexity, the number of possible repairs grows
 1006 approximately in the triangular number series— i.e. $\frac{n(n+1)}{2}$, a nested loop over previous
 1007 words as *n* gets incremented, which in terms of a complexity class is a quadratic $O(n^2)$.

³³We constrain the problem not to include embedded deletes which may share their *rp_{start}* word with another repair – these are in practice very rare.

References

1008

- 1009 Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5–32.
- 1010 Clark, A., Giorgolo, G., & Lappin, S. (2013, August). Statistical representation of
1011 grammaticality judgements: the limits of n-gram models. In *Proceedings of the*
1012 *fourth annual workshop on cognitive modeling and computational linguistics*
1013 *(cmcl)* (pp. 28–36). Sofia, Bulgaria: Association for Computational Linguistics.
1014 Retrieved from <http://www.aclweb.org/anthology/W13-2604>
- 1015 Collobert, R., & Weston, J. (2008). A unified architecture for natural language
1016 processing: Deep neural networks with multitask learning. In *Proceedings of the*
1017 *25th international conference on machine learning* (pp. 160–167).
- 1018 Domingos, P. (1999). Metacost: A general method for making classifiers cost-sensitive.
1019 In *Proceedings of the 5th ACM SIGKDD international conference on knowledge*
1020 *discovery and data mining* (pp. 155–164).
- 1021 Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., & Witten, I. H. (2009).
1022 The WEKA data mining software: An update. *SIGDKDD Explorations*, 11(1),
1023 10–18.
- 1024 Heeman, P., & Allen, J. (1999). Speech repairs, intonational phrases, and discourse
1025 markers: modeling speakers' utterances in spoken dialogue. *Computational*
1026 *Linguistics*, 25(4), 527–571.
- 1027 Hough, J., & Purver, M. (2014). Strongly incremental repair detection. In A. Moschitti,
1028 B. Pang, & W. Daelemans (Eds.), *Proceedings of the 2014 conference on empirical*
1029 *methods in natural language processing (emnlp)*. Doha, Qatar: Association for
1030 Computational Linguistics.
- 1031 Kneser, R., & Ney, H. (1995). Improved backing-off for m-gram language modeling. In
1032 *Acoustics, speech, and signal processing, 1995. icassp-95., 1995 international*
1033 *conference on* (Vol. 1, pp. 181–184).
- 1034 Mikolov, T., Yih, W.-t., & Zweig, G. (2013). Linguistic regularities in continuous space
1035 word representations. In *Proceedings of naacl-hlt* (pp. 746–751).
- 1036 Mitchell, J., & Lapata, M. (2010). Composition in distributional models of semantics.

- 1037 *Cognitive Science*, 34(8), 1388–1439.
- 1038 Turian, J., Ratinov, L.-A., & Bengio, Y. (2010, July). Word representations: A simple
1039 and general method for semi-supervised learning. In *Proceedings of the 48th*
1040 *annual meeting of the association for computational linguistics* (pp. 384–394).
1041 Uppsala, Sweden: Association for Computational Linguistics. Retrieved from
1042 <http://www.aclweb.org/anthology/P10-1040>

PRE-PUBLICATION DRAFT

Measure	Description
NumWords	Number of words in turn
OpenClassRepair	Number of Open Class Repair Initiator words (e.g. <i>pardon, huh</i>)
WhWords	Number of wh-words (e.g. <i>what, who, when</i>)
Backchannel	Number of backchannels (e.g. <i>uh-huh, yeah</i>)
FillerWords	Number of fillers (e.g. <i>er, um</i>)
MarkedPauses	Number of pauses transcribed
OverlapAny	Number of portions of overlapping talk
OverlapAll	Entirely overlapping another turn
RepeatedWords	Number of words repeated from preceding turn
RepeatedPos	Number of PoS-tags repeated from preceding turn
W2vSim	Cosine similarity with preceding turn (word2vec, Mikolov et al., 2013)
TeaSim	Cosine similarity with preceding turn (Turian et al., 2010)

Table 1

Turn-level features for NTRI detection

Corpus	Features	P	R	F	AUC-PRC	AUC-ROC
PCC patient	OCRProportion	.86	.23	.36	-	-
PCC patient	High-level	.43	.41	.41	-	-
PCC patient	All	.45	.44	.44	-	-
PCC all	(baseline)	.01	1.0	.03	.01	.50
PCC all	qMarkProportion	.65	.14	.24	.11	.57
PCC all	High-level	.44	.43	.44	.40	.90
PCC all	All	.46	.47	.46	.43	.73
BNC-CH	(baseline)	.04	1.0	.08	.04	.49
BNC-CH	qMarkProportion	.62	.31	.41	.22	.65
BNC-CH	High-level	.55	.55	.55	.52	.90
BNC-CH	All	.57	.62	.60	.61	.80
BNC-PGH	(baseline)	.04	1.0	.08	.04	.50
BNC-PGH	OCRProportion	.70	.09	.16	.10	.55
BNC-PGH	High-level	.52	.53	.52	.51	.92
BNC-PGH	All	.61	.52	.56	.56	.75
MapTask	(baseline)	.11	1.0	.20	.11	.50
MapTask	TeaSimSum	.29	.03	.06	.12	.51
MapTask	High-level	.38	.38	.38	.34	.81
MapTask	All	.41	.63	.50	.55	.76
Switchboard	(baseline)	.002	1.0	.005	.002	.50
Switchboard	OCRProportion	0	0	0	0	.50
Switchboard	High-level	.54	.52	.53	.50	.98
Switchboard	All	.52	.60	.56	.58	.80

Table 2

NTRI detection: Precision, Recall, F-score and Area Under Curve (AUC) metrics for NTRI utterances, using 10-fold cross-validation. We show AUC for the precision-recall curve for NTRIs (AUC-PRC) as well as the more usual receiver-operator curve (AUC-ROC); AUC-PRC is more informative with unbalanced data.

Corpus	Features	P	R	F	Correl.	F <i>rm</i>
PCC all	words	.648	.555	.598	.798**	-
PCC all	words+POS	.660	.585	.620	.805**	-
BNC-CH	words	.350	.446	.392	.530**	-
BNC-CH	words+POS	.397	.438	.417	.583**	-
Switchboard	words	.910	.758	.827	.962**	.749
Switchboard	words+POS	.928	.785	.850	.956**	.781

Table 3

*Self-repair detection: STIR's per-utterance performance on our corpora in terms of rp_{start} (repair onset) detection and the Spearman's rank correlation between STIR and the annotators' repair rates (rp_{start} per utterance) per speaker (**= $p < 0.001$). The reparandum word detection accuracy is also given for Switchboard.*

System (evaluation)	F rm (word)	F rp_{start}	Correl.
RNN (+ partial)	0.631	0.751	0.948**
RNN (- partial)	0.668	0.790	0.956**
STIR +POS (+ partial)	0.781	0.850	0.956**
STIR +POS (- partial)	0.768	0.833	0.937**

Table 4

The effect of partial words: Comparison of STIR's performance to an RNN disfluency tagger testing on Switchboard heldout data with and without partial words. STIR improves whilst the RNN suffers with partial words.

PRE-PUBLICATION DRAFT

(a)

Reparandum + Interregnum length	(support)	RNN	STIR (-POS)	STIR (+POS)
1	(1254)	.756	.852	.874
2	(531)	.590	.730	.782
3	(227)	.397	.600	.688
4	(106)	.286	.533	.559
5	(50)	.098	.370	.430
6	(25)	.000	.308	.500
7	(11)	.000	.000	.154
8	(6)	.000	.250	.286

(b)

Repair Type	(support)	RNN	STIR (-POS)	STIR (+POS)
repetition	(1022)	.923	.970	.969
substitution	(1061)	.536	.708	.759
delete	(132)	.366	.453	.407

Table 5

Self-repair detection error analysis: (a) F-score for detecting the correct repair start word and reparandum start word of repairs with different combined reparandum and interregnum lengths; (b) F-score for detecting repair onset word of different types. Compared with off-the-shelf RNN disfluency tagger on the SWBD held-out data.

Rank	BNC-PGH	SWBD	MAPTASK
1	qMarkProportion	qMarkProportion	BackChProportion
2	qMarks	qMarks	BackChWholeTurn
3	WhProportion	WhProportion	NumWords
4	WhWords	OCRProportion	TeaSimSum
5	OCRProportion	OpenClassRepair	RepeatedProportion
6	NumWords	SelfW2vSimSum	RepeatedPos
7	OpenClassRepair	NumWords	RepeatedWords
8	RepeatedProportion	BackChProportion	W2vSimSum
9	SelfW2vSimSum	NumBackchannel	TeaSimMax
10	NumBackchannel	TeaSimMean	W2vSimMax
11	SelfTeaSimSum	W2vSimMean	RepeatedPosProportion
12	BackChProportion	W2vSimSum	RepeatedSelfPosProportion
13	SelfTeaSimMax	RepeatedSelfPos	W2vSimMin
14	SelfW2vSimMax	TeaSimMax	TeaSimMin
15	RepeatedPosProportion	W2vSimMax	NumBackchannel
16	RepeatedSelfPos	SelfTeaSimMax	W2vSimMean
17	RepeatedWords	SelfW2vSimMax	TeaSimMean
18	RepeatedPos	RepeatedSelfWords	RepeatedSelfPos
19	SelfW2vSimMin	WhWords	SelfTeaSimSum
20	SelfTeaSimMin	RepeatedPosProportion	SelfW2vSimSum
21	BackChWholeTurn	RepeatedProportion	RepeatedSelfProportion
22	TeaSimSum	RepeatedSelfProportion	SelfTeaSimMax
23	W2vSimMean	SelfTeaSimSum	SelfW2vSimMax
24	RepeatedSelfWords	SelfW2vSimMin	SelfTeaSimMin
25	TeaSimMean	SelfTeaSimMin	SelfW2vSimMin

Table 6

Feature ranker (Information Gain) for other-repair (NTRI) detection: top 15 features in order, using 10-fold cross-validation on BNC-PGH, SWBD and MAPTASK datasets.

Note that question marks (qMarks, qMarkProportion) are not transcribed in MAPTASK.

average merit	average rank	attribute
0.139 (+- 0.002)	1 (+- 0.00)	H^{pos}
0.131 (+- 0.001)	2 (+- 0.00)	WML^{pos}
0.126 (+- 0.001)	3.4 (+- 0.66)	WML^{lex}
0.125 (+- 0.003)	4 (+- 1.10)	s^{pos}
0.122 (+- 0.001)	5.9 (+- 0.94)	$w_{i-1} = w_i$
0.122 (+- 0.001)	5.9 (+- 0.70)	BestWMLboost ^{lex}
0.122 (+- 0.002)	5.9 (+- 1.22)	InformationGain ^{pos}
0.119 (+- 0.001)	7.9 (+- 0.30)	BestWMLboost ^{pos}
0.098 (+- 0.002)	9 (+- 0.00)	H^{lex}
0.08 (+- 0.001)	10.4 (+- 0.49)	ΔWML^{pos}
0.08 (+- 0.003)	10.6 (+- 0.49)	ΔH^{pos}
0.072 (+- 0.001)	12 (+- 0.00)	$POS_{i-1} = POS_i$
0.066 (+- 0.003)	13.1 (+- 0.30)	s^{lex}
0.059 (+- 0.000)	14.2 (+- 0.40)	ΔWML^{lex}
0.058 (+- 0.005)	14.7 (+- 0.64)	BestEntropyReduce ^{pos}
0.049 (+- 0.001)	16.3 (+- 0.46)	InformationGain ^{lex}
0.047 (+- 0.004)	16.7 (+- 0.46)	BestEntropyReduce ^{lex}
0.035 (+- 0.004)	18 (+- 0.00)	ΔH^{lex}
0.024 (+- 0.000)	19 (+- 0.00)	$w_{i-2} = w_i$
0.013 (+- 0.000)	20 (+- 0.00)	$POS_{i-2} = POS_i$
0.01 (+- 0.000)	21 (+- 0.00)	$w_{i-3} = w_i$
0.009 (+- 0.000)	22 (+- 0.00)	edit
0.006 (+- 0.000)	23 (+- 0.00)	$POS_{i-3} = POS_i$

Table 7

Feature ranker (Information Gain) for rp_{start} detection- 10-fold x -validation on Switchboard heldout data.

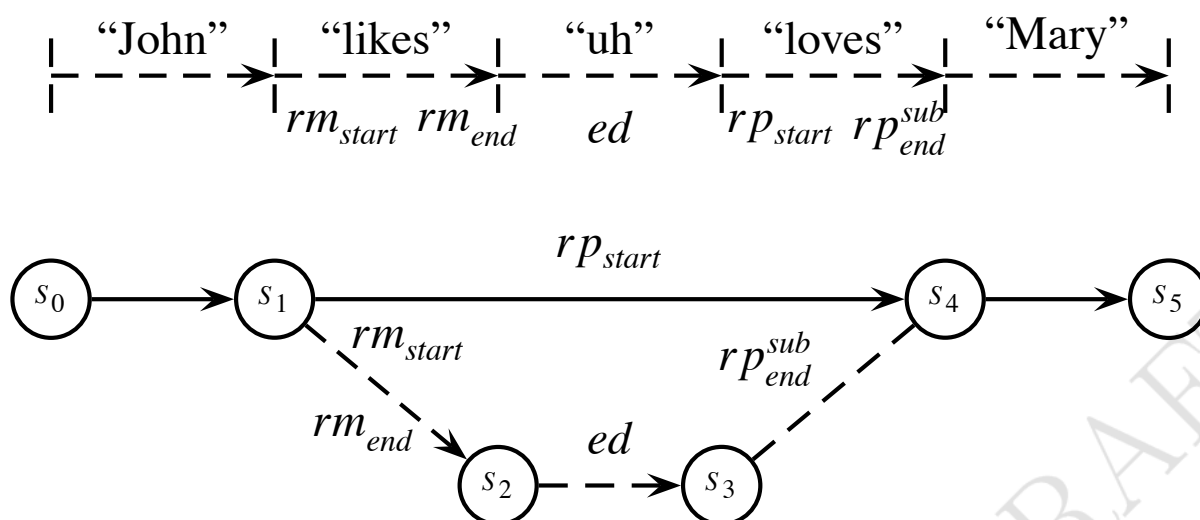


Figure 1. Strongly Incremental Repair detection (STIR); application to the utterance "John likes, uh, loves Mary", with incoming words and STIR's output tags at top.

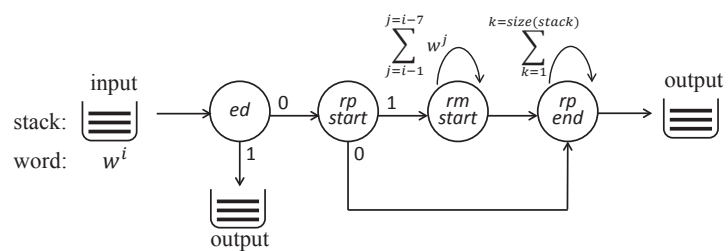


Figure 2. STIR's pipeline of classifiers

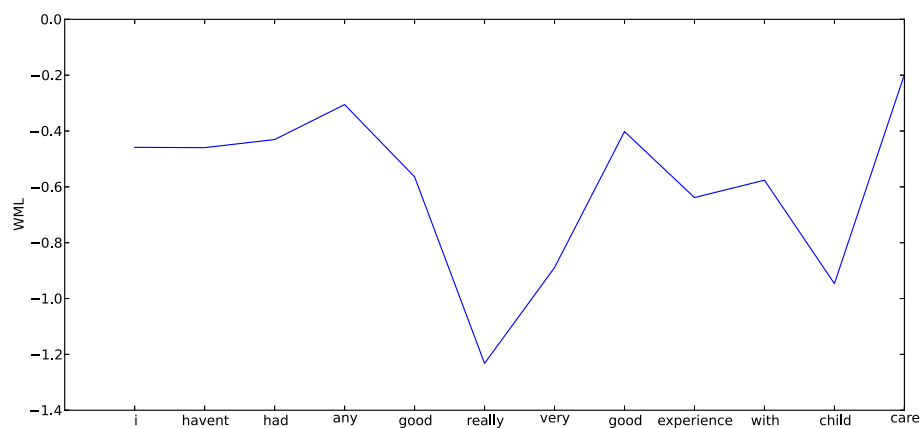


Figure 3. WML^{lex} values for trigrams for a repaired utterance exhibiting the drop at the repair onset