

# Split Utterances in Dialogue: a Corpus Study

**Matthew Purver, Christine Howes,  
and Patrick G. T. Healey**

Department of Computer Science  
Queen Mary University of London  
Mile End Road, London E1 4NS, UK  
{mpurver, chrizba, ph}@dcs.qmul.ac.uk

**Eleni Gregoromichelaki**

Department of Philosophy  
King's College London  
Strand, London WC2R 2LS, UK  
eleni.gregor@kcl.ac.uk

## Abstract

This paper presents a preliminary English corpus study of *split utterances (SUs)*, single utterances split between two or more dialogue turns or speakers. It has been suggested that SUs are a key phenomenon of dialogue, which this study confirms: almost 20% of utterances were found to fit this general definition, with nearly 3% being the between-speaker case most often studied. Other claims/assumptions in the literature about SUs' form and distribution are investigated, with preliminary results showing: splits can occur within syntactic constituents, apparently at any point in the string; it is unusual for the separate parts to be complete units in their own right; explicit repair of the antecedent does not occur very often. The theoretical consequences of these results for claims in the literature are pointed out. The practical implications for dialogue systems are mentioned too.

## 1 Introduction

*Split utterances (SUs)* – single utterances split between two or more dialogue turns/speakers – have been claimed to occur regularly in dialogue, especially according to the observations reported in the Conversational Analysis (CA) literature, which is based on the analysis of naturally occurring dialogues. SUs are of interest to dialogue theorists as they are a clear sign of how turns cohere with each other at all levels – syntactic, semantic and pragmatic. They also indicate the radical context-dependency of conversational contributions. Turns can, in general, be highly elliptical and nevertheless not disrupt the flow of the

dialogue. SUs are the most dramatic illustration of this: contributions spread across turns/speakers rely crucially on the dynamics of the unfolding context, linguistic and extra-linguistic, in order to guarantee successful processing and production.

Utterances that are split across speakers also present a canonical example of participant coordination in dialogue. The ability of one participant to continue another interlocutor's utterance coherently, both at the syntactic and the semantic level, suggests that both speaker and hearer are highly coordinated in terms of processing and production. The initial speaker must be able to switch to the role of hearer, processing and integrating the continuation of their utterance, whereas the initial hearer must be closely monitoring the grammar and content of what they are being offered so that they can take over and continue in a way that respects the constraints set up by the first part of the utterance. In fact there is (anecdotal) evidence that such constraints are fully respected across speaker and hearer in such utterances (see e.g. Gregoromichelaki et al. (2009)). A large proportion of the CA literature on SUs tries to identify the conditions under which SUs usually occur (see section 2). However, this emphasis seems to miss the important generalisation, confirmed by the present study, that, syntactically, a speaker switch may be able to occur anywhere in a string.

From a theoretical point of view, the implications of the above are that, if such observations have an empirical foundation, the grammar employed by the interlocutors must be able to license and the semantics interpret chunks much smaller than the usual sentence/proposition units. Moreover, these observations have implications for the nature of the grammar itself: dynamic, incremental formalisms seem more amenable to the mod-

elling of this phenomenon as the switch of roles while syntactic/semantic dependencies are pending can be taken as evidence for direct involvement of the grammar in the successful processing/production of such utterances. Indeed, Poesio and Rieser (to appear) claim that “[c]ollaborative completions ... are among the strongest evidence yet for the argument that dialogue requires *coordination* even at the sub-sentential level” (italics original).

From a psycholinguistic point of view, the phenomenon of SUs is compatible with mechanistic approaches as exemplified by the Interactive Alignment model of Pickering and Garrod (2004) where it is claimed that it should be as easy to complete someone else’s sentence as one’s own (Pickering and Garrod, 2004, p186). According to this model, speaker and listener ought to be interchangeable at any point. This is also the stance taken by the grammatical framework of Dynamic Syntax (DS) (Kempson et al., 2001; Cann et al., 2005). In DS, parsing and production are taken to employ the same mechanisms, leading to a prediction that split utterances ought to be strikingly natural (Purver et al., 2006). However, from a pragmatic point of view, utterance continuation by another speaker might involve some kind of guessing<sup>1</sup> or preempting the other interlocutor’s intended content. It has therefore been claimed that a full account of this phenomenon requires a complete model of pragmatics that can handle intention recognition and formation. Indeed, Poesio and Rieser (to appear) claim that “the study of sentence completions ... may be used to compare competing claims about coordination – i.e. whether it is best explained with an intentional model like Clark (1996)’s ... or with a model based on simpler alignment models like Pickering and Garrod (2004)’s.” They conclude that a model which includes modelling of intentions better captures the data.

For computational models of dialogue, however, SUs pose a challenge. While Poesio and Rieser (to appear) and Purver et al. (2006) provide general foundational models for various parts of the phenomenon, there are many questions that remain if we are to begin automatic processing. A computational dialogue system must be able to identify SUs, match up their two (or more)

<sup>1</sup>Note that this says nothing about whether such a continuation is the same as the initial speaker’s intended continuation.

parts (which may not necessarily be adjacent), integrate them into some suitable syntactic and/or semantic representation, and determine the overall pragmatic contribution to the dialogue context. SUs also have implications for the organisation of *turn-taking* in such models (see e.g. Sacks et al. (1974)), as regards what conditions (if any) allow or prevent successful turn transfer. Additionally, from a socio-linguistic point of view, turn-taking operates (according to Schegloff (1995)) not on individual conversational participants, but on ‘parties’. Lerner (1991) suggests that split utterances can clarify the formation of such parties in that they reveal evidence of how syntax can be employed to organise participants into ‘groups’.

Analysis of SUs, when they can or cannot occur, and what effects they have on the coordination of agents in dialogue, is therefore an area of interest not only for conversational analysts wishing to characterise systematic interactions in dialogue, but also for linguists trying to formulate grammars of dialogue, psychologists and sociolinguists interested in alignment mechanisms and social interaction, and those interested in building automatic dialogue processing systems. In this paper we present and examine empirical corpus data in order to shed light on some of the questions and controversies around this phenomenon.

## 2 Related Work

Most previous work on what we call SUs has examined specific sub-cases, generally of the cross-speaker type, and have referred to these variously as *collaborative turn sequences* (Lerner, 1996; Lerner, 2004), *collaborative completions* (Clark, 1996; Poesio and Rieser, to appear), *co-constructions* (Sacks, 1992), *joint productions* (Helasvuo, 2004), *co-participant completions* (Hayashi, 1999; Lerner and Takagi, 1999), *collaborative productions* (Szczepek, 2000) and *anticipatory completions* (Fox and others, 2007) (amongst others). Here we discuss some of these views.

**Conversation Analysis** Lerner (1991) identifies various structures typical of SUs which contain characteristic split points. Firstly he gives a number of ‘compound’ *turn-constructive units* (TCUs), i.e., structures that include an initial constituent that hearers can identify as introducing some later final component. Examples include the IF X-THEN Y, WHEN X-THEN Y and INSTEAD

OF X-Y constructions:

(1) A: Before that then if they were ill

G: They get nothing. [BNC H5H 110-111]

Other cues for potential *anticipatory completions* include quotation markers (e.g. SHE SAID), parenthetical inserts and lists, as well as non-syntactic cues such as contrast stress or prefaced disagreements. Rühlemann (2007) uses corpus analysis to examine *sentence relatives* as typical expansions of another interlocutor's turn (see also (16)):

(2) A: profit for the group is a hundred and ninety thousand pounds.

B: Which is superb. [BNC FUK 2460-2461]

**Opportunistic Cases** Although Lerner focuses on these projectable turn completions, he also mentions that splits can occur at other points such as "intra-turn silence", hesitations etc. which he terms *opportunistic completions*:

(3) A: Well I do know last week that=uh Al was certainly very < pause 0.5)

B: pissed off [(Lerner, 1996, p260)]

As he makes no claims regarding the frequency of such devices for SUs, it would be interesting to know how common these are (insomuch as they occur at all and can be accordingly classified), especially as studies on SUs in Japanese (Hayashi, 1999) show that although SUs do occur, they do not rely on compound TCUs.

**Expansions vs. Completions** Other classifications of SUs often distinguish between *expansions* and *completions* (Ono and Thompson, 1993). Expansions are continuations which add, e.g., an adjunct, to an already complete syntactic element:

(4) T: It'll be an E sharp.

G: Which will of course just be played as an F. [BNC G3V 262-263]

whilst completions involve the addition of syntactic material which is required to make the whole utterance complete:

(5) A: ... and then we looked along one deck, we were high up, and down below there were rows of, rows of lifeboats in case you see

B: There was an accident.

A: of an accident [BNC HDK 63-65]

In terms of frequency, the only estimate we know of is Szczepek (2000), where there are apparently 200 cross-person SUs in 40 hours of English conversation (there is no mention of the number of sentences or turns this equates to), of which

75% are completions.<sup>2</sup> As briefly outlined above, CA analyses of SUs tend to be broadly descriptive of what they reveal for conversational practices. Because such analyses present real examples they establish that the phenomenon is a genuine one; however, there is no discussion of its scale (with the exception of Szczepek (2000), which offers extremely limited figures). Even though as a genuine phenomenon it is of theoretical interest, the lack of frequency statistics prevents generalisability. Therefore, any claims that SUs are pervasive in dialogue need empirical backing.

**Linguistic Models** Purver et al. (2006) present a grammatical model for split utterances, using an inherently incremental grammar formalism, Dynamic Syntax (Kempson et al., 2001; Cann et al., 2005). This model shows how syntactic and semantic processing can be accounted for no matter where the split occurs in a sentence; however, as their interest is in grammatical processing, they give no account of any higher-level inferences which may be required. Poesio and Rieser (to appear) present a general model for *collaborative completions* based in the PTT framework, using an incremental LTAG-based grammar and an information-state-based approach to context modelling. While many parts of their model are compatible with a simple alignment-based communication model like Pickering and Garrod (2004)'s, they see intention recognition as crucial to dialogue management. They conclude that an intention-based model, more like Clark (1996)'s, is more suitable. Their primary concern is to show how such a model can account for the hearer's ability to infer a suitable continuation, but their use of an incremental interpretation method also allows an explanation of the low-level utterance processing required. Nevertheless, the use of an essentially head-driven grammar formalism suggests that some syntactic splits that appear in our corpus might be more problematic than others.

**Corpus Studies** Skuplik (1999), as reported by Poesio and Rieser (to appear), collected data from German two-party task-oriented dialogue, and annotated for split utterance phenomena. She found that *expansions* (cases where the part before the split can be considered already complete) were

<sup>2</sup>However, this could be affected by her decision not to include what she calls *appendor questions* in her data which could also be argued to be expansion SUs.

more common than *completions* (where the first part is incomplete as it stands). Given that this study focuses on task-oriented dialogue, it needs to be shown that its results can be replicated in naturally occurring dialogue. In addition, de Ruiter and van Dinst (in preparation) are also in the process of studying other-initiated completions, in the above sense, and their effect on the progressivity of dialogue turns; however no results are available to us at this point in time.

**Dialogue Models** We are not aware of any system/model which treats other-person splits, but same-person ones are now being looked at. Skantze and Schlangen (2009) present an incremental system design (for a limited domain) which can react to user feedback, e.g., backchannels, and resume with utterance completion if interrupted. Some related empirical work regarding the issue of turn-switch addressed here is also presented in Schlangen (2006) but the emphasis there centered mostly on prosodic rather than grammar/theory-based factors.

### 3 Method

#### 3.1 Terminology

In this paper, as our interest is general, we use the term **split utterances** (*SUs*) to cover all instances where an utterance is spread across more than one dialogue contribution – whether the contributions are by the same or different speakers. We therefore use the term **split point** to refer to the point at which the utterance is split (rather than e.g. *transition point* which is associated with a speaker change). Cases where speaker does change across the split will be called **other-person** splits; otherwise **same-person** splits. One of the reasons for including same-person splits is that there are claims in the literature that the initial speaker may strategically continue completing their own utterance, after another person's intervention, as an alternative to acceptance or rejection of this intervention (*delayed completion*, (Lerner, 1996)). In addition, both grammatical formalisms (Purver et al., 2006) and psycholinguistic models (Pickering and Garrod, 2004) predict that SUs should be equally natural in both the same- and other- person conditions.

As not all cases will lead to complete contributions, and not all will be split over exactly two contributions, we also avoid terms like *first-half*,

*second-half* and *completion*: instead the contributions on either side of a split point will be referred to as the **antecedent** and the **continuation**. In cases where an utterance has more than one split point, some portions may therefore act as the continuation for one split point, and the antecedent for the next.

#### 3.2 Questions

**General** Our first interest is in the general statistics regarding SUs: how often do they occur, and what is the balance between same- and other-person splits? Do they usually fall into the specific categories (with specific preferred split points) examined by e.g. Lerner (1991), or can the split point be anywhere?

**Completeness** For a grammatical treatment of SUs, as well as for implementing parsing/production mechanisms for their processing, we need to know about the likely completeness of antecedent and continuation (if they are always complete in their own right, a standard head-driven grammar may be suitable; if not, something more fundamentally incremental may be required). In addition, CA and other strategic analyses of dialogue phenomena predict that split utterances should occur at turn-transfer points that are foreseeable by the participants. Complete syntactic units serve this purpose from this point of view and lack of such completeness will seem to weaken this general claim. We therefore ask how often antecedents and continuations are themselves complete,<sup>3</sup> and look at the syntactic and lexical categories which occur either side of the split.

**Repair and Overlap** Thirdly, we look at how often splits involve explicit repair of antecedent material, and how this depends on antecedent completeness. Although, sometimes, repair might be attributed to overlap or speaker uncertainty, it also might indicate issues regarding preemptive tactics on the part of the current speaker who needs to reformulate the original contribution in order to accommodate their novel offering or take into account feedback offered while constructing their utterance. Amount of repair also indicates the degree of attempt the current speaker is making to

<sup>3</sup>For antecedents, we are more interested in whether they *end* in a way that seems complete (they may have started irregularly due to overlap or another split); for continuations, whether they *start* in such a way (they may not get finished for some other reason, but we want to know if they would be complete if they do get finished).

Tag	Value	Explanation
end-complete	y/n	For all sentences: does this sentence end in such a way as to yield a complete proposition or speech act?
continues	sentence ID	For all sentences: does this sentence continue the proposition or speech act of a previous sentence? If so, which one?
repairs	number of words	For continuations: does this continuation explicitly repair words in the antecedent? If so, how many?
start-complete	y/n	For continuations: does this continuation start in such a way as to be able to stand alone as a complete proposition or speech act?

Table 1: Annotation Tags

integrate syntactically their contribution with the antecedent. However, we also examine how often continuations involve overlap, which also has implications for turn-taking management, and how this depends on antecedent completeness.

### 3.3 Corpus

For this exercise we used the portion of the BNC (Burnard, 2000) annotated by Fernández and Ginzburg (2002), chosen to maintain a balance between context-governed dialogue (tutorials, meetings, doctor’s appointments etc.) and general conversation. This portion comprises 11,469 sentences taken from 200-turn sections of 53 separate dialogues.

The BNC transcripts are already annotated for overlapping speech, for non-verbal noises (laughter, coughing etc.) and for significant pauses. Punctuation is included, based on the original audio and the transcribers’ judgements; as the audio is not available, we allowed annotators to use punctuation where it aided interpretation. The BNC transcription protocol provides a sentence-level annotation as well as an utterance (turn)-level one, where turns may be made of several sentences by the same speaker. We annotated at a sentence-level, to allow self-continuations within a turn to be examined. The BNC also forces turns to be presented in linear order, which is vital if we are to accurately assess whether turns are continuations of one another; however, this has a side-effect of forcing long turns to appear split into several shorter turns when interrupted by intervening backchannels. We will discuss this further below.

**Annotation Scheme** The initial stage of manual annotation involved 4 tags: `start-complete`, `end-complete`, `continues` and `repairs` – these are explained in Table 1 above. Sentences which somehow *require* continuation (whether

they receive it or not) are therefore those marked `end-complete=n`; sentences which act as continuations are those marked with non-empty `continues` tags; and their antecedents are the values of those `continues` tags. Further specific information about the syntactic or lexical nature of antecedent or continuation components could then be extracted (semi-)automatically, using the BNC transcript and part-of-speech annotations.

**Inter-Annotator Agreement** Three annotators were used, all linguistically knowledgeable. First, all three annotators annotated one dialogue independently, then compared results and discussed differences. They then annotated 3 further dialogues independently to assess inter-annotator agreement; kappa statistics (Carletta, 1996) are shown in Table 2 below.

Tag	KND	KBG	KB0
end-complete	.86-.92	.80-1.0	.73-.90
continues (y/n)	.89-.81	.76-.85	.77-.89
continues (ant)	.90-.82	.74-.85	.76-.86
repairs	1.0-1.0	.55-.81	1.0-1.0

Table 2: Inter-Annotator  $\kappa$  statistic (min-max)

With the exception of the `repairs` tag for one annotator pair for one dialogue, all are above 0.7; the low figure results from a few disagreements in a dialogue with only a very small number of `repairs` instances. The remaining dialogues were divided evenly between the three annotators.

## 4 Results and Discussion

The 11,469 sentences annotated yielded 2,228 SUs, of which 1,902 were same-person and 326 other-person splits; 111 examples involved an explicit repair by the continuation of some part of the antecedent.

person:	same	other
overlapping	0	17
adjacent	840	260
sep. by overlap	320	10
sep. by backchnl	460	17
sep. by 1 sent	239	16
sep. by 2 sents	31	4
sep. by 3 sents	5	1
sep. by 4 sents	4	0
sep. by 5 sents	1	0
sep. by 6 sents	2	1
Total	1902	326

Table 3: Antecedent/continuation separation

**General** Same-person splits are much more common than other-person; however, this is partly an artefact of the BNC transcription protocol (which forces contributions to be linearly ordered) and our choice to annotate at the sentence level. Around 44% of same-person cases are splits between sentences within the same-speaker turn; and a further 17% are separated only by other-speaker material which entirely overlaps with the antecedent and therefore does not necessarily actually interrupt the turn. Both of these might be considered as single utterances under some views. However, we believe that splits between same-turn sentences must be investigated in that the transcription into separate sentences does indicate some pause or other separating prosody and, from a processing/psycholinguistic point of view, it should be determined whether other-person splits occur in the same places as same-person split boundaries. Even in cases of overlap, one cannot exclude the fact that the shape of the current speaker’s utterance is influenced by receipt of the feedback. Nevertheless, we will examine these issues in further research and hence we exclude within-turn splits of this type from here on.

Many splits are non-adjacent (see Table 3), with the antecedent and continuation separated by at least one intervening sentence. In same-person cases, once we have excluded the within-turn splits described above, this must in fact always be the case; the intervening material is usually a backchannel (62% of remaining cases) or a single other sentence (32%, often e.g. a clarification question), but two intervening sentences are possible (4%) with up to six being seen. In other-person cases, 88% are adjacent or separated only by overlapping material, but again up to six intervening

person:	same	other
and/but/or	748	116
so/whereas	257	39
because	77	3
(pause)	56	5
which/who/etc	26	4
instead of	4	1
said/thought/etc	14	0
if_then	1	0
when_then	1	1
(other)	783	161

Table 4: Continuation categories

sentences were seen, with a single sentence most common (10%, in half of which the intervening sentence was a backchannel).

Many utterances have more than one split. In same-person cases, a single utterance can be split over as many as thirteen individual sentence contributions; although such extreme cases occur generally within one-sided dialogues such as tutorials, many multi-split cases are also seen in general conversation. Only 63% of cases consisted of only two contributions. Antecedents can also receive more than one competing continuation, although this is rare: two continuations are seen in 2% of cases.

**CA Categories** We searched for examples which match CA categories (Lerner, 1991; Rühlemann, 2007) by looking for particular lexical items on either side of the split. Matching was done loosely, to allow for the ungrammatical nature of dialogue – for example, an instance was taken to match the IF X-THEN Y pattern if the continuation began with ‘then’ (modulo filled pauses and non-verbal material) and the antecedent contained ‘if’ at any point) – so the counts may be over-estimates. For Lerner (1996)’s *opportunistic* cases, we looked for filled pauses (‘er/erm’ etc.) or pauses explicitly annotated in the transcript, so counts in this case may be underestimates.<sup>4</sup> We also chose some other broad categories based on our observations of the most common cases. Results are shown in Table 4.<sup>5</sup>

The most common of the CA categories can be

<sup>4</sup>In further research we will examine other features as specialised laugh tokens, repetitions etc. as well as their particular positioning

<sup>5</sup>Note that the categories in Table 4 are not all mutually exclusive (e.g. an example may have both an ‘and’-initial continuation and an antecedent ending in a pause), so column sums will not match Table 3.

seen to be Lerner (1996)'s hesitation-related *opportunistic* cases, which make up at least 2-3% of both same- and other-person splits. Rühlemann (2007)'s *sentence relative* clause cases are next, with over 1%; the others make up only small proportions.

In contrast, by far the most common pattern (for both same- and other-) is the addition of an extending clause, either a conjunction introduced by 'and/but/or/nor' (35-40%), or other clause types with 'so/whereas/nevertheless/because'. Other less obviously categorisable cases make up 40-50% of continuations, with the most common first words being 'you', 'it', 'I', 'the', 'in' and 'that'.

**Completeness and repair** Examination of the end-complete annotations shows that about 8% of sentences in general are incomplete, but that (perhaps surprisingly) only 63% of these get continued. For both same- and other-person continuations, the vast majority (72% and 74%) continue an already complete antecedent, with only 26-28% therefore being *completions* in the sense of e.g. de Ruiter and van Dienst (in preparation). This does, however, mean that continuations are significantly more likely than other sentences to follow an incomplete antecedent ( $p < 0.001$  using  $\chi^2_{(1)}$ ). Interestingly, though, continuations are no more likely than other sentences to be complete themselves.

The frequent clausal categories from Table 4 are all more likely to continue complete antecedents than incomplete ones, with the exception of the (other) category; this suggests that split points often occur at random points in a sentence, without regard to particular clausal constructions (see also A.1 for more examples and context):

- (6) D: you know what the actual variations  
 U: entails  
 D: entails. you know what the actual quality of the variations are.  
 [BNC G4V 114-117]

For the less frequent (e.g. 'if/then', 'instead of') categories, the counts are too low to be sure.

Excluding all the clausal constructions (i.e. looking only at the general (other) category), and looking only at other-person cases, we see that antecedents often end in a complete way (53%) but that continuations do not often start in a complete way (24%). Continuations are more than twice as likely to start in a non-complete as opposed

to complete way, even after complete antecedents. Explicit repair of some portion of the antecedent is not common, only occurring in just under 5% of splits. As might be expected, incomplete antecedents are more likely to be repaired (13% vs. 2%,  $p < 0.001$  using  $\chi^2_{(1)}$ ). Other-continuations are also significantly more likely to repair their antecedents than same-person cases (10% vs. 4%,  $p < 0.001$  using  $\chi^2_{(1)}$ ).

**Problematic cases** Examination of the data shows that SUs is not necessarily an autonomous well-defined category independent of other fragment classifications in the literature. Besides cases where it is not easy to identify whether a fragment is a continuation or not or what the antecedent is (see A.2), there are also cases where, as has already been pointed out in the literature (Gregoromichelaki et al., 2009; Bunt, 2009), fragments exhibit multifunctionality. This can be illustrated by the following where the continuation could be taken also as request for confirmation/question (7) or a reply to a clarification request (8):

- (7) M: It's generated with a handle and  
 J: Wound round?  
 M: Yes [BNC K69 109-112]
- (8) S: Quite a good word processor.  
 J: A word processor?  
 S: Which is vag- it's basically a subset of Word. [BNC H61 37-39]

In this respect, an interesting category is Lerner's *delayed completions* where often the continuation also serves as some kind of repair or reformulation (see e.g. (6) and A.3 (26)).

## 5 Conclusions

Although most of Lerner (1991)'s categories appear, they are not necessarily the most frequent. On the other hand, the general results seem to indicate that splits can occur anywhere in a string, both in the same- or other- conditions. Both these are consistent with models that advocate highly coordinated resources between interlocutors and, moreover, the need for highly incremental means of processing (Purver et al., 2006; Skantze and Schlangen, 2009). From a computational modelling point of view, the results also indicate that start-completeness of continuations is rare, which means that a dialogue system has a chance of spotting continuations from surface characteristics of

the input. This is hampered though by the fact that the split can occur within any type of syntactic constituent, hence no reliable grammatical features can be employed securely. On the other hand, end-incompleteness of antecedents is not as common as would be expected and long distances between antecedent and continuation are possible. In this respect, locating the antecedent is not a straightforward task for automated systems, especially again as this can be any type of constituent.

## References

- H. Bunt. 2009. Multifunctionality and multidimensional dialogue semantics. In *Proceedings of Dia-Holmia, 13th SEMDIAL Workshop*.
- L. Burnard. 2000. *Reference Guide for the British National Corpus (World Edition)*. Oxford University Computing Services <http://www.natcorp.ox.ac.uk/docs/userManual/>.
- R. Cann, R. Kempson, and L. Marten. 2005. *The Dynamics of Language*. Elsevier, Oxford.
- J. Carletta. 1996. Assessing agreement on classification tasks: The kappa statistic. *Computational Linguistics*, 22(2):249–255.
- H. Clark. 1996. *Using Language*. Cambridge University Press.
- J. de Ruiter and M. van Dinst. in preparation. Completing other people's utterances: evidence for forward modeling in conversation. ms.
- R. Fernández and J. Ginzburg. 2002. Non-sentential utterances: A corpus-based study. *Traitement Automatique des Langues*, 43(2).
- A. Fox et al. 2007. Principles shaping grammatical practices: an exploration. *Discourse Studies*, 9(3):299.
- E. Gregoromichelaki, Y. Sato, R. Kempson, A. Gargett, and C. Howes. 2009. Dialogue modelling and the remit of core grammar. In *Proceedings of IWCS*.
- M. Hayashi. 1999. Where Grammar and Interaction Meet: A Study of Co-Participant Completion in Japanese Conversation. *Human Studies*, 22(2):475–499.
- M. Helasvuo. 2004. Shared syntax: the grammar of co-constructions. *Journal of Pragmatics*, 36(8):1315–1336.
- R. Kempson, W. Meyer-Viol, and D. Gabbay. 2001. *Dynamic Syntax: The Flow of Language Understanding*. Blackwell.
- G. Lerner and T. Takagi. 1999. On the place of linguistic resources in the organization of talk-in-interaction: A co-investigation of English and Japanese grammatical practices. *Journal of Pragmatics*, 31(1):49–75.
- G. Lerner. 1991. On the syntax of sentences-in-progress. *Language in Society*, pages 441–458.
- G. Lerner. 1996. On the semi-permeable character of grammatical units in conversation: Conditional entry into the turn space of another speaker. In E. Ochs, E. A. Schegloff, and S. A. Thompson, editors, *Interaction and grammar*, pages 238–276. Cambridge University Press.
- G. Lerner. 2004. Collaborative turn sequences. In *Conversation analysis: Studies from the first generation*, pages 225–256. John Benjamins.
- T. Ono and S. Thompson. 1993. What can conversation tell us about syntax. In P. Davis, editor, *Alternative Linguistics: Descriptive and Theoretical Modes*. Benjamin.
- M. Pickering and S. Garrod. 2004. Toward a mechanistic psychology of dialogue. *Behavioral and Brain Sciences*, 27:169–226.
- M. Poesio and H. Rieser. to appear. Completions, coordination, and alignment in dialogue. Ms.
- M. Purver, R. Cann, and R. Kempson. 2006. Grammars as parsers: Meeting the dialogue challenge. *Research on Language and Computation*, 4(2-3):289–326.
- C. Rühlemann. 2007. *Conversation in context: a corpus-driven approach*. Continuum.
- H. Sacks, E. A. Schegloff, and G. Jefferson. 1974. A simplest systematics for the organization of turn-taking for conversation. *Language*, 50(4):696–735.
- H. Sacks. 1992. *Lectures on Conversation*. Blackwell.
- E. Schegloff. 1995. Parties and talking together: Two ways in which numbers are significant for talk-in-interaction. *Situated order: Studies in the social organization of talk and embodied activities*, pages 31–42.
- D. Schlangen. 2006. From reaction to prediction: Experiments with computational models of turn-taking. In *Proceedings of the 9th International Conference on Spoken Language Processing (INTERSPEECH - ICSLP)*.
- G. Skantze and D. Schlangen. 2009. Incremental dialogue processing in a micro-domain. In *Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009)*.
- K. Skuplik. 1999. Satzkooperationen. definition und empirische untersuchung. SFB 360 1999/03, Bielefeld University.
- B. Szczepek. 2000. Formal Aspects of Collaborative Productions in English Conversation. *Interaction and Linguistic Structures (InLiSt)*, <http://www.uni-potsdam.de/u/inlist/issues/17/>.



## A Examples

### A.1 Split points

- (6) D: Yeah I mean if you're looking at quantitative things it's really you know how much actual- How much variation happens whereas qualitative is ⟨pause⟩ you know what the actual variations

U: entails

D: entails. you know what the actual quality of the variations are.

[BNC G4V 114-117]

- (9) A: All the machinery was

G: [[All steam.]]<sup>6</sup>

A: [[operated]] by steam

[BNC H5G 177-179]

- (10) K: I've got a scribble behind it, oh annual report I'd get that from.

S: Right.

K: And the total number of [[sixth form students in a division.]]

S: [[Sixth form students in a division.]] Right.

[BNC H5D 123-127]

- (11) M: 292 And another sixteen percent is the other Ne- Nestle coffee ⟨pause⟩ erm Blend Thirty Seven which I used to drink a long time ago and others ⟨laugh⟩ and twenty two percent is er ⟨pause⟩

U: Maxwell.

M: Maxwell House, which has become the other local brand now seeing as how Maxwell House is owned by Kraft, and Kraft now own Terry's.

[BNC G3U 292-294]

- (12) A: Erm because as Moira said that Kraft is erm ⟨pause⟩ now what was she saying, what was she saying Kraft is the same as ⟨pause⟩

M: Craft? [BNC G3U 412-413]

- (13) J: And I couldn't remember whether she said at the end of the three months or

A: End of the month. [BNC H4P 17-18]

- (14) G: Had their own men

A: unload the boats?

G: unload the boats, yes. [BNC H5H 91-93]

- (15) G: That's right they had to go on a rota.

A: Run by the Dock Commission?

G: Run by the Dock Commission.

[BNC H5H 100-102]

- (16) A: So I thought, oh, I think I'll put lace over it, it'll tone the lilac [[down.]]

B: [[down.]] Yes.  
Which it is has done

[BNC KBC 3195-3198]

### A.2 Uncertain antecedents

- (17) C: Look you're cleaning this ⟨pause⟩ [[with erm]]

G: [[That box.]]

C: [[This.]]

G: [[With]] this. [[And this.]]

C: [[And this.]] [[And this.]]

G: [[And this.]]

Whoops! [BNC KSR 9-17]

- (18) S: You're trying to be everything ⟨pause⟩ and they're pushing it away cos it's not what they really want ⟨pause⟩ and they, I mean, all, all you can get from him is how marvellous, you're right, how marvellous his brothers are ⟨pause⟩ and yet, what I've heard of the brothers they're not

C: Not much, [[yeah.]]

S: [[they're]] not all that marvellous, they're not really that much to look [[up]]

C: [[Ah]].

S: to.

C: No [BNC KBG 76-81]

- (19) S: Well this is why I think he'd be better off, hi- his needs ⟨pause⟩ are not met by a class teacher. And I don't think they have been for this last

C: Mm, we need a support teacher [[to go there.]]

S: [[for the last]] year. But yo-, you need somebody who's gonna work with him every day ⟨pause⟩ and ⟨pause⟩ with an individual programme and you just can't offer that ⟨pause⟩ in a class. [BNC KBG 56-60]

<sup>6</sup>Overlapping material is shown in double square brackets, aligned with the material with which it co-occurs.

(20) M: I might be a bit biased, I think they still do that but I think erm <pause>

J: The television has <pause>

M: the television has made a difference. I think not only just at fire stations, I think in the whole of life, hasn't it?

[BNC K69 51-54]

(21) A5: I'll definitely use that

U: <reading>:[ Get a headache ]?

A5: [[in getting to know ]]

A2: [[Year seven ]]

A5: new [[year seven]]

A2: [[Oh yeah]] for year seven

[BNC J8D 190-195]

(22) G: Well a chain locker is where all the spare chain used to like coil up

A: So it <unclear> came in and it went round

G: round the barrel about three times round the barrel then right down into the chain locker but if you kept, let it ride what we used to call let it ride well <unclear> well now it get so big then you have to run it all off cos you had one lever, that's what you had and the steam valve could have all steamed.

[BNC H5G 174:176]

### A.3 Multifunctionality of fragments

(7) *Completion and confirmation request:*

J: How does it generate?

M: It's generated with a handle and

J: Wound round?

M: Yes, wind them round and this should, should generate a charge which rang bells and sounded bells and then er you lift up a telephone and plug in a jack and, and take a message in that way.

[BNC K69 109-112]

(23) *Completion and confirmation request:*

G: Had their own men

A: unload the boats?

G: unload the boats, yes. [BNC H5H 91-93]

(24) *Late completion and (repetitive) confirmation:*

N: Alistair [last or full name] erm he's, he's made himself er he has made himself co-ordinator.

U: And section engineer.

N: And section engineer.

N: I didn't sign it as coordinator.

[BNC H48 141-144]

(25) *Completion and clarification reply:*

John: If you press N

Sarah: N?

John: N for name, it'll let you type in the docu document name. [BNC G4K 84-86]

(26) *Expansion and reformulation/repair:*

S: Secondly er

J: We guarantee P five.

S: We we are we're guaranteeing P five plus a noise level.

J: Yeah. [BNC JP3 167-170]

(27) *Expansion and question:*

I: I can't remember exactly who lived on the right hand side, I've forgotten but th I know the Chief Clerk lived just a little way down [address], you see, er

A: In one of those little red brick cottages?

[BNC HDK 124-125]

(28) *Answer and expansion:*

A: We could hear it from outside <unclear>.

R: Oh you could hear it?

A: Occasionally yeah. [BNC J8D 13-15]

(29) *Answer/reformulation and expansion:*

G: [address], that was in the middle, more or less in the middle of the town.

A: And you called that the manual?

G: The manual school, yes.

[BNC H5G 96-98]