



# Attention, Intentions, and the Structure of Discourse

## Citation

Grosz, Barbara J. and Candace L. Sidner. 1986. Attention, intentions, and the structure of discourse. *Computational Linguistics* 12(3): 175-204.

## Published Version

<http://www.aclweb.org/anthology-new/J/J86/J86-3001.pdf>

## Permanent link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:2579648>

## Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA>

## Share Your Story

The Harvard community has made this article openly available.  
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

# ATTENTION, INTENTIONS, AND THE STRUCTURE OF DISCOURSE

Barbara J. Grosz

Artificial Intelligence Center *and*  
Center for the Study of Language and Information  
SRI International  
Menlo Park, CA 94025

Candace L. Sidner

BBN Laboratories Inc.  
Cambridge, MA 02238

In this paper we explore a new theory of discourse structure that stresses the role of purpose and processing in discourse. In this theory, discourse structure is composed of three separate but interrelated components: the structure of the sequence of utterances (called the linguistic structure), a structure of purposes (called the intentional structure), and the state of focus of attention (called the attentional state). The linguistic structure consists of segments of the discourse into which the utterances naturally aggregate. The intentional structure captures the discourse-relevant purposes, expressed in each of the linguistic segments as well as relationships among them. The attentional state is an abstraction of the focus of attention of the participants as the discourse unfolds. The attentional state, being dynamic, records the objects, properties, and relations that are salient at each point of the discourse. The distinction among these components is essential to provide an adequate explanation of such discourse phenomena as cue phrases, referring expressions, and interruptions.

The theory of attention, intention, and aggregation of utterances is illustrated in the paper with a number of example discourses. Various properties of discourse are described, and explanations for the behavior of cue phrases, referring expressions, and interruptions are explored.

This theory provides a framework for describing the processing of utterances in a discourse. Discourse processing requires recognizing how the utterances of the discourse aggregate into segments, recognizing the intentions expressed in the discourse and the relationships among intentions, and tracking the discourse through the operation of the mechanisms associated with attentional state. This processing description specifies in these recognition tasks the role of information from the discourse and from the participants' knowledge of the domain.

## 1 INTRODUCTION

This paper presents the basic elements of a computational theory of discourse structure that simplifies and expands upon previous work. By specifying the basic units a discourse comprises and the ways in which they can relate, a proper account of discourse structure provides the basis for an account of discourse meaning. An account of discourse structure also plays a central role in language processing because it stipulates constraints on those portions of a discourse to which any given utterance in the discourse must be related.

An account of discourse structure is closely related to two questions: What individuates a discourse? What makes it coherent? That is, faced with a sequence of utterances, how does one know whether they constitute a single discourse, several (perhaps interleaved) discourses, or none? As we develop it, the theory of discourse structure will be seen to be intimately connected with two nonlinguistic notions: intention and attention. Attention is an essential factor in explicating the processing of utterances in discourse. Intentions play a primary role in explaining discourse structure, defining discourse coherence, and providing a coherent conceptualization of the term "discourse" itself.

Copyright 1986 by the Association for Computational Linguistics. Permission to copy without fee all or part of this material is granted provided that the copies are not made for direct commercial advantage and the *CL* reference and this copyright notice are included on the first page. To copy otherwise, or to republish, requires a fee and/or specific permission.

0362-613X/86/030175-204\$03.00

The theory is a further development and integration of two lines of research: work on focusing in discourse (Grosz 1978a, 1978b, 1981) and more recent work on intention recognition in discourse (Sidner and Israel 1981; Sidner 1983; 1985; Allen 1983, Litman 1985; Pollack 1986). Our goal has been to generalize these constructs properly to a wide range of discourse types. Grosz (1978a) demonstrated that the notions of focusing and task structure are necessary for understanding and producing task-oriented dialogue. One of the main generalizations of previous work will be to show that discourses are generally in some sense "task-oriented," but the kinds of "tasks" that can be engaged in are quite varied – some are physical, some mental, others linguistic. Consequently, the term "task" is misleading; we therefore will use the more general terminology of *intentions* (e.g., when speaking of discourse purposes) for most of what we say.

Our main thesis is that the structure of any discourse is a composite of three distinct but interacting components:

- the structure of the actual sequence of utterances in the discourse;
- a structure of intentions;
- an attentional state.

The distinction among these components is essential to an explanation of interruptions (see Section 5), as well as to explanations of the use of certain types of referring expressions (see Section 4.2) and various other expressions that affect discourse segmentation and structure (see Section 6). Most related work on discourse structure (including Reichman-Adar 1984, Linde 1979, Linde and Goguen 1978, Cohen 1983) fails to distinguish among some (or, in some cases, all) of these components. As a result, significant generalizations are lost, and the computational mechanisms proposed are more complex than necessary. By carefully distinguishing these components, we are able to account for significant observations in this related work while simplifying both the explanations given and computational mechanisms used.

In addition to explicating these linguistic phenomena, the theory provides an overall framework within which to answer questions about the relevance of various segments of discourse to one another and to the overall purposes of the discourse participants. Various properties of the intentional component have implications for research in natural-language processing in general. In particular, the intentions that underlie discourse are so diverse that approaches to discourse coherence based on selecting discourse relationships from a fixed set of alternative rhetorical patterns (e.g., Hobbs 1979, Mann and Thompson 1983, Reichman 1981) are unlikely to suffice. The intentional structure introduced in this paper depends instead on a small number of structural relations that can hold between intentions. This study also reveals several problems that must be confronted in expanding speech-act-related theories (e.g., Allen and Perrault 1980, Cohen and Levesque 1980, Allen 1983) from

coverage of individual utterances to coverage of extended sequences of utterances in discourse.

Although a definition of discourse must await further development of the theory presented in this paper, some properties of the phenomena we want to explain must be specified now. In particular, we take a *discourse* to be a piece of language behavior that typically involves multiple utterances and multiple participants. A discourse may be produced by one or more of these participants as speakers or writers; the audience may comprise one or more of the participants as hearers or readers. Because in multi-party conversations more than one participant may speak (or write) different utterances within a segment, the terms *speaker* and *hearer* do not differentiate the unique roles that the participants maintain in a segment of a conversation. We will therefore use the terms *initiating conversational participant* (ICP) and *other conversational participant(s)* (OCP) to distinguish the initiator of a discourse segment from its other participants. The ICP speaks (or writes) the first utterance of a segment, but an OCP may be the speaker of some subsequent utterances. By speaking of ICPs and OCPs, we can highlight the purposive aspect of discourse. We will use the terms *speaker* and *hearer* only when the particular speaking/hearing activity is important for the point being made.

In most of this paper, we will be concerned with developing an abstract model of discourse structure; in particular, the definitions of the components will abstract away from the details of the discourse participants. Whether one constructs a computer system that can participate in a discourse (i.e., one that is a language user) or defines a psychological theory of language use, the task will require the appropriate projection of this abstract model onto properties of a language user, and specification of additional details (e.g., specifying memory for linguistic structure, means for encoding attentional state, and appropriate representations of intentional structure). We do, however, address ourselves directly to certain processing issues that are essential to the computational validity of the [abstract] model and to its utilization for a language-processing system or psychological theory.

Finally, it is important to note that although discourse *meaning* is a significant, unsolved problem, we will not address it in this paper. An adequate theory of discourse meaning needs to rest at least partially on an adequate theory of discourse structure. Our concern is with providing the latter.

The next section examines the basic theory of discourse structure and presents an overview of each of the components of discourse structure. Section 3 analyzes two sample discourses – a written text and a fragment of task-oriented dialogue – from the perspective of the theory being developed; these two examples are also used to illustrate various points in the remainder of the paper. Section 4 investigates various processing

issues that the theory raises. The following two sections describe the role of the discourse structure components in explaining various properties of discourse, thereby corroborating the necessity of distinguishing among its three components. Section 7 describes the generalization from utterance-level to discourse-level intentions, establishes certain properties of the latter, and contrasts them with the rhetorical relations of alternative theories. Finally, Section 8 poses a number of outstanding research questions suggested by the theory.

## 2 THE BASIC THEORY

Discourse structure is a composite of three interacting constituents: a linguistic structure, an intentional structure, and an attentional state. These three constituents of discourse structure deal with different aspects of the utterances in a discourse. Utterances – the actual saying or writing of particular sequences of phrases and clauses – are the linguistic structure's basic elements. Intentions of a particular sort and a small number of relationships between them provide the basic elements of the intentional structure. Attentional state contains information about the objects, properties, relations, and discourse intentions that are most salient at any given point. It is an abstraction of the focus of attention of the discourse participants; it serves to summarize information from previous utterances crucial for processing subsequent ones, thus obviating the need for keeping a complete history of the discourse.

Together the three constituents of discourse structure supply the information needed by the CPs to determine how an individual utterance fits with the rest of the discourse – in essence, enabling them to figure out why it was said and what it means. The context provided by these constituents also forms the basis for certain expectations about what is to come; these expectations play a role in accommodating new utterances. The attentional state serves an additional purpose: namely, it furnishes the means for actually using the information in the other two structures in generating and interpreting individual utterances.

### 2.1 LINGUISTIC STRUCTURE

The first component of discourse structure is the structure of the sequence of utterances that comprise a discourse.<sup>1</sup> Just as the words in a single sentence form constituent phrases, the utterances in a discourse are naturally aggregated into **discourse segments**. The utterances in a segment, like the words in a phrase, serve particular roles with respect to that segment. In addition, the discourse segments, like the phrases, fulfill certain functions with respect to the overall discourse. Although two consecutive utterances may be in the same discourse segment, it is also common for two consecutive utterances to be in different segments. It is also possible for two utterances that are nonconsecutive to be in the same segment.

The factoring of discourses into segments has been observed across a wide range of discourse types. Grosz (1978a) showed this for task-oriented dialogues. Linde (1979) found it valid for descriptions of apartments; Linde and Goguen (1978) describe such structuring in the Watergate transcripts. Reichman-Adar (1984) observed it in informal debates, explanations, and therapeutic discourse. Cohen (1983) found similar structures in essays in rhetorical texts. Polanyi and Scha (1986) discuss this feature of narratives.

Although different researchers with different theories have examined a variety of discourse types and found discourse-level segmentation, there has been very little investigation of the extent of agreement about where the segment boundaries lie. There have been no psychological studies of the consistency of recognition of section boundaries. However, Mann (Mann et al. 1975) asked several people to segment a set of dialogues. He has reported [personal communication] that his subjects segmented the discourses approximately the same; their disagreements were about utterances at the boundaries of segments.<sup>2</sup> Several studies of spontaneously produced discourses provide additional evidence of the existence of segment boundaries, as well as suggesting some of the linguistic cues available for detecting boundaries. Chafe (1979, 1980) found differences in pause lengths at segment boundaries. Butterworth (1975) found speech rate differences that correlated with segments; speech rate is slower at start of a segment than toward the end.

The linguistic structure consists of the discourse segments and an embedding relationship that can hold between them. As we discuss in Sections 2.2 and 5, the embedding relationships are a surface reflection of relationships among elements of the intentional structure. It is important to recognize that the linguistic structure is not strictly decompositional. An individual segment may include a combination of subsegments and utterances only in that segment (and not members of any of its embedded subsegments). Both of the examples in Section 3 exhibit such nonstrict decompositionality. Because the linguistic structure is not strictly decompositional, various properties of the discourse (most notably the intentional structure) are functions of properties of individual utterances and properties of segments.

There is a two-way interaction between the discourse segment structure and the utterances constituting the discourse: linguistic expressions can be used to convey information about the discourse structure; conversely, the discourse structure constrains the interpretation of expressions (and hence affects what a speaker says and how a hearer will interpret what is said). Not surprisingly, linguistic expressions are among the primary indicators of discourse segment boundaries. The explicit use of certain words and phrases (e.g., *in the first place*) and more subtle cues, such as intonation or changes in tense and aspect, are included in the repertoire of linguistic devices that function, wholly or in part, to indicate these

boundaries (Grosz 1978a, Reichman-Adar 1984, Cohen 1983, Polanyi and Scha 1983, Hirschberg and Pierrehumbert 1986). Reichman (1981) discusses some words that function in this way and coined the term **clue words**. We will use the term **cue phrases** to generalize on her observation as well as many others because each one of these devices cue the hearer to some change in the discourse structure.

As discussed in Section 6, these linguistic boundary markers can be divided according to whether they explicitly indicate changes in the intentional structure or in the attentional state of the discourse. The differential use of these linguistic markers provides one piece of evidence for considering these two components to be distinct. Because these linguistic devices function explicitly as indicators of discourse structure, it becomes clear that they are best seen as providing information at the discourse level, and not at the sentence level; hence, certain kinds of questions (e.g., about their contribution to the truth conditions of an individual sentence) do not make sense. For example, in the utterance *Incidentally, Jane swims every day*, the *incidentally* indicates an interruption of the main flow of discourse rather than affecting in any way the meaning of *Jane swims every day*. Jane's swimming every day could hardly be fortuitous.

Just as linguistic devices affect structure, so the discourse segmentation affects the interpretation of linguistic expressions in a discourse. Referring expressions provide the primary example of this effect.<sup>3</sup> The segmentation of discourse constrains the use of referring expressions by delineating certain points at which there is a significant change in what entities (objects, properties, or relations) are being discussed. For example, there are different constraints on the use of pronouns and reduced definite-noun phrases within a segment than across segment boundaries. While discourse segmentation is obviously not the only factor governing the use of referring expressions, it is an important one.

## 2.2 INTENTIONAL STRUCTURE

A rather straightforward property of discourses, namely, that they (or, more accurately, those who participate in them) have an overall purpose, turns out to play a fundamental role in the theory of discourse structure. In particular, some of the purposes that underlie discourses, and their component segments, provide the means of individuating discourses and of distinguishing discourses that are coherent from those that are not. These purposes also make it possible to determine when a sequence of utterances comprises more than one discourse.

Although typically the participants in a discourse may have more than one aim in participating in the discourse (e.g., a story may entertain its listeners as well as describe an event; an argument may establish a person's brilliance as well as convince someone that a claim or

allegation is true), we distinguish one of these purposes as foundational to the discourse. We will refer to it as the **discourse purpose** (DP). From an intuitive perspective, the discourse purpose is the intention that underlies engaging in the particular discourse. This intention provides both the reason a discourse (a linguistic act), rather than some other action, is being performed and the reason the particular content of this discourse is being conveyed rather than some other information. For each of the discourse segments, we can also single out one intention – the **discourse segment purpose** (DSP). From an intuitive standpoint, the DSP specifies how this segment contributes to achieving the overall discourse purpose. The assumption that there are single such intentions will in the end prove too strong. However, this assumption allows us to describe the basic theory more clearly. We must leave to future research (and a subsequent paper) the exploration and discussion of the complications that result from relaxing this assumption.

Typically, an ICP will have a number of different kinds of intentions that lead to initiating a discourse. One kind might include intentions to speak in a certain language or to utter certain words. Another might include intentions to amuse or to impress. The kinds of intentions that can serve as discourse purposes or discourse segment purposes are distinguished from other intentions by the fact that they are intended to be recognized (cf. Allen and Perrault 1980, Sidner 1985), whereas other intentions are private; that is, the recognition of the DP or DSP is *essential* to its achieving its intended effect. Discourse purposes and discourse segment purposes share this property with certain utterance-level intentions that Grice (1969) uses in defining utterance meaning (see Section 7).

It is important to distinguish intentions that are intended to be recognized from other kinds of intentions that are associated with discourse. Intentions that are intended to be recognized achieve their intended effect only if the intention is recognized. For example, a compliment achieves its intended effect only if the intention to compliment is recognized; in contrast, a scream of *boo* typically achieves its intended effect (scaring the hearer) without the hearer having to recognize the speaker's intention.

Some intention that is private and not intended to be recognized may be the primary motivation for an ICP to begin a discourse. For example, the ICP may intend to impress someone or may plan to teach someone. In neither case is the ICP's intention necessarily intended to be recognized. Quite the opposite may be true in the case of impressing, as the ICP may not want the OCP to be aware of his intention. When teaching, the ICP may not care whether the OCP knows the ICP is teaching him or her. Thus, the intention that motivates the ICP to engage in a discourse may be private. By contrast, the discourse segment purpose is always intended to be recognized.

DPs and DSPs are basically the same sorts of intentions. If an intention is a DP, then its satisfaction is a main purpose of the discourse, whereas if it is a DSP, then its satisfaction contributes to the satisfaction of the DP. The following are some of the types of intentions that could serve as DP/DSPs, followed by one example of each type.

1. Intend that some agent intend to perform some physical task. Example: *Intend that Ruth intend to fix the flat tire.*
2. Intend that some agent believe some fact. Example: *Intend that Ruth believe the campfire has started.*
3. Intend that some agent believe that one fact supports another. Example: *Intend that Ruth believe the smell of smoke provides evidence that the campfire is started.*
4. Intend that some agent intend to identify an object (existing physical object, imaginary object, plan, event, event sequence). Example: *Intend that Ruth intend to identify my bicycle.*
5. Intend that some agent know some property of an object. Example: *Intend that Ruth know that my bicycle has a flat tire.*

We have identified two structural relations that play an important role in discourse structure: **dominance** and **satisfaction-precedence**. An action that satisfies one intention, say DSP1, may be intended to provide part of the satisfaction of another, say DSP2. When this is the case, we will say that DSP1 **contributes to** DSP2; conversely, we will say that DSP2 **dominates** DSP1 (or DSP2 *DOM* DSP1). The dominance relation invokes a partial ordering on DSPs that we will refer to as the **dominance hierarchy**. For some discourses, including task-oriented ones, the order in which the DSPs are satisfied may be significant, as well as being intended to be recognized. We will say that DSP1 **satisfaction-precedes** DSP2 (or, DSP1 *SP* DSP2) whenever DSP1 must be satisfied before DSP2.<sup>4</sup>

Any of the intentions on the preceding list could be either a DP or a DSP. Furthermore, a given instance of any one of them could contribute to another, or to a different, instance of the same type. For example, the intention that someone intend to identify some object might dominate several intentions that she or he know some property of that object; likewise, the intention to get someone to believe some fact might dominate a number of contributing intentions that that person believe other facts.

As the above list makes clear, the range of intentions that can serve as discourse, or discourse segment, purposes is open-ended (cf. Wittgenstein 1953: paragraph 23), much like the range of intentions that underlie more general purposeful action. There is no finite list of discourse purposes, as there is, say, of syntactic categories. It remains an unresolved research question whether there is a finite description of the open-ended set of such intentions. However, even if there were finite descriptions, there would still be no finite list of

intentions from which to choose. Thus, a theory of discourse structure cannot depend on choosing the DP/DSPs from a fixed list (cf. Reichman-Adar 1984, Schank et al. 1982, Mann and Thompson 1983), nor on the particulars of individual intentions. Although the particulars of individual intentions, like a wide range of common sense knowledge, are crucial to understanding any discourse, such particulars cannot serve as the basis for *determining* discourse structure.

What is essential for discourse structure is that such intentions bear certain kinds of structural relationships to one another. Since the CPs can never know the whole set of intentions that might serve as DP/DSPs, what they must recognize is the relevant structural relationships among intentions. Although there is an infinite number of intentions, there are only a small number of relations relevant to discourse structure that can hold between them.

In this paper we distinguish between the *determination* of the DSP and the *recognition* of it. We use the term **determination** to refer to a semantic-like notion, namely, the complete specification of what is intended by whom; we use the term **recognition** to refer to a processing notion, namely, the processing that leads a discourse participant to identify what the intention is. These are obviously related concepts; the same information that determines a DSP may be used by an OCP to recognize it. However, some questions are relevant to only one of them. For example, the question of when the information becomes available is not relevant to determination but is crucial to recognition. An analogous distinction has been drawn with respect to sentence structure; the parse tree (determination) is differentiated from the parsing process (recognition) that produces the tree.

### 2.3 ATTENTIONAL STATE

The third component of discourse structure, the attentional state, is an abstraction of the participants' focus of attention as their discourse unfolds. The attentional state is a property of the discourse itself, not of the discourse participants. It is inherently dynamic, recording the objects, properties, and relations that are salient at each point in the discourse. The attentional state is modeled by a set of **focus spaces**; changes in attentional state are modeled by a set of transition rules that specify the conditions for adding and deleting spaces. We call the collection of focus spaces available at any one time the **focusing structure** and the process of manipulating spaces **focusing**.

The focusing process associates a focus space with each discourse segment; this space contains those entities that are salient – either because they have been mentioned explicitly in the segment or because they became salient in the process of producing or comprehending the utterances in the segment (as in the original work on focusing: Grosz 1978a). The focus space also includes the DSP; the inclusion of the purpose reflects the

fact that the CPs are focused not only on what they are talking about, but also on why they are talking about it.

To understand the attentional state component of discourse structure, it is important not to confuse it with two other concepts. First, the attentional state component is not equivalent to cognitive state, but is only one of its components. Cognitive state is a richer structure, one that includes at least the knowledge, beliefs, desires, and intentions of an agent, as well as the cognitive correlates of the attentional state as modeled in this paper. Second, although each focus space contains a DSP, the focus structure does *not* include the intentional structure as a whole.

Figure 1 illustrates how the focusing structure, in addition to modeling attentional state, serves during processing to coordinate the linguistic and intentional structures. The discourse segments (to the left of the figure) are tied to focus spaces (drawn vertically down the middle of the figure). The focusing structure is a stack. Information in lower spaces is usually accessible from higher ones (but less so than the information in the higher spaces); we use a line with intersecting hash marks to denote when this is not the case. Subscripted terms are used to indicate the relevant contents of the focus spaces because the spaces contain representations of entities (i.e., objects, properties, and relations) and not linguistic expressions.

Part one of Figure 1 shows the state of focusing when discourse segment DS2 is being processed. Segment DS1 gave rise to FS1 and had as its discourse purpose DSP<sub>1</sub>. The properties, objects, relations, and purpose represented in FS1 are accessible but less salient than those in FS2. DS2 yields a focus space that is stacked relative to FS1 because DSP<sub>1</sub> of DS1 dominates DS2's DSP, DSP<sub>2</sub>. As a result of the relationship between FS1 and FS2, reduced noun phrases will be interpreted differently in DS2 than in DS1. For example, if some red balls exist in the world one of which is represented in DS2 and another in FS1, then *the red ball* used in DS2 will be understood to mean the particular red ball that is represented in DS2. If, however, there is also a green truck (in the world) and it is represented only in FS1, *the green truck* uttered in DS2 will be understood as referring to that green truck.

Part two of Figure 1 shows the state of focusing when segment DS3 is being processed. FS2 has been popped from the stack and FS3 has been pushed onto it because the DSP of DS3, DSP<sub>3</sub>, is dominated solely by DSP<sub>1</sub>, not by DSP<sub>2</sub>. In this example, the intentional structure includes only dominance relationships, although, it may, in general, also include satisfaction-precedence relationships.

The stacking of focus spaces reflects the relative salience of the entities in each space during the corresponding segment's portion of the discourse. The stack relationships arise from the ways in which the various DSPs relate; information about such relationships is

represented in the dominance hierarchy (depicted on the right in the figure). The spaces in Figure 1 are snapshots illustrating the results of a sequence of operations, such as pushes onto and pops from a stack. A push occurs when the DSP for a new segment contributes to the DSP for the immediately preceding segment. When the DSP contributes to some intention higher in the dominance hierarchy, several focus spaces are popped from the stack before the new one is inserted.

Two essential properties of the focusing structure are now clear. First, the focusing structure is parasitic upon the intentional structure, in the sense that the relationships among DSPs determine pushes and pops. Note however, that the relevant operation may sometimes be indicated in the language itself. For example, the cue word *first* often indicates the start of a segment whose DSP contributes to the DSP of the preceding segment. Second, the focusing structure, like the intentional and linguistic structures, evolves as the discourse proceeds. None of them exists a priori. Even in those rare cases in which an ICP has a complete plan for the discourse prior to uttering a single word, the intentional structure is constructed by the CPs as the discourse progresses. This discourse-time construction of the intentional structure may be more obviously true for speakers and hearers of spoken discourse than for readers and writers of texts, but, even for the writer, the intentional structure is developed as the text is being written.

Figure 1 illustrates some fundamental distinctions between the intentional and attentional components of discourse structure. First, the dominance hierarchy provides, among other things, a complete record of the discourse-level intentions and their dominance (as well as, when relevant, satisfaction-precedence) relationships, whereas the focusing structure at any one time can essentially contain only information that is relevant to purposes in a portion of the dominance hierarchy. Second, at the conclusion of a discourse, if it completes normally, the focus stack will be empty, while the intentional structure will have been fully constructed. Third, when the discourse is being processed, only the attentional state can constrain the interpretation of referring expressions directly.

We can now also clarify some misinterpretations of focus-space diagrams and task structure in our earlier work (Grosz 1978a, 1981, 1974). The focus-space hierarchies in that work are best seen as representing attentional state. The task structure was used in two ways:

1. to represent common knowledge about the task;
2. as a special case of the intentional structure we posit in this paper.

Although the same representational scheme was used for encoding the focus-space hierarchies and the task structure (partitioned networks: Hendrix 1979), the two structures were distinct.

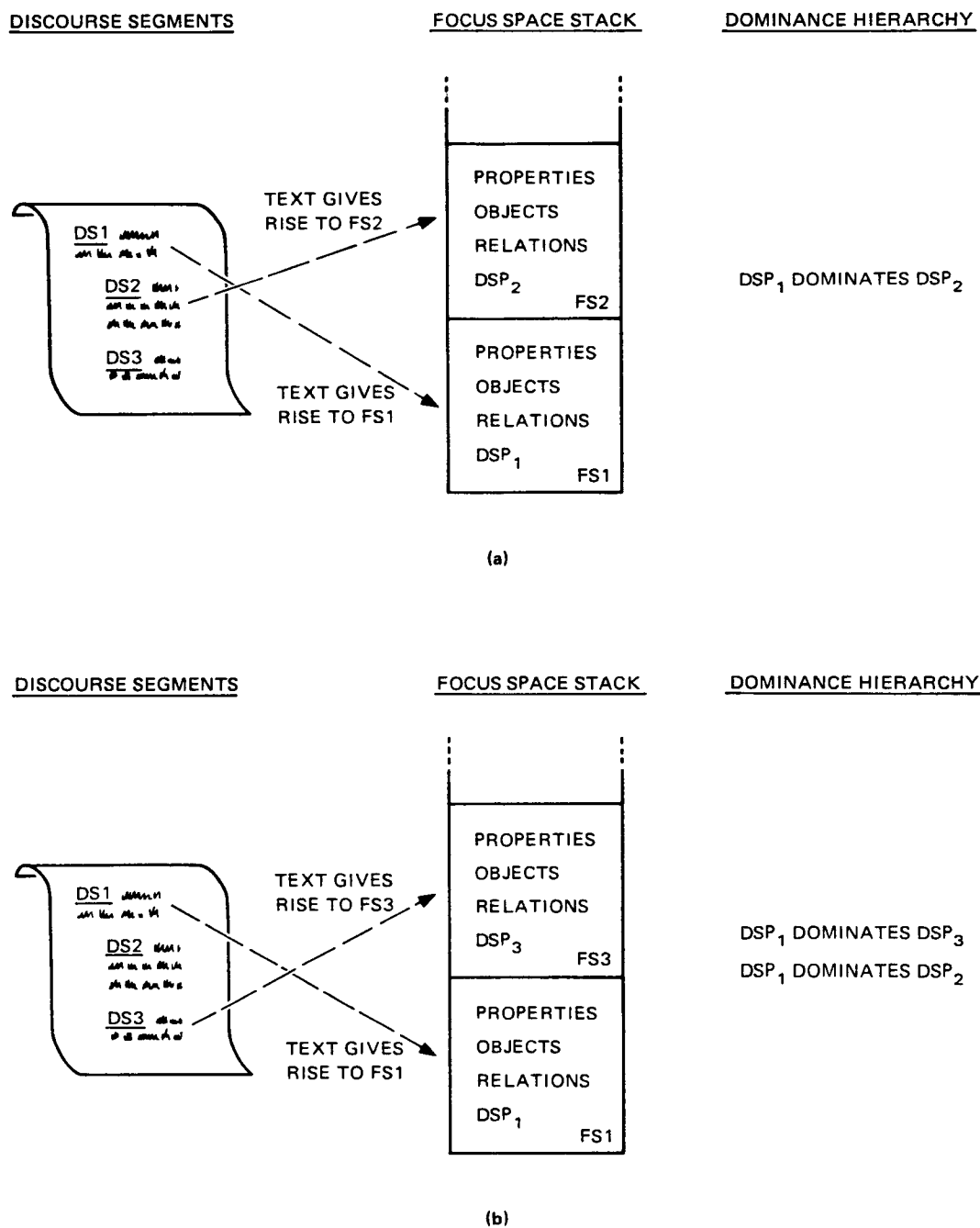


Figure 1. Discourse Segments, Focus Spaces and Dominance Hierarchy.



Several researchers (e.g., Linde and Goguen 1978, Reichman-Adar 1984) misinterpreted the original research in an unfortunate and unintended way: they took the focus-space hierarchy to include (or be identical to) the task structure. The conflation of these two structures forces a single structure to contain information about attentional state, intentional relationships, and general task knowledge. It prevents a theory from accounting adequately for certain aspects of discourse, including interruptions (see Section 5).

A second instance of confusion was to infer (incorrectly) that the task structure was necessarily a prebuilt tree. If the task structure is taken to be a special case of intentional structure, it becomes clear that the tree structure is simply a more constrained structure than one might require for other discourses; the nature of the task related to the task-oriented discourse is such that the dominance hierarchy of the intentional structure of the dialogue has both dominance and satisfaction-precedence relationships,<sup>5</sup> while other discourses may not exhibit significant precedence constraints among the DSPs. Furthermore, there has never been any reason to assume that the task structures in task-oriented dialogues are prebuilt, any more than the intentional structure of any other kind of discourses. It is rather that one objective of discourse theory (not a topic considered here, however) is to explain how the OCP builds up a model of the task structure by using information supplied in the discourse.

However, it is important to note that conflating the aforementioned two roles of information about the task itself (as a portion of general commonsense knowledge and as a special case of intentional structure) was regrettable, as it fails to make an important distinction. Furthermore, as is clear when intentional structures are considered more generally, such a conflation of roles does not allow for differences between what one knows about a task and one's intentions for (or what one makes explicit in discourse about) performing a task.

In summary, the focusing structure is the central repository for the contextual information needed to process utterances at each point in the discourse. It distinguishes those objects, properties, and relations that are most salient at that point and, moreover, has links to relevant parts of both the linguistic and intentional structures. During a discourse, an increasing amount of infor-

mation, only some of which continues to be needed for the interpretation of subsequent utterances, is discussed. Hence, it becomes more and more necessary to be able to identify relevant discourse segments, the entities they make salient, and their DSPs. The role of attentional state in delineating the information necessary for understanding is thus central to discourse processing.

### 3 TWO EXAMPLES

To illustrate the basic theory we have just sketched, we will give a brief analysis of two kinds of discourse: an argument from a rhetoric text and a task-oriented dialogue. For each example we discuss the segmentation of the discourse, the intentions that underlie this segmentation, and the relationships among the various DSPs. In each case, we point out some of the linguistic devices used to indicate segment boundaries as well as some of the expressions whose interpretations depend on those boundaries. The analysis is concerned with specifying certain aspects of the behavior to be explicated by a theory of discourse; the remainder of the paper provides a partial account of this behavior.

#### 3.1 AN ARGUMENT

Our first example is an argument taken from a rhetoric text (Holmes and Gallagher 1917<sup>6</sup>). It is an example used by Cohen (1983) in her work on the structure of arguments. Figure 2 shows the dialogue and the eight discourse segments of which it is composed. The division of the argument into separate (numbered) clauses is Cohen's, but our analysis of the discourse structure is different, since in Cohen's analysis, every utterance is directly subordinated to another utterance, and there is only one structure to encode linguistic segmentation and the purposes of utterances. Although both analyses segment utterance (4) separately from utterances (1-3), some readers place this utterance in DS1 with utterances (1) through (3); this is an example of the kind of disagreement about boundary utterances found in Mann's data (as discussed in Section 2.1). The two placements lead to slightly different DSPs, but not to radically different intentional structures. Because the differences do not affect the major thrust of the argument, we will discuss only one segmentation.

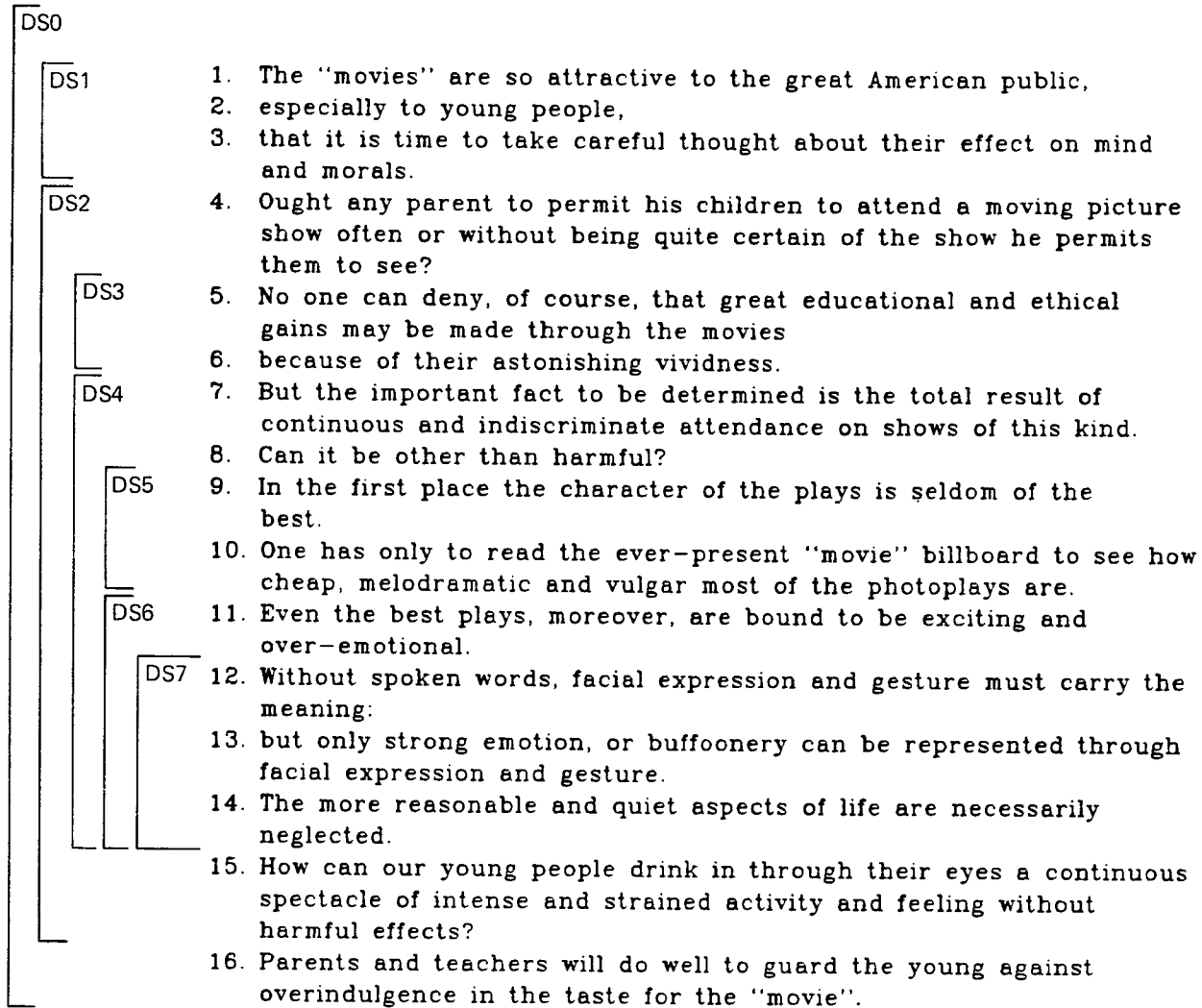


Figure 2. The Movies Essay.

Figure 3 lists the primary component of the DSP for each of these segments and Figure 4 shows the dominance relationships that hold among these intentions. In Section 7 we discuss additional components of the discourse segment purpose; because these additional components are more important for completeness of the theory than for determining the essential dominance and

satisfaction-precedence relationships between DSPs, we omit such details here. Rather than commit ourselves to a formal language in which to express the intentions of the discourse, we will use a shorthand notation and English sentences that are intended to be a gloss for a formal statement of the actual intentions.

I0: (Intend ICP (Believe OCP P0))

where P0 = the proposition that parents and teachers should guard the young from overindulgence in the movies.

I1: (Intend ICP (Believe OCP P1))

where P1 = the proposition that it is time to consider the effect of movies on mind and morals.

I2: (Intend ICP (Believe OCP P2))

where P2 = the proposition that young people cannot drink in through their eyes a continuous spectacle of intense and strained activity without harmful effects.

I3: (Intend ICP (Believe OCP P3))

where P3 = the proposition that it is undeniable that great educational and ethical gains may be made through the movies.

I4: (Intend ICP (Believe OCP P4))

where P4 = the proposition that although there are gains, the total result of continuous and indiscriminate attendance at movies is harmful.

I5: (Intend ICP (Believe OCP P5))

where P5 = the proposition that the content of movies (i.e., the character of the plays) is not the best.

I6: (Intend ICP (Believe OCP P6))

where P6 = the proposition that the stories (i.e., the plays) in movies are exciting and over-emotional.

I7: (Intend ICP (Believe OCP P7))

where P7 = the proposition that movies portray strong emotion and buffoonery while neglecting the quiet and reasonable aspects of life.

**Figure 3.** Primary intentions of the DSPs for Movies essay.

I0 DOM I1  
I0 DOM I2  
I2 DOM I3  
I2 DOM I4  
I4 DOM I5  
I4 DOM I6  
I6 DOM I7

**Figure 4.** Dominance relationships for the DSPs of the Movies essay.

All the primary intentions for this essay are intentions that the reader (OCP) come to believe some proposition. Some of these propositions, such as P5 and P6, can be read off the surface utterances directly. Other propositions and the intentions of which they are part, such as P2 and I2, are more indirect. Like the Gricean utterance-level intentions (the analogy with these will be explored in Section 7), DSPs may or may not be directly expressed in the discourse. In particular, they may be expressed in any of the following ways:

1. *explicitly* as in *I intend for you to believe that it's time to consider the effects of movies on mind and morals.* [which would produce I1]
2. *directly, in one utterance*, as in (3) [which does produce I1]
3. *directly, through multiple utterances*, as in using (7) and the utterance *It can only be harmful* to produce I4,
4. *by derivation, in one or more utterances with an associated context*, as in (15) to produce I2.

Not only may information about the DSP be conveyed by a number of features of the utterances in a discourse, but it also may come in any utterance in a segment. For example, although I0 is the DP, it is stated directly only in the last utterance of the essay. This leads to a number of questions about the ways in which OCPs can recognize discourse purposes, and about those junctures at which they need to do so. We turn to these matters directly in Subsection 4.1.

This discourse also provides several examples of the different kinds of interactions that can hold between the linguistic expressions in a discourse and the discourse structure. It includes examples of the devices that may be used to mark overtly the boundaries between discourse segments – examples of the use of aspect, mood, and particular cue phrases – as well as of the use of referring expressions that are affected by discourse segment boundaries.

The use of cue phrases to indicate discourse boundaries is illustrated in utterances (9) and (11); in (9) the phrase *in the first place* marks the beginning of DS5 while in (11) *moreover* ends DS5 and marks the start of DS6. These phrases also carry information about the intentional structure, namely, that DSP5 and DSP6 are dominated by DSP4. In some cases, cue phrases have multiple functions; they convey propositional content as well as marking discourse segment boundaries. The *but* in utterance (7) is an example of such a multiple function use.

The boundaries between DS1 and DS2, DS4 and DS5, and DS4 and DS2 reflect changes of aspect and mood. The switch from declarative, present tense to interrogative modal aspect does not in itself seem to signal the boundary (for recognition purposes) in this discourse unambiguously, but it does indicate a possible line of demarcation which, in fact, is valid.

The effect of segmentation on referring expressions is shown by the use of the generic noun phrase *a moving*

*picture show* in (4). Although a reference to the movies was made with a pronoun (*their*) in (3), a full noun phrase is used in (4). This use reflects, and perhaps in part marks, the boundary between the segments DS1 and DS2.

Finally, this discourse has an example of the trade-off between explicitly marking a discourse boundary, as well as the relationship between the associated DSPs, and reasoning about the intentions themselves. There is no overt linguistic marker of the beginning of DS7; its separation must be inferred from DSP7 and its relationship to DSP6.

### 3.2 A TASK-ORIENTED DIALOGUE

The second example is a fragment of a task-oriented dialogue taken from Grosz (1981; it is from the same corpus that was used by Grosz 1974). Figure 5 contains the dialogue fragment and indicates the boundaries for its main segments.<sup>7</sup> Figure 6 gives the primary component of the DSPs for this fragment and shows the dominance relationships between them.

In contrast with the movies essay, the primary components of the DSPs in this dialogue are mostly intentions of the segment's ICP that the OCP intend to perform some action. Also, unlike the essay, the dialogue has two agents initiating the different discourse segments. In this particular segment, the expert is the ICP of DS1 and DS5, while the apprentice is the ICP of DS2-4. To furnish a complete account of the intentional structure of this discourse, one must be able to say how the satisfaction of one agent's intentions can contribute to satisfying the intentions of another agent. Such an account is beyond the scope of this paper, but in Section 7 we discuss some of the complexities involved in providing one (as well as its role in discourse theory).

For the purposes of discussing this example, though, we need to postulate two properties of the relationships among the participants' intentions. These properties seem to be rooted in features of cooperative behavior and depend on the two participants' sharing some particular knowledge of the task. First, it is a shared belief that, unless he states otherwise, the OCP will adopt the intention to perform an action that the ICP intended him to. Second, in adopting the intention to carry out that action, the OCP also intends to perform whatever subactions are necessary. Thus, once the apprentice intends to remove the flywheel, he also commits himself to the collateral intentions of loosening the setscrews and pulling the wheel off. Note, however, that not all the subactions need to be introduced explicitly into the discourse. The apprentice may do several actions that are never mentioned, and the expert may assume that these are being undertaken on the basis of other information that the apprentice obtains. The partiality of the intentional structure stems to some extent from these characteristics of intentions and actions.

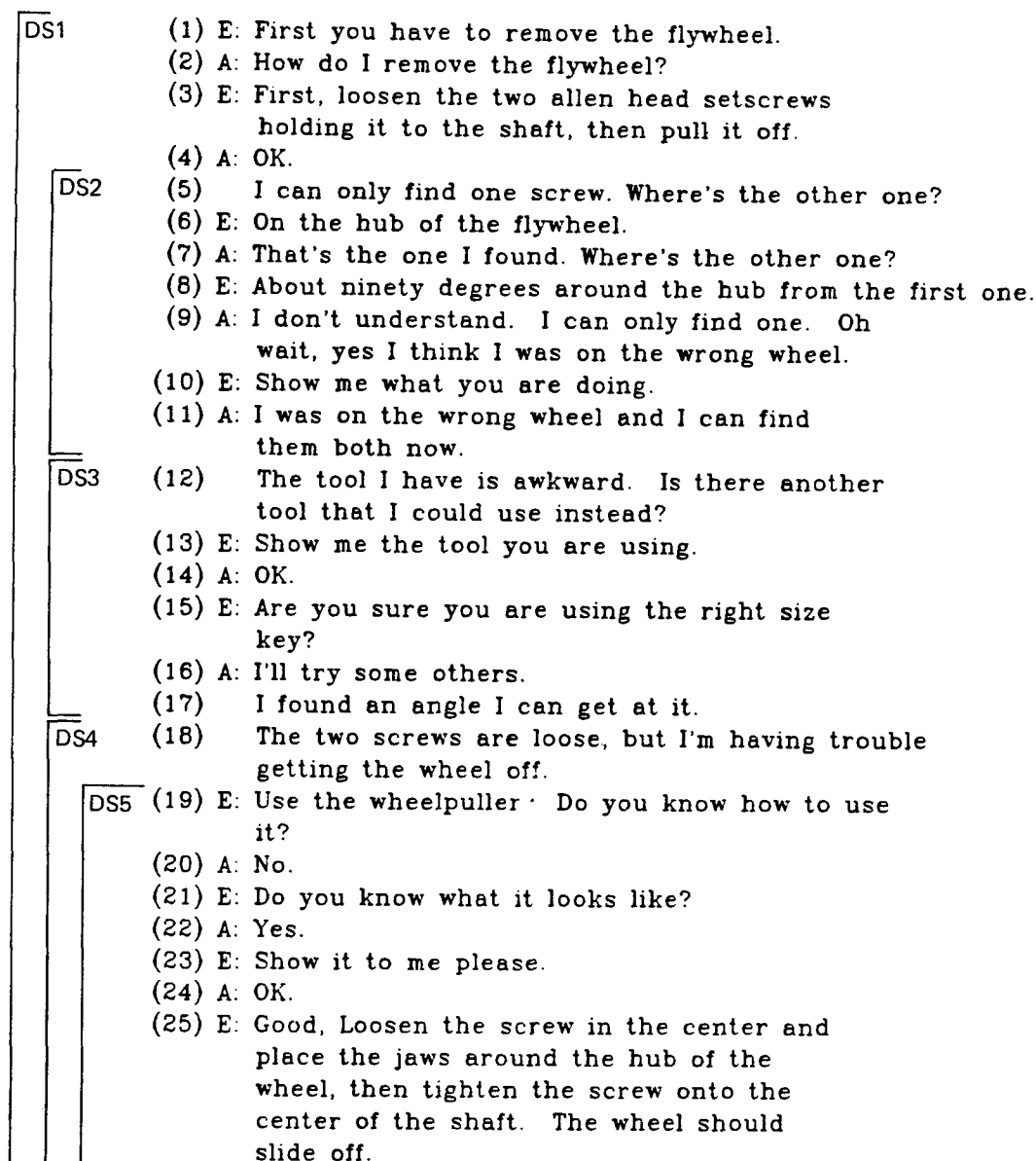


Figure 5. A segment of a task-oriented dialogue.

As in the movies essay, some of the DSPs for this dialogue are expressed directly in utterances. For instance, utterances (1), (5), and (12) directly express the primary components of DSP1, DSP2 and DSP3, respectively. The primary component of DSP4 is a derived intention. The surface intention of *but I'm having trouble getting the wheel off* is that the apprentice intends the expert to believe that the apprentice is having trouble taking off the flywheel. I4 is derived from the utterance and its surface intention, as well as from features of discourse, conventions about what intentions are associated with the *I am having trouble doing X* type

of utterance, and what the ICP and OCP know about the task they have undertaken.

The dominance relationship that holds between I1 and I2, as well as the one that holds between I1 and I3, may seem problematic at first glance. It is not clear how locating any single setscrew contributes to removing the flywheel. It is even less clear how, in and of itself, identifying another tool does. Two facts provide the link: first, that the apprentice (the OCP of DS1) has taken on the task of removing the flywheel; second, that the apprentice and expert share certain knowledge about the task. Some of this shared task knowledge comes from the discourse per se [e.g., utterance (3)], but some of it

comes from general knowledge, perceptual information, and the like. Thus, a combination of information is relevant to determining I2 and I3 and their relationships to I1, including all of the following: the fact that I1 is part of the intentional structure, the fact that the apprentice is currently working on satisfying I1, the utterance-level intentions of utterances (5) and (12), and general knowledge about the task.

The satisfaction-precedence relations among I2, I3, and I4 are not communicated directly in the dialogue, but, like dominance relations, depend on domain knowledge. One piece of relevant knowledge is that a satisfaction precedence relation exists between loosening the setscrews and pulling off the flywheel. That relation is shared knowledge that is stated directly (*First loosen ..., then pull*). The relation, along with the fact that both I2 and I3 contribute to loosening the setscrews, and that I4 contributes to pulling off the flywheel, makes it possible to conclude I3 SP I4 and I2 SP I4. To conclude that I2 SP I3, the apprentice must employ knowledge of how to go about loosening screw-like objects.

The dominance and satisfaction-precedence relations for this task-oriented fragment form a tree of intentions rather than just a partial ordering. In general, however, for any fragment, task-oriented or otherwise, this is not necessary.

It is essential to notice that the intentional structure is neither identical to nor isomorphic to a general plan for removing the flywheel. It is not identical because a plan encompasses more than a collection of intentions and relationships between them (compare Pollack's (1986) critique of AI planning formalisms as the basis for inferring intentions in discourse). It is not isomorphic because the intentional structure has a different substructure from the general plan for removing the flywheel. In addition to the intentions arising from steps in the plan, the inten-

tional structure typically contains DSPs corresponding to intentions generated by the particular execution of the task and the dialogue. For example, the general plan for the disassembly of a flywheel includes subplans for loosening the setscrews and pulling off the wheel; it might also include subplans (of the loosening step) for finding the setscrews, finding a tool with which to loosen the screws, and loosening each screw individually. However, this plan would not contain contingency subplans for what to do when one cannot find the screws or realizes that the available tool is unsatisfactory. Intentions I2 and I3 stem from difficulties encountered in locating and loosening the setscrews. Thus, the intentional structure for this fragment is not isomorphic to the general plan for removing the flywheel.

Utterance (18) offers another example of the difference between the intentional structure and a general plan for the task. This utterance is part of DS4 – not just part of DS1 – even though it contains references to more than one single part of the overall task (which is what I1 is about). It functions to establish a new DSP, I4, as most salient. Rather than being regarded as a report on the overall status of the task, the first clause is best seen as modifying the DSP.<sup>8</sup> With it, the apprentice tells the expert that the trouble in removing the wheel is not with the screws. Thus, although general task knowledge is used in determining the intentional structure, it is not identical to it.

In this dialogue, there are fewer instances in which cue phrases are employed to indicate segment boundaries than occur in the movies essay. The primary example is the use of *first* in (1) to mark the start of the segment and to indicate that its DSP is the first of several intentions whose satisfaction will contribute to satisfying the larger discourse of which they are a part.

#### Primary Intentions:

- I1: (Intend E<sub>xpert</sub> (Intend A<sub>pprentice</sub> (Remove A flywheel)))
- I2: (Intend A (Intend E (Tell E A (Location other setscrew))))
- I3: (Intend A (Intend E (Identify E A another tool)))
- I4: (Intend A (Intend E (Tell E A (How (Getoff A wheel)))))
- I5: (Intend E (Know-How-to A (Use A wheelpuller)))

#### Dominance Relationships:

- I1 DOM I2
- I1 DOM I3
- I1 DOM I4
- I4 DOM I5

#### Satisfaction-Precedence Relationships:

- I2 SP I3
- I2 SP I4
- I3 SP I4

Figure 6. Intentional structure for the task-oriented dialogue segment.

The dialogue includes a clear example of the influence of discourse structure on referring expressions. The phrase *the screw in the center* is used in (25) to refer to the center screw of the wheelpuller, not one of the two setscrews mentioned in (18). This use of the phrase is possible because of the attentional state of the discourse structure at the time the phrase is uttered.

#### 4 PROCESSING ISSUES

In previous sections of the paper, we abstracted from the cognitive states of the discourse participants. The various components of discourse structure discussed so far are properties of the discourse itself, not of the discourse participants. To use the theory in constructing computational models requires determining how each of the individual components projects onto the model of an individual discourse participant. In this regard, the principal issues include specifying

1. how the ICP indicates and the OCP recognizes the beginning and end of a discourse segment,
2. how the OCP recognizes the discourse segment purposes, and
3. how the focus space stack operates.

In essence, the OCP must judge for each utterance whether it starts a new segment, ends the current one (and possibly some of its embedding segments), or contributes to the current one. The information available to the OCP for recognizing that an utterance starts a new segment includes any explicit linguistic cues contained in the utterance (see Section 6<sup>9</sup>) as well as the relationship between its utterance-level intentions and the active DSPs (i.e., those in some focus space that is still on the stack). Likewise, the fact that an utterance ends a segment may be indicated explicitly by linguistic cues or implicitly from its utterance-level intentions and their relationship to elements of the intentional structure. If neither of these is the case, the utterance is part of the current segment. Thus, intention recognition and focus space management play key roles in processing. Moreover, they are also related: the intentional structure is a primary factor in determining focus space changes, and the focus space structure helps constrain the intention recognition process.

##### 4.1 INTENTION RECOGNITION

The recognition of DP/DSPs is the central issue in the computational modeling of intentional structure. If, as we have claimed, for the discourse to be coherent and comprehensible, the OCP must be able to recognize both the DP/DSPs<sup>10</sup> and relationships (dominance and satisfaction-precedence) between them, then the question of how the OCP does so is a crucial issue.

For the discourse as a whole, as well as for each of its segments, the OCP must identify both the intention that serves as the discourse segment purpose and its relationship to other discourse-level intentions. In particular, the OCP must be able to recognize which other DSPs that

specific intention dominates and is dominated by, and, where relevant, with which other DSPs it has satisfaction-precedence relationships. Two issues that are central to the recognition problem are what information the OCP can utilize in effecting the recognition and at what point in the discourse that information becomes available.

An adequate computational model of the recognition process depends critically on an adequate theory of intention and action; this, of course, is a large research problem in itself and one not restricted to matters of discourse. The need to use such a model for discourse, however, adds certain constraints on the adequacy of any theory or model. Pollack (1986) describes several properties such theories and models must possess if they are to be adequate for supporting recognition of intention in single-utterance queries; she shows how current AI planning models are inadequate and proposes an alternative planning formalism. The need to enable recognition of discourse-level intentions leads to yet another set of requirements.

As will become clear in what follows, the information available to the OCP comes from a variety of sources. Each of these can typically provide partial information about the DSPs and their relationships. These sources are each partially constraining, but only in their ensemble do they constrain in full. To the extent that more information is furnished by any one source, commensurately less is needed from the others. The overall processing model must be one of constraint satisfaction that can operate on partial information. It must allow for incrementally constraining the range of possibilities on the basis of new information that becomes available as the segment progresses.

##### 4.1.1 INFORMATION CONSTRAINING THE DSP

At least three different kinds of information play a role in the determination of the DSP: specific linguistic markers, utterance-level intentions, and general knowledge about actions and objects in the domain of discourse. Each plays a part in the OCP's recognition of the DSP and can be utilized by the ICP to facilitate this recognition.

Cue phrases are the most distinguished linguistic means that speakers have for indicating discourse segment boundaries and conveying information about the DSP. Recent evidence by Hirschberg and Pierrehumbert (1986) suggests that certain intonational properties of utterances also provide partial information about the DSP relationships. Because some cue phrases may be used as clausal connectors, there is a need to distinguish their discourse use from their use in conveying propositional content at the utterance level. For example, the word *but* functions as a boundary marker in utterance (7) of the discourse in Section 3.1, but it can also be used solely (as in the current utterance) to convey propositional content (e.g., the conjunction of two propositions) and serve to connect two clauses within a segment.

As discussed in Section 6, cue phrases can provide information about dominance and satisfaction-precedence relationships between segments' DSPs. However, they may not completely specify which DSP dominates or satisfaction-precedes the DSP of the segment they start. Furthermore, cue phrases that explicitly convey information only about the attentional structure (see Section 6) may be ambiguous about the state to which attention is to shift. For example, if there have been several interruptions (see Section 5), the phrase *but anyway* indicates a return to some previously interrupted discourse, but does not specify which one. Although cue phrases do not completely specify a DSP, the information they provide is useful in limiting the options to be considered.

The second kind of information the OCP has available is the utterance-level intention of each utterance in the discourse. As the discussion of the movies example (Section 3.1) pointed out, the DSP may be identical to the utterance-level intention of some utterance in the segment. Alternatively, the DSP may combine the intentions of several utterances, as is illustrated in the following discourse segment:

I want you to arrange a trip for me to Palo Alto.  
It will be for two weeks.  
I only fly on TWA.

The DSP for this segment is, roughly, that the ICP intends for the OCP to make (complete) trip arrangements for the ICP to go to Palo Alto for two weeks, under the constraint that any flights be on TWA. The Gricean intentions for these three utterances are as follows:

- Utterance1:** ICP intends that OCP believe that ICP intends that OCP intend to make trip plans for ICP to go to Palo Alto
- Utterance2:** ICP intends that OCP believe that ICP intends OCP to believe that the trip will last two weeks
- Utterance3:** ICP intends that OCP believe that ICP intends OCP to believe that ICP flies only on TWA

These intentions must be combined in some way to produce the DSP. The process is quite complex, since the OCP must recognize that the reason for utterances 2 and 3 is not simply to have some new beliefs about the ICP, but to use those beliefs in arranging the trip. While this example fits the schema of a request followed by two informings, schemata will not suffice to represent the behavior as a general rule. A different sequence of utterances with different utterance-level intentions can have the same DSP; this is the case in the following segment:

- S1: Have I told you yet to arrange my trip to Palo Alto?  
Remember that I will fly only on TWA. OK?
- S2: OK.
- S3: I'm planning on staying for two weeks.

It is possible for a sequence that consists of a request followed by two informings not to result in a modifica-

tion of the trip plans. For example, in the following sequence the third utterance results in changing the way the arrangements are made, rather than constraining the nature of the arrangements themselves.

I want you to arrange a two-week trip for me to Palo Alto. I fly only on TWA. The rates go up tomorrow, so you'll want to call today.

Not only is the contribution of utterance-level intentions to DSPs complicated, but in some instances the DSP for a segment may both constrain and be partially determined by the Gricean intention for some utterance in the segment. For example, the Gricean-intention for utterance (15) in the movies example (Section 3.1) is derived from a combination of facts about the utterance itself, and from its place in the discourse. On the surface, (15) appears to be a question addressed to the OCP; its intention would be roughly that the ICP intends the OCP to believe that the ICP wants to know how young people, etc. But (15) is actually a rhetorical question and has a very different intention associated with it – namely, that the ICP intends the OCP to believe proposition P2 (namely, that young people cannot drink in through their eyes a continuous spectacle of intense and strained activity without harmful effects). In this example, this particular intention is also the primary component of the DSP.

The third kind of information that plays a role in determining the DP/DSPs is shared knowledge about actions and objects in the domain of discourse. This shared knowledge is especially important when the linguistic markers and utterance-level intentions are insufficient for determining the DSP precisely.

In Section 7 we introduce two relations, a **supports** relation between propositions and a **generates** relation between actions, and present two rules stating equivalences; one links a dominance relation between two DSPs with a supports relation between propositions and the other links a dominance relation between DSPs to a generates relation between actions. Use of these rules in one direction allows for (partially) determining what supports or generates relationship holds from the dominance relationship. But the rules can be used in the opposite direction also: if, from the content of utterances and reasoning about the domain of discourse, a supports or generates relationship can be determined, then the dominance relationship between DSPs can be determined. In such cases it is important to derive the dominance relationship so that the appropriate intentional and attentional structures are available for processing or determining the interpretation of the subsequent discourse.

From the perspective of recognition, a trade-off implicit in the two equivalences is important. If the ICP makes the dominance relationship between two DSPs explicit (e.g., with cue phrases), then the OCP can use this information to help recognize the (ICP's beliefs about the) supports relationship. Conversely, if the ICP's utterances make clear the (ICP's beliefs about the) supports or generates relationship, then the OCP can use



this information to help recognize the dominance relationship. Although it is most helpful to use the dominance relationships to constrain the search for appropriate supports and generates relationships, sometimes these latter relationships can be inferred reasonably directly from the utterances in a segment using general knowledge about the objects and actions in the domain of discourse. It remains an open question what inferences are needed and how complex it will be to compute supports and generates relationships if the dominance relationship is not directly indicated in a discourse.

Utterances from the movies essay illustrate this trade-off. In utterance (9), the phrase *in the first place* expresses the dominance relationship between DSPs of the new segment DS5 and the parent segment DS4 directly. Because of the dominance relationship (as well as the intentions expressed in the utterances), the OCP can determine that the ICP believes that the proposition that the content of the plays is not the best provides support for the proposition that the result of indiscriminate movie going is harmful. Hence determining dominance yields the support relation. The support relation can also yield dominance. Utterances (12)-(14), which comprise DS7, are not explicitly marked for a dominance relation. It can be inferred from the fact that the propositions in (12)-(14) provide support for the proposition embedded in DSP6 (that is, that the stories in movies are exciting and over-emotional) that DSP6 dominates DSP7.

Finally, the more information an ICP supplies explicitly in the actual utterances of a discourse, the less reasoning about domain information an OCP has to do to achieve recognition. Cohen (1983) has made a similar claim regarding the problem of recognizing the relationship between one proposition and another.

#### 4.1.2 WHEN IS THE INTENTION RECOGNIZED?

As discussed in Section 2.2, the intentional structure evolves as the discourse does. By the same token, the discourse participants' mental-state correlates of the intentional structure are not prebuilt; neither participant may have a complete model of the intentional structure "in mind" until the discourse is completed. The dominance relationships that actually shape the intentional structure cannot be known a priori, because the specific intentions that will come into play are not known (never by the OCP, hardly ever by the ICP) until the utterances in the discourse have been made. Although it is assumed that the participants' common knowledge includes<sup>11</sup> enough information about the domain to determine various relationships such as supports and generates, it is not assumed that, prior to a discourse, they actually had inferred and are aware of all the relationships they will need for that discourse.

Because any of the utterances in a segment may contribute information relevant to a complete determination of the DSP, the recognition process is not complete until the end of the segment. However, the OCP

must be able to recognize at least a generalization of the DSP so that he can make the proper moves with respect to the attentional structure. That is, some combination of explicit indicators and intentional and propositional content must allow the OCP to ascertain where the DSP will fit in the intentional structure at the beginning of a segment, even if the specific intention that is the DSP cannot be determined until the end of the segment.

Utterance (15) in the movies example illustrates this point. The author writes, "How can our young people drink in through their eyes a continuous spectacle of intense and strained activity and feeling without harmful effects?" The primary intention I2 is derived from this utterance, but this cannot be done until very late in the discourse segment [since (15) occurs at the end of DS2]. Furthermore, the segment for which I2 is primary has complex embedding of other segments. Utterance (16), intention I0, and DS0 constitute another example of the expression of a primary intention late in a discourse segment. In that case, I0 cannot be computed until (16) has been read, and (16) is not only the last utterance in DS0, but is one that covers the entire essay. If an OCP must recognize a DSP to understand a segment, then we ask: how does the OCP recognize a DSP when the utterance from which its primary intention is derived comes so late in the segment?

We conjecture with regard to such segments as D2 of the movies essay that the primary intention (e.g., I2) may be determined partially (and hence a generalized version become recognizable) before the point at which it is actually expressed in the discourse. While the DP/DSP may not be expressed early, there is still *partial information* about it. This partial information often suffices to establish dominance (or satisfaction-precedence) relationships for additional segments. As these latter are placed in the hierarchy, their DSPs can provide further partial information for the underspecified DSP. For example, even though the intention I0 is expressed directly only in the last utterance of the movies essay, utterance (4) expresses an intention to know whether *p* or  $\sim p$  is true (i.e., whether or not parents should let children see movies often and without close monitoring). I0 is an intention to believe, whose proposition is a generalization of the  $\sim p$  expressed in (4). Consider also the primary intention I4. It occurs in a segment embedded within DS2, is more general than I2, but is an approximation to it. It would not be surprising to discover that OCPs can in fact predict something close to I2 on the basis of I4, utterances (9)-(14), and the partial dominance hierarchy available at each point in the discourse.

#### 4.2 USE OF THE ATTENTIONAL STATE MODEL

The focus space structure enables certain processing decisions to be made locally. In particular, it limits the information that must be considered in recognizing the DSP as well as that considered in identifying the referents of certain classes of definite noun phrases.

A primary role of the focus space stack is to constrain the range of DSPs considered as candidates for domination or satisfaction-precedence of the DSP of the current segment. Only those DSPs in some space on the focusing stack are viable prospects. As a result of this use of the focusing structure, the theory predicts that this decision will be a local one with respect to attentional state. Because two focus spaces may be close to each other in the attentional structure without the discourse segments they arise from necessarily being close to one another and vice versa, this prediction corresponds to a claim that locality in the focusing structure is what matters to determination of the intentional structure.

A second role of the focusing structure is to constrain the OCP's search for possible referents of definite noun phrases and pronouns. To illustrate this role, we will consider the phrase *the screw in the center* in utterance (25) of the task-oriented dialogue of Section 3. The focus stack configuration when utterance (25) is spoken is shown in Figure 7. The stack contains (in bottom-to-top order) focus spaces FS1, FS4, and FS5 for segments DS1, DS4, and DS5, respectively. For DS5 the wheelpuller is a focused entity, while for DS4 the two setscrews are (because they are explicitly mentioned). The entities in FS5 are considered before those in FS4 as potential referents. The wheelpuller has three screws: two small screws fasten the side arms, and a large screw in the center is the main functioning part. As a result, this large screw is implicitly in focus in FS5 (Grosz 1977) and thus identified as the referent without the two setscrews ever being considered.

Attentional state also constrains the search for referents of pronouns. Because pronouns contain less explicit information about their referents than definite descriptions, additional mechanisms are needed to account for what may and may not be pronominalized in the discourse. One such mechanism is centering (which we previously called immediate focusing; Grosz, Joshi, and Weinstein 1983; Sidner 1979).

Centering, like focusing, is a dynamic behavior, but is a more local phenomenon. In brief, a backward-looking center is associated with each utterance in a discourse segment; of all the focused elements the backward-looking center is the one that is central in that utterance (i.e., the uttering of the particular sequence of words at that point in the discourse). A combination of syntactic, semantic, and discourse information is used to identify the backward-looking center. The fact that some entity is the backward-looking center is used to constrain the search for the referent of a pronoun in a subsequent utterance. Note that unlike the DSP, which is constant for a segment, the backward-looking center may shift: different entities may become more salient at different points in the segment.

The presence of both centers and DSPs in this theory leads us to an intriguing conjecture: that "topic" is a concept that is used ambiguously for both the DSP of a segment and the center. In the literature the concept of "topic" has appeared in many guises. In syntactic form it is used to describe the preposing of syntactic constituents in English and the "wa" marking in Japanese. Researchers have used it to describe the sentence topic (i.e., what the sentence is about; Firbas 1971, Sgall, Hajičová, and

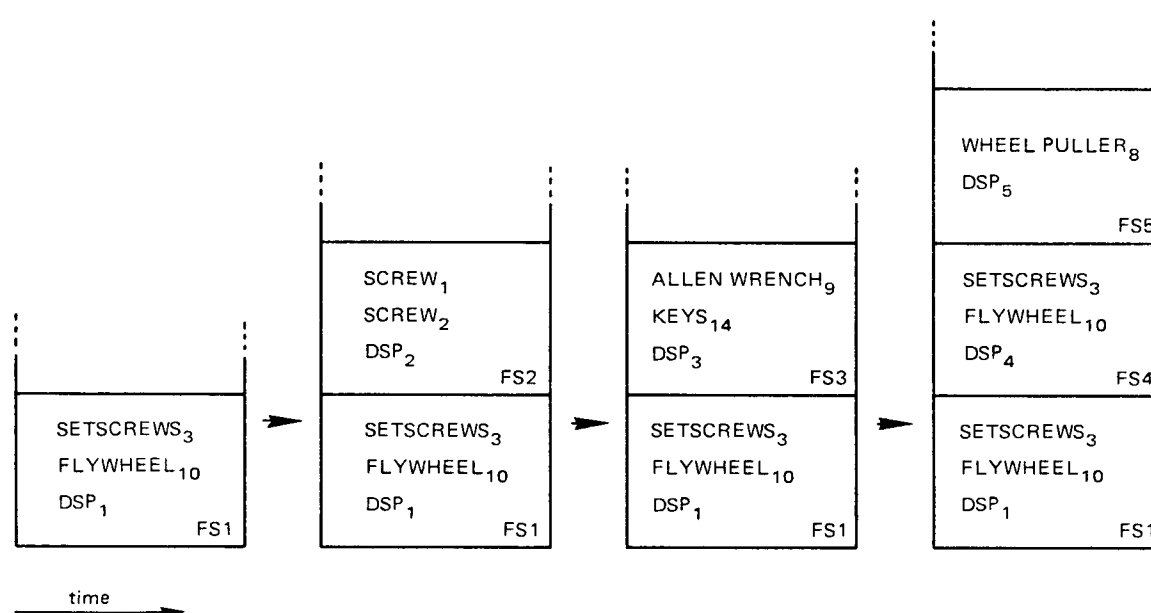


Figure 7. Focus Stack Transitions Leading up to Utterance (25).

Benesova 1973), and as a pragmatic notion (Reinhart 1981); others want to use the term for discourse topic, either to mean what the discourse is about, or to be defined as those proposition(s) the ICP provides or requests new information about (see Reinhart (1981) for a review of many of the notions of aboutness and topic). It appears that many of the descriptions of sentence topic correspond (though not always) to centers, while discourse topic corresponds to the DSP of a segment or of the discourse.

## 5 APPLICATION OF THE THEORY: INTERRUPTIONS

Interruptions in discourses pose an important test of any theory of discourse structure. Because processing an utterance requires ascertaining how it fits with previous discourse, it is crucial to decide which parts of the previous discourse are relevant to it, and which cannot be. Interruptions, by definition, do not fit; consequently their treatment has implications for the treatment of the normal flow of discourse. Interruptions may take many forms – some are not at all relevant to the content and flow of the interrupted discourse, others are quite relevant, and many fall somewhere in between these extremes. A theory must differentiate these cases and explain (among other things) what connections exist between the main discourse and the interruption, and how the relationship between them affects the processing of the utterances in both.

The importance of distinguishing between intentional structure and attentional state is evident in the three examples considered in Subsections 5.2, 5.3, and 5.4. The distinction also permits us to explain a type of behavior deemed by others to be similar – so-called semantic returns – an issue we examine in Subsection 5.5.

These examples do not exhaust the types of interruptions that can occur in discourse. There are other ways to vary the explicit linguistic (and nonlinguistic) indicators used to indicate boundaries, the relationships between DSPs, and the combinations of focus space relationships present. However, the examples provide illustrations of interruptions at different points along the spectrum of relevancy to the main discourse. Because they can be explained more adequately by the theory of discourse structure presented here than by previous theories, they support the importance of the distinctions we have drawn.

### 5.1 PRELIMINARY DEFINITIONS

From an intuitive view, we observe that interruptions are pieces of discourse that break the flow of the preceding discourse. An interruption is in some way distinct from the rest of the preceding discourse; after the break for the interruption, the discourse returns to the interrupted piece of discourse. In the example below, from Polanyi and Scha (forthcoming), there are two (separate) discourses, D1 indicated in normal type, and D2 in italics.

D2 is an interruption that breaks the flow of D1 and is distinct from D1.

D1: John came by  
and left the groceries  
D2: *Stop that*  
*you kids*  
D1: and I put them away  
after he left

Using the theory described in previous sections, we can capture the above intuitions about the nature of interruptions with two slightly different definitions. The strong definition holds for those interruptions we classify as “true interruptions” and digressions, while the weaker form holds for those that are flashbacks. The two definitions are as follows:

**Strong definition:** An interruption is a discourse segment whose DSP is not dominated nor satisfaction-preceded by the DSP of any preceding segment.

**Weak definition:** An interruption is a discourse segment whose DSP is not dominated nor satisfaction-preceded by the DSP of the immediately preceding segment.

Neither of the above definitions includes an explicit mention of our intuition that there is a “return” to the interrupted discourse after an interruption. The return is an effect of the normal progress of a conversation. If we assume a focus space is normally popped from the focus stack if and only if a speaker has satisfied the DSP of its corresponding segment, then it naturally follows both that the focus space for the interruption will be popped after the interruption, and that the focus space for the interrupted segment will be at the top of the stack because its DSP is yet to be satisfied.

There are other kinds of discourse segments that one may want to consider in light of the interruption continuum and these definitions. Clarification dialogues (Allen 1979) and debugging explanations (Sidner 1983) are two such possibilities. Both of them, unlike the interruptions discussed here, share a DSP with their preceding segment and thus do not conform to our definition of interruption. These kinds of discourses may constitute another general class of discourse segments that, like interruptions, can be abstractly defined.

### 5.2 TYPE 1: TRUE INTERRUPTIONS

The first kind of interruption is the true interruption, which follows the strong definition of interruptions. It is exemplified by the interruption given in the previous subsection. Discourses D1 and D2 have distinct, unrelated purposes and convey different information about properties, objects, and relations. Since D2 occurs within D1, one expects the discourse structures for the two segments to be somehow embedded as well. The theory described in this paper differs from Polanyi and Scha’s (1984; and other more radically different proposals as well; e.g., Linde and Goguen 1978, Cohen 1983, Reich-

man-Adar 1984) because the “embedding” occurs *only* in the attentional structure. As shown in Figure 8, the focus space for D2 is pushed onto the stack above the focus space for D1, so that the focus space for D2 is more salient than the one for D1, until D2 is completed. The intentional structures for the two segments are distinct. There are two DP/DSP structures for the utterances in this sequence – one for those in D1 and the other for those in D2. It is not necessary to relate these two; indeed, from an intuitive point of view, they are not related.

The focusing structure for true interruptions is different from that for the normal embedding of segments, because the focusing boundary between the interrupted discourse and the interruption is impenetrable.<sup>12</sup> (This is depicted in the figure by a line with intersecting hash marks between focus spaces). The impenetrable boundary between the focus spaces prevents entities in the spaces below the boundary from being available to the spaces above it. Because the second discourse shifts attention totally to a new purpose (and may also shift the identity of the intended hearers), the speaker cannot use any

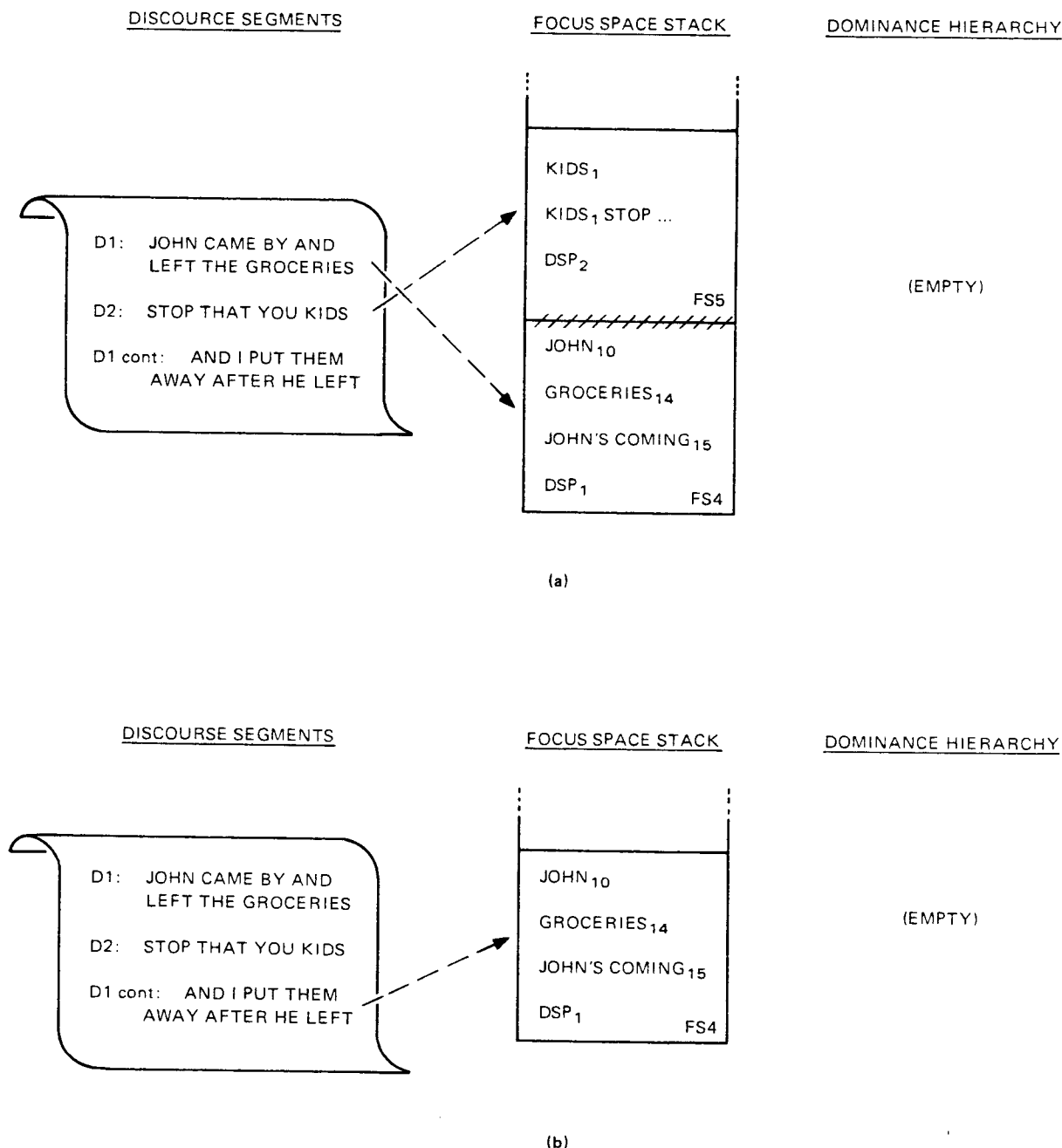


Figure 8. The structures of a true interruption.

referential expressions during it that depend on the accessibility of entities from the first discourse. Since the boundary between the focus space for D1 and the one for D2 is impenetrable, if D2 were to include an utterance such as *put them away*, the pronoun would have to refer deictically, and not anaphorically, to the groceries.

In this sample discourse, however, D1 is resumed almost immediately. The pronoun *them* in *and I put them away* cannot refer to the children (the focus space for D2 has been popped from the stack), but only to the groceries. For this to be clear to the OCP, the ICP must indicate a return to D1 explicitly. One linguistic indicator in this example is the change of mood from imperative. Indicators that the *stop that* utterance is an interruption include the change to imperative mood and the use of the vocative (Polanyi and Scha 1983). Two other indicators may be assumed to have been present at the time of the discourse – a change of intonation (imagine a slightly shrill tone of command with an undercurrent of annoyance) and a shift of gaze (toward and then away from the kids). It is also possible that the type of pause present in such cases is evidence of the interruption, but further research is needed to establish whether this is indeed the case.

In contrast to previous accounts, we are not forced to integrate these two discourses into a single grammatical structure, or to answer questions about the specific relationship between segments D2 and D1, as in Reichman's model (Reichman-Adar 1984). Instead, the intuition that readers have of an embedding in the discourse structure is captured in the attentional state by the stacking of focus spaces. In addition, a reader's intuitive impression of the distinctness of the two segments is captured in their different intentional (DP/DSP) structures.

### 5.3 TYPE 2: FLASHBACKS AND FILLING IN MISSING PLACES

Sometimes an ICP interrupts the flow of discussion because some purposes, propositions, or objects need to be brought into the discourse but have not been: the ICP forgot to include those entities first, and so must now go back and fill in the missing information. A flashback segment occurs at that point in the discourse. The flashback is defined as a segment whose DSP satisfaction-precedes the interrupted segment and is dominated by some other segment's DSP. Hence, it is a specialization of the weak definition of interruptions. This type of interruption differs from true interruptions both intentionally and linguistically: the DSP for the flashback bears some relationship to the DP for the whole discourse. The linguistic indicator of the flashback typically includes a comment about something going wrong. In addition the audience always remains the same, whereas it may change for a true interruption (as in the example of the previous section).

In the example below, taken from Sidner (1982), the ICP is instructing a mock-up system (mimicked by a person) about how to define and display certain informa-

tion in a particular knowledge-representation language. Again the interruption is indicated by italics.

OK. Now how do I say that Bill is

*Whoops I forgot about ABC.*

*I need an individual concept for the company ABC*

*...[remainder of discourse segment on ABC]...*

Now back to Bill. How do I say that Bill is an employee of ABC?

The DP for the larger discourse from which this sequence was taken is to provide information about various companies (including ABC) and their employees. The outer segment in this example –  $D_{\text{Bill}}$  – has a DSP –  $\text{DSP}_{\text{Bill}}$  – to tell about Bill, while the inner segment –  $D_{\text{ABC}}$  – has a DSP –  $\text{DSP}_{\text{ABC}}$  – to convey certain information about ABC. Because of the nature of the information being told, there is order in the final structure of the DP/DSPs: information about ABC must be conveyed before all of the information about Bill can be. The ICP in this instance does not realize this constraint until after he begins. The “flashback” interruption allows him to satisfy  $\text{DSP}_{\text{ABC}}$  while suspending satisfaction of  $\text{DSP}_{\text{Bill}}$  (which he then resumes). Hence, there is an intentional structure rooted at DP and with  $\text{DSP}_{\text{ABC}}$  and  $\text{DSP}_{\text{Bill}}$  as *ordered* sister nodes. The following three relationships hold between the different DSPs:<sup>14</sup>

DP DOM  $\text{DSP}_{\text{ABC}}$

DP DOM  $\text{DSP}_{\text{Bill}}$

$\text{DSP}_{\text{ABC}}$  SP  $\text{DSP}_{\text{Bill}}$

This kind of interruption is distinct from a true interruption because there is a connection, although indirect, between the DSPs for the two segments. Furthermore, the linguistic features of the start of the interruption signify that there is a precedence relation between these DSPs (and hence that the correction is necessary). Flashbacks are also distinct from normally embedded discourses because of the precedence relationship between the DSPs for the two segments and the order in which the segments occur.

The available linguistic data permit three possible attentional states as appropriate models for flashback-type interruptions: one is identical to the state that would ensue if the flashback segment were a normally embedded segment, the second resembles the model of a true interruption, and the third differs from the others by requiring an auxiliary stack. An example of the stack for a normally embedded sequence is given in Section 4.2

Figure 9 illustrates the last possibility. The focus space for the flashback –  $\text{FS}_{\text{ABC}}$  – is pushed onto the stack after an appropriate number of spaces, including the focus space for the outer segment –  $\text{FS}_{\text{Bill}}$ , have been popped from the main stack and pushed onto an auxiliary stack. All of the entities in the focus spaces remaining on the main stack are normally accessible for reference, but none of those on the auxiliary stack are. In the example in the figure, entities in the spaces from  $\text{FS}_A$  to  $\text{FS}_B$  are accessible as well (though less salient than) those in

space  $FS_{ABC}$ . Evidence for this kind of stack behavior could come from discourses in which phrases in the segment about ABC could refer to entities represented in  $FS_B$ , but not to those in  $FS_{Bill}$  or  $FS_C$ . After an explicit indication that there is a return to  $DSP_{Bill}$  (e.g., the *Now back to Bill* used in this example), any focus spaces left on the stack from the flashback are popped off, and all spaces on the auxiliary stack (including  $FS_{Bill}$ ) are returned to the main stack. Note, however, that this model does not preclude the possibility of a return to some space between  $FS_A$  and  $FS_B$  before popping the auxiliary stack. Whether there are discourses that include such a return and are deemed coherent is an open question.

The auxiliary stack model differs from the other two models by the references permitted and by the spaces that can be popped to. Given the initial configuration in Figure 9, if the segment with  $DSP_{ABC}$  were normally embedded,  $FS_{ABC}$  would just be added to the top of the stack. If it were a true interruption, the space would also be added to the stack, but with an impenetrable boundary between it and  $FS_{Bill}$ . In the normal stack model, entities in the spaces lower in the stack would be accessible; in the true interruption they would not. In either of these two models, however,  $FS_{Bill}$  would be the space returned to first. The auxiliary stack model is obviously more

complicated than the other two alternatives. Whether it (or some equivalent alternative) is necessary depends on facts of discourse behavior that have not yet been determined.

#### 5.4 TYPE 3: DIGRESSIONS

The third type of interruption, which we call a digression, is defined as a strong interruption that contains a reference to some entity that is salient in both the interruption and the interrupted segment. For example, if while discussing Bill's role in company ABC, one conversational participant interrupts with, *Speaking of Bill, that reminds me, he came to dinner last week*, Bill remains salient, but the DP changes. Digressions commonly begin with phrases such as *speaking of John* or *that reminds me*, although no cue phrase need be present, and *that reminds me* may also signal other stack and intention shifts.

In the processing of digressions, the discourse-level intention of the digression forms the base of a separate intentional structure, just as in the case of true interruptions. A new focus space is formed and pushed onto the stack, but it contains at least one – and possibly other – entities from the interrupted segment's focus space. Like the flashback-type interruption, the digression must usually be closed with an explicit utterance such as *getting back to ABC... or anyway*.

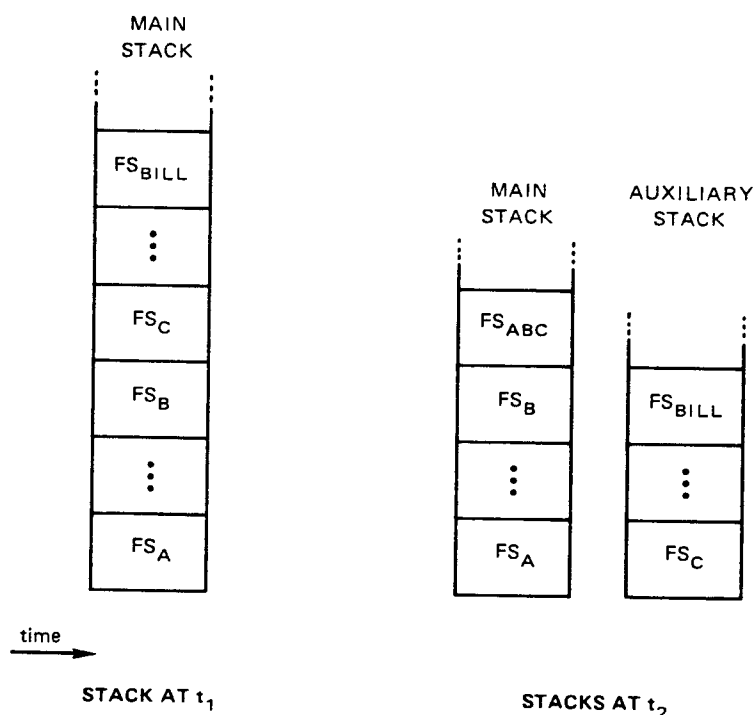


Figure 9. The auxiliary stack model for flashbacks.

### 5.5 NONINTERRUPTIONS – “SEMANTIC RETURNS”

One case of discourse behavior that we must distinguish comprises the so-called “semantic returns” observed by Reichman (1981) and discussed by Polanyi and Scha (1983). In all the interruptions we have considered so far, the stack must be popped when the interruption is over and the interrupted discourse is resumed. The focus space for the interrupted segment is “returned to.” In the case of semantic returns, entities and DSPs that were salient during a discourse in the past are taken up once again, but are explicitly reintroduced. For example, suppose that yesterday two people discussed how badly Jack was behaving at the party; then today one of them says *Remember our discussion about Jack at the party? Well, a lot of other people thought he acted just as badly as we thought he did.* The utterances today recall, or return to, yesterday’s conversation to help satisfy the intention that more be said about Jack’s poor behavior.

Anything that can be talked about once can be talked about again. However, if there is no focus space on the stack corresponding to the segment and DSP being discussed further, then, as Polanyi and Scha (1983) point out, there is no popping of the stack. There need not be any discourse underway when a semantic return occurs; in such cases, the focus stack will be empty. Thus, unlike the returns that follow normal interruptions, semantic returns involve a push onto the stack of a new space containing, among other things, representations of the reintroduced entities.

The separation of attentional state from intentional structure makes clear not only what is occurring in such cases, but also the intuitions underlying the term semantic return. In reintroducing some entities from a previous discourse, conversational participants are establishing some connection between the DSP of the new segment and the intentional structure of the original discourse. It is not a return to a previous focus space because the focus space for the original discourse is gone from the stack, and the items to be referred to must be re-established explicitly. For example, the initial reference to Jack in the preceding example cannot be accomplished with a pronoun; with no prior mention of Jack in the current discussion, one cannot say, *Remember our discussion about him at the party.* The intuitive impression of a return in the strict sense is only a return to a previous intentional structure.

## 6 APPLICATION OF THE THEORY: CUE WORDS

Both attentional state and intentional structure change during a discourse. ICPs rarely change attention by directly and explicitly referring to attentional state (e.g., using the phrase *Now let’s turn our attention to...*). Likewise, discourses only occasionally include an explicit reference to a change in purpose (e.g., with an utterance such as *Now I want to explain the theory of dynamic programming*). More typically, ICPs employ indirect

means of indicating that a change is coming and what kind of change it is. Cue phrases provide abbreviated, indirect means of indicating these changes.

In all discourse changes, the ICP must provide information that allows the OCP to determine all of the following:

1. that a change of attention is imminent;
2. whether the change returns to a previous focus space or creates a new one;
3. how the intention is related to other intentions;
4. what precedence relationships, if any, are relevant;
5. what intention is entering into focus.

Cue phrases can pack in all of this information, except for (5). In this section, we explore the predictions of our discourse structure theory about different uses of these phrases and the explanations the theory offers for their various roles.

We use the configuration of attentional state and intentional structure illustrated in Figure 10 as the starting point of our analysis. In the initial configuration, the focus space stack has a space with DSP X at the bottom and another space with DSP A at the top. The intentional structure includes the information that X dominates A. From this initial configuration, a wide variety of moves may be made. We examine several changes and the cue phrases that can indicate each of them. Because these phrases and words in isolation may ambiguously play either discourse or other functional roles, we also discuss the other uses whenever appropriate. Furthermore, cue phrases do not function unambiguously with respect to a particular discourse role. Thus for example, *first* can be used for two different moves that we discuss below.

First, consider what happens when the ICP shifts to a new DSP, B, that is dominated by A (and correspondingly by X). The dominance relationship between A and B becomes part of the intentional structure. In addition, the change in DSP results in a change in the focus stack. The focus stack models this change, which we call **new dominance**, by having new space pushed onto the stack with B as the DSP of that space (as illustrated in Figure 11). The space containing A is salient, but less so than the space with B. Cue phrase(s) to signal this case, and only this one, must communicate two pieces of information: that there is a change to some new purpose (resulting in a new focus space being created in the attentional state model rather than a return to one on the stack) and that the new purpose (DSP B) is dominated by DSP A. Typical cue phrases for this kind of change are *for example* and *to wit*, and sometimes *first* and *second*.

Cue phrases can also exhibit the existence of a satisfaction-precedence relationship. If B is to be the first in a list of DSPs dominated by A, then words such as *first* and *in the first place* can be used to communicate this fact. Later in the discourse, cue phrases such as *second*, *third*, and *finally* can be used to indicate DSPs that are dominated by A and satisfaction-preceded by B. In these cases, the focus space containing B would be popped

from the stack and the new focus space inserted above the one containing A.

There are three other kinds of discourse segments that change the intentional structure with a resulting push of new focus spaces onto the stack: the true-interruption, where B is not dominated by A; the flashback, where B satisfaction-precedes A; and the digression, where B is not dominated by A, but some entity from the focus space containing A is carried over to the new focus space.

One would expect that there might be cue phrases that would distinguish among all four of these kinds of chang-

es. Just that is so. There are cue phrases that announce one and only one kind of change. The cue phrases mentioned above for new dominance are never used for the three kinds of discourse interruption pushes. The cue phrases for true-interruptions express the intention to interrupt (e.g. *Excuse me a minute*, or *I must interrupt*) while the distinct cue phrase for flashbacks (e.g. *Oops, I forgot about ...*) indicates that something is out of order. The typical opening cue phrases of the digression mention the entity that is being carried forward (e.g. *Speaking of John ...* or *Did you hear about John?*).

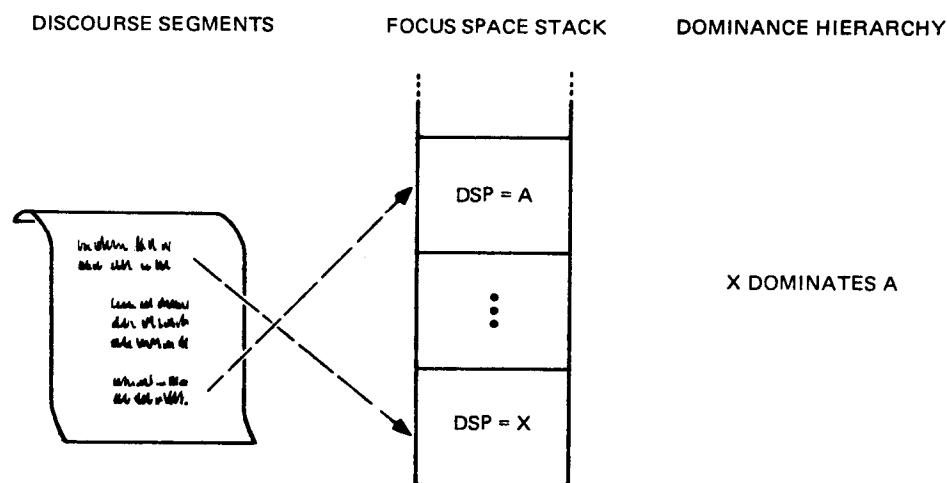


Figure 10. An initial discourse structure configuration.

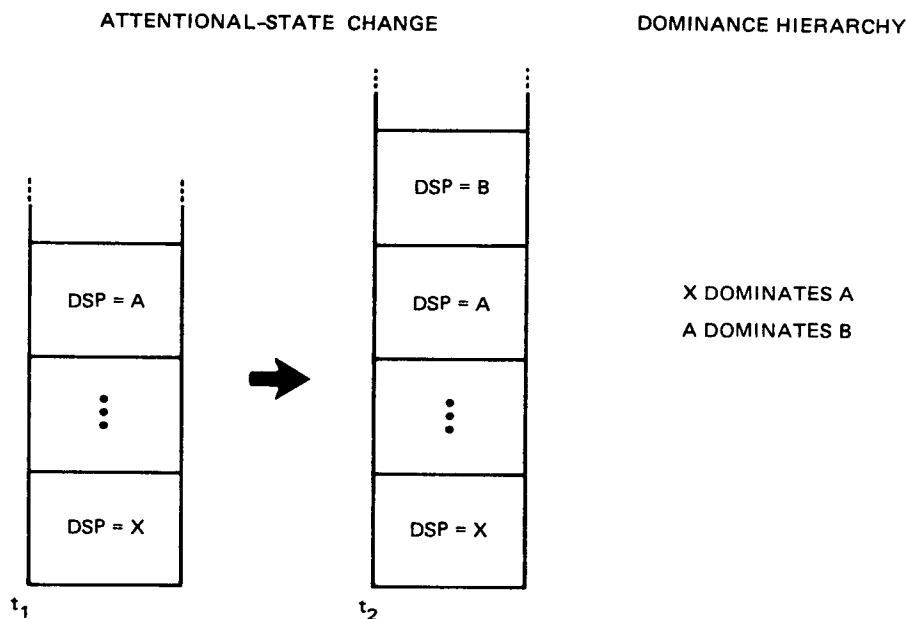


Figure 11. Attentional and intentional structures for a new subsegment.



Cue phrases can also exhibit the satisfaction of a DSP, and hence the completion of a discourse segment. The completion of a segment causes the current space to be popped from the stack. There are many means of linguistically marking completions. In texts, paragraph and chapter boundaries and explicit comments (e.g. *The End*) are common. In conversations, completion can be indicated either with cue phrases such as *fine* or *OK*<sup>15</sup> or with more explicit references to the satisfaction of the intention (e.g., *That's all for point 2*, or *The ayes have it*).

Most cue phrases that communicate changes to attentional state announce pops of the focus stack. However, at least one cue phrase can be construed to indicate a push, namely, *That reminds me*. By itself, this phrase does not specify any particular change in intentional structure, but merely shows that there will be a new DSP. Since this is equivalent to indicating that a new focus space is to be pushed onto the stack, this cue phrase is best seen as conveying attentional information.

Cue phrases that indicate pops to some other space back in the stack include *but anyway*, *anyway*, *in any case*, and *now back to...* When the current focus space is popped from the stack, a space already on the stack becomes most salient. From the configuration in Figure 10, the space with A is popped from the stack, perhaps with others, and another space on the stack becomes the top of the stack. Popping back changes the stack without creating a new DSP, or a dominance or satisfaction-precedence relationship. The pop entails a return to an old DSP; no change is effected in the intentional structure.

There are cue phrases, such as *now* and *next*, that signal a change of attentional state, but do not distinguish between the creation of a new focus space and the return to an old one. These words can be used for either move. For example, in a task-oriented discourse during which some task has been mentioned but put aside to ask a question, the use of *now* indicates a change of focus. The utterance following *now*, however, will either return the discussion to the deferred task or will introduce some new task for consideration.

Note, finally, that a pop of the focus stack may be achieved without the use of cue phrases as in the following fragment of a task-oriented dialogue (Grosz 1974):

- A: One bolt is stuck. I'm trying to use both the pliers and the wrench to get it unstuck, but I haven't had much luck.  
 E: Don't use pliers. Show me what you are doing.  
 A: I'm pointing at the bolts.  
 E: Show me the 1/2" combination wrench, please.  
 A: OK.  
 E: Good, now show me the 1/2" box wrench.  
 A: I already got it loosened.

The last utterance in this fragment returns the discourse to the discussion of the unstuck bolt. The pop can be inferred only from the content of the main portion of the utterance. The pronoun (or, more accurately, the fact that it cannot be referring to the wrench) is a cue that a pop is needed, but only the reference to the loosening action allows the OCP to recognize to which discourse segment this utterance belongs, as discussed by Sidner (1979) and Robinson (1981). A summary of the uses of cue phrases is given in Figure 12.

#### Attentional Change

- (push) now, next, that reminds me, and, but
- (pop to) anyway, but anyway, in any case, now back to
- (complete) the end, ok, fine, (paragraph break)

#### True interruption

- I must interrupt, excuse me

#### Flashbacks

- Oops, I forgot.

#### Digressions

- By the way, incidentally, speaking of,
- Did you hear about..., That reminds me

#### Satisfaction-precedes

- in the first place, first, second, finally, moreover,
- furthermore

#### New dominance

- for example, to wit, first, second, and, moreover,
- furthermore, therefore, finally

**Figure 12.** The uses of cue phrases.

The cases listed here do not exhaust the changes in focus spaces and in the dominance hierarchy that can be represented – nor have we furnished a set of rules that specify when cue phrases are necessary. Additional cases, especially special subcases of these, may be possible. When discourse is viewed in terms of intentional structure and attentional state, it is clearer just what kinds of information linguistic expressions and intonation convey to the hearer about the discourse structure. Furthermore, it is clear that linguistic expressions can function as cue phrases, as well as sentential connections; they can tell the hearer about changes in the discourse structure and be carriers of discourse, rather than sentence-level semantic, meaning.

## 7 SOME PROPERTIES AND PROBLEMS OF DISCOURSE-LEVEL INTENTIONS

The intentions that serve as DP/DSPs are natural extensions of the intentions Grice (1969) considers essential to developing a theory of utterer's meaning. There is a crucial difference, however, between our use of discourse-level intentions in this paper (and the theory, as developed so far) and Grice's use of utterance-level intentions. We are not yet addressing the issue of discourse *meaning*, but are concerned with the role of DP/DSPs in determining discourse *structure* and in specifying how these intentions can be recognized by an OCP. Although the intentional structure of a discourse plays a role in determining discourse meaning, the DP/DSPs do not in and of themselves constitute discourse segment meaning. The connection between intentional structure and discourse meaning is similar to that between attentional and cognitive states; the attentional state plays a role in a hearer's understanding of what the speaker means by a given sequence of utterances in a discourse segment, but it is not the only aspect of cognitive state that contributes to this understanding.

We will draw upon some particulars of Grice's definition of utterer's meaning to explain DSPs more fully. His initial definition is as follows:

*U* meant something by uttering *x* is true iff [for some audience *A*]:

1. *U* intended, by uttering *x*, to induce a certain response in *A*
2. *U* intended *A* to recognize, at least in part from the utterance of *x*, that *U* intended to produce that response
3. *U* intended the fulfillment of the intention mentioned in (2) to be at least in part *A*'s reason for fulfilling the intention mentioned in (1).

Grice refines this definition to address a number of counterexamples. The following portion of his final definition<sup>16</sup> is relevant to this paper:

By uttering *x* *U* meant that  $\ast\psi p$  is true iff

( $\exists A$ )( $\exists f$  [features of the utterance]) ( $\exists c$  [ways of correlating *f* with utterances<sup>17</sup>]):

(a) *U* uttered *x* intending

1. *A* to think *x* possesses *f*
2. *A* to think *f* correlated in way *c* with  $\psi$ -ing that *p*
3. *A* to think, on the basis of fulfillment of (1) and (2) that *U* intends *A* to think that *U*  $\psi$ s that *p*
4. *A* on the basis of fulfillment of (3) to think that *U*  $\psi$ s that *p*
5. and (in some cases), *A* on the basis of fulfillment of (4) himself to  $\psi$  that *p*

Grice takes  $\ast\psi p$  to be the meaning of the utterance, where  $\ast\psi$  is a mood indicator associated with the propositional attitude  $\psi$  (e.g.,  $\ast\psi$ =assert and  $\psi$ =believe). He considers attitudes like believing that ICP is a German soldier and intending to give the ICP a beer as examples of the kinds of  $\psi$ -ing that *p* that utterance intentions can embed. For expository purposes, we use the following notation to represent these utterance-level intentions:

Intend(ICP, Believe(OCP, ICP is a German soldier))  
Intend(ICP, Intend(OCP, OCP give ICP a beer))

To extend Grice's definition to discourses, we replace the utterance *x* with a discourse segment DS, the utterer *U* with the initiator of a discourse segment ICP, and the audience *A* with the OCP. To complete this extension, the following problems must be resolved:

1. specifying the discourse-level intentions and attitudes that correspond to the utterance-level intentions and  $\psi$ 's that *p*;
2. identifying the kinds of *f*s that contribute to determining discourse-level intentions;
3. identifying the modes of correlation (the *c*'s) between features of the discourse segments and types of discourse-level intentions;
4. specifying how the discourse-level intentions can be recognized by an OCP.

Although each of these issues is an unresolved problem in discourse theory, this paper has provided partial answers. The examples presented illustrate the range of discourse-level intentions; these intentions appear to be similar to utterance-level intentions in kind, but differ in that they occur in a context in which several utterances may be required to ensure their comprehension and satisfaction. The features so far identified as conveying information about DSPs are: specific linguistic markers (e.g., cue phrases, intonation), utterance-level intentions, and propositional content of the utterances. We have not explored the problem of identifying modes of correlation in any detail, but it is clear that those modes that operate at the utterance level also function at the discourse level.

As discussed previously, the proper treatment of the recognition of discourse-level intentions is especially necessary for a computationally useful account of discourse. At the discourse level, just as at the utterance level, the intended recognition of intentions plays a

central role. The DSPs are intended to be recognized: they achieve their effects, in part, because the OCP recognizes the ICP's intention for the OCP to  $\psi$  that  $p$ . The OCP's recognition of this intention is crucial to its achieving the desired effect. In Section 4 we described certain constraints on the recognition process.

### 7.1 THE BASIC GENERALIZATION

In extending Grice's analysis to the discourse level, we have to consider not only individual beliefs and intentions, but also the relationships among them that arise because of the relationships among various discourse segments (and utterances within a segment) and the purposes the segments serve with respect to the entire discourse. To clarify these relationships, consider an analogous situation with nonlinguistic actions.<sup>18</sup> An action may divide into several subactions; for example, the planting of a rose bush divides into preparing the soil, digging a hole, placing the rose bush in the hole, filling the rest of the hole with soil, and watering the ground around the bush. The intention to perform the planting action includes several subsidiary intentions (one for each of the subactions – namely, to do it).

In discourse, in a manner that is analogous to nonlinguistic actions, the DP (and some DSPs) includes several subsidiary intentions related to the DSPs it dominates. For purposes of exposition, we will use the term **primary intention** to distinguish the overall intention of the DP from the subsidiary intentions of the DP. For example in the movies argument of Section 3.1, the primary intention is for the reader to come to believe that parents and teachers should keep children from seeing too many movies; in the task dialogue of Section 3.2, the intention is that the apprentice remove the flywheel. Subsidiary intentions include, respectively, the intention that the reader believe that it is important to evaluate movies and the intention that the expert help the apprentice locate the second setscrew.

Because the beliefs and intentions of at least two different participants are involved in discourse, two properties of the general-action situation (assuming a single agent performs all actions) do not carry over. First, in a discourse, the ICP intends the OCP to recognize the ICP's beliefs about the connections among various propositions and actions. For example, in the movies argument, the reader (OCP) is intended to recognize that the author (ICP) believes some propositions provide support for others; in the task dialogue the expert (ICP) intends the apprentice (OCP) to recognize that the expert believes the performance of certain actions contributes to the performance of other actions. In contrast, in the general-action situation in which there is no communication, there is no need for recognition of another agent's beliefs about the interrelationship of various actions and intentions.

The second difference concerns the extent to which the subsidiary actions or intentions specify the overall action or intention. To perform some action, the agent

must perform each of the subactions involved; by performing all of these subactions the agent performs the action. In contrast in a discourse, the participants share the assumption of **discourse sufficiency**: it is a convention of the communicative situation that the ICP believes the discourse is sufficient to achieve the primary intention of the DP. Discourse sufficiency does not entail logical sufficiency or action completeness. It is not necessarily the case that satisfaction of all of the DSPs is sufficient in and of itself for satisfaction of the DP. Rather, there is an assumption that the information conveyed in the discourse will suffice *in conjunction with other information the ICP believes the OCP has (or can obtain)* to allow for satisfaction of the primary intention of the DP. Satisfaction of all of the DSPs, in conjunction with this additional information, is enough for satisfaction of the DP. Hence, in discourse the intentional structure (the analogue of the action hierarchy) need not be complete.

For example, the propositions expressed in the movies essay do not provide a logically sufficient proof of the claim. The author furnishes information he believes to be adequate for the reader to reach the desired conclusion and assumes the reader will supplement what is actually said with appropriate additional information and reasoning. Likewise, the task dialogue does not mention all the subtasks explicitly. Instead, the expert and apprentice discuss explicitly only those subtasks for which some instruction is needed or in connection with which some problem arises.

To be more concrete, we shall look at the extension of the Gricean analysis for two particular cases, one involving a belief, the other an intention to perform some action. We shall consider only the simplest situations, in which the primary intentions of the DP/DSPs are about either beliefs or actions, but not a mixture. Although the task dialogue obviously involves a mixture, this is an extremely complicated issue that demands additional research.

### 7.2 THE BELIEF CASE

In the belief case, the primary intention of the DP is to get the OCP to believe some proposition, say  $p$ . Each of the discourse segments is also intended to get the OCP to believe a proposition, say  $q_i$  for some  $i=1,\dots,n$  (where there are  $n$  discourse segments). In addition to the primary intention – i.e., that the OCP should come to believe  $p$  – the DP includes an intention that the OCP come to believe each of the  $q_i$  and, in addition, an intention that the OCP come to believe the  $q_i$  provide support for  $p$ . We can represent this schematically as:<sup>19</sup>

$$\begin{aligned} \forall i=1,\dots,n \text{ Intend}(\text{ICP}, \text{Believe}(\text{OCP}, p) \wedge \\ \text{Believe}(\text{OCP}, q_i) \wedge \\ \text{Believe}(\text{OCP}, \text{Supports}(p, q_1 \wedge \dots \wedge q_n))) \end{aligned}$$

There are several things to note here. To begin with, the first intention, ( $\text{Intend ICP}(\text{Believe}(\text{OCP } p))$ ), is the primary component of the DSP. Second, each of the

intended beliefs in the second conjunct corresponds to the primary component of the DSP of some embedded discourse segment. Third, the *supports* relation is not implication. The OCP is not intended to believe that the  $q_i$  imply  $p$ , but rather to believe that the  $q_i$  in conjunction with other facts and rules that the ICP assumes the OCP has available or can obtain and thus come to believe are sufficient for the OCP to conclude  $p$ . Fourth, the DP/DSP may only be completely determined at the end of the discourse (segment), as we discussed in Section 4.

Finally, to determine how the discourse segments corresponding to the  $q_i$  are related to the one corresponding to  $p$ , the OCP only has to believe that the ICP believes a supports relationship holds. Hence, for the purpose of recognizing the discourse structure, it would be sufficient for the third clause to be

... Believe(OCP, Believe(ICP,  
Supports ( $p, q_1 \wedge \dots \wedge q_n$ )))

However, the DP of a belief-case discourse is not merely to get the OCP to believe  $p$ , but to get the OCP to believe  $p$  by virtue of believing the  $q_i$ . That this is so can be seen clearly by considering situations in which the OCP already believes  $p$  and is known by the ICP to do so, but does not have a good reason for believing  $p$ . This last property of the belief case is not shared by the action case.

There is an important relationship between the supports relation and the dominance relation that can hold between DP/DSPs; it is captured in the following rule (using the same notation as above):

$$\begin{aligned} \forall i=1, \dots, n \text{ Intend}(CP_1, \text{Believe}(CP_2, p)) \wedge \\ \text{Intend}(CP_1, \text{Believe}(CP_2, q_i)) \wedge \\ \text{Believe}(CP_1, \text{Supports}(p, q_1 \wedge \dots \wedge q_n)) \iff \\ \text{DOM}(\text{Intend}(CP_1, \text{Believe}(CP_2, p)) \\ \text{Intend}(CP_1, \text{Believe}(CP_2, q_i))) \end{aligned}$$

The implication in the forward direction states that if a conversational participant ( $CP_1$ ) believes that the proposition  $p$  is supported by the proposition  $q_i$ , and he intends another participant ( $CP_2$ ) to adopt these beliefs, then his intention that  $CP_2$  believe  $p$  dominates his intention that  $CP_2$  believe  $q_i$ . Viewed intuitively,  $CP_1$ 's belief that  $q_i$  provides support for  $p$ , underlies his intention to get  $CP_2$  to believe  $p$  by getting him to believe  $q_i$ . The satisfaction of  $CP_1$ 's intention that  $CP_2$  should believe  $q_i$  will help satisfy  $CP_1$ 's intention that  $CP_2$  believe  $p$ . This relationship plays a role in the recognition of DSPs.

### 7.3 THE ACTION CASE

An analogous situation holds for a discourse segment comprising utterances intended to get the OCP to perform some set of actions directed at achieving some overall task (e.g., some segments in the task-oriented dialogue of Section 3.2). The full specification of the DP/DSP contains a *generates* relation that is derived from a relation defined by Goldman (1970). For this case, the DP/DSPs are of the following form:

$$\begin{aligned} \forall i=1, \dots, n \text{ Intend}(ICP, \text{Intend}(OCP, \text{Do}(A)) \wedge \\ \text{Intend}(OCP, \text{Do}(a_i)) \wedge \\ \text{Believe}(OCP, \text{Believe}(ICP, \\ \text{Generates}(A, a_1 \wedge \dots \wedge a_n)))) \end{aligned}$$

Each intention to act represented in the second conjunct corresponds to the primary intention of some discourse segment.

Like supports, the generates relation is partial (its partiality distinguishes it in part from Goldman's relation). Thus, the OCP is not intended to believe that the ICP believes that performance of  $a_i$  alone is sufficient for performance of  $A$ , but rather that doing all of the  $a_i$  and other actions that the OCP can be expected to know or figure out constitutes a performance of  $A$ . In the task dialogue of Section 3.2, many actions that are essential to the task (e.g., the apprentice picking up the Allen wrench and applying it correctly to the setscrews) are never even mentioned in the dialogue.

Note that it is unnecessary for the ICP or OCP to have a complete plan relating all of the  $a_i$  to  $A$  at the start of the discourse (or discourse segment). All that is required is that, for any given segment, the OCP be able to determine what intention to act the segment corresponds to and which other intentions dominate that intention. Finally, unlike the belief case, the third conjunct here requires only that the OCP recognize that the ICP believes a generates relationship holds. The OCP can do  $A$  by virtue of doing the  $a_i$  without coming himself to believe anything about the relationships between  $A$  and the  $a_i$ .

As in the belief case, there is an equivalence that links the generates relation among actions to the dominance relation between intentions. Schematically, it is as follows:

$$\begin{aligned} \forall i=1, \dots, n \text{ Intend}(CP_1, \text{Intend}(CP_2, \text{Do}(A))) \wedge \\ \text{Intend}(CP_1, \text{Intend}(CP_2, \text{Do}(a_i))) \wedge \\ \text{Believe}(CP_1, \text{Generates}(A, a_1 \wedge \dots \wedge a_n)) \iff \\ \text{DOM}(\text{Intend}(CP_1, \text{Intend}(CP_2, \text{Do}(A))) \\ \text{Intend}(CP_1, \text{Intend}(CP_2, \text{Do}(a_i)))) \end{aligned}$$

This equivalence states that, if an agent ( $CP_1$ ) believes that the performance of some action ( $a_i$ ) contributes in part to the performance of another action ( $A$ ), and if  $CP_1$  intends for  $CP_2$  to (intend to) do both of these actions, then his intention that  $CP_2$  (intend to) perform  $a_i$  is dominated by his intention that  $CP_2$  (intend to) perform  $A$ . Viewed intuitively,  $CP_1$ 's belief that doing  $a_i$  will contribute to doing  $A$  underlies his intention to get  $CP_2$  to do  $A$  by getting  $CP_2$  to do  $a_i$ . The satisfaction of  $CP_1$ 's intention for  $CP_2$  to do  $a_i$  will help satisfy  $CP_1$ 's intention for  $CP_2$  to do  $A$ .

So, for example, in the task-oriented dialogue of Section 3.2, the expert knows that using the wheelpuller is a necessary part of removing the flywheel. His intention that the apprentice intend to use the wheelpuller is thus dominated by his intention that the apprentice intend to take off the flywheel. Satisfaction of the intention to use the wheelpuller will contribute to satisfying

the intention to remove the flywheel. In general, the action  $a_i$  does not have to be a necessary action though it is in this example (at least if the task is done correctly).

A definitive statement characterizing primary and subsidiary intentions for task-oriented dialogues awaits further research not only in discourse theory, but also in the theory of intentions and actions. In particular, a clearer statement of the interactions among the intentions of the various discourse participants (with respect to both linguistic and nonlinguistic actions) awaits the formulation of a better theory of cooperation and multiagent activity.

#### 7.4 RHETORICAL RELATIONS

We are now in a position to contrast the role of DP/DSPs, supports, generates, DOM, and SP in our theory with the rhetorical relations that, according to a number of alternative theories (e.g., Grimes 1975, Hobbs 1979, Mann and Thompson 1983, Reichman-Adar 1984, McKeown 1985), are claimed to underlie discourse structure. Among the various rhetorical relations that have been investigated are elaboration, summarization, enablement, justification, and challenge. Although the theories each identify different specific relations, they all use such relations as the basis for determining discourse structure.

These rhetorical relations apply specifically to linguistic behavior and most of them implicitly incorporate intentions (e.g., the intention to summarize, the intention to justify). The intentions that typically serve as DP/DSPs in our theory are more basic than those that underlie such rhetorical relations in that they are not specialized for linguistic behavior; in many cases, their satisfaction can be realized by nonlinguistic actions as well as linguistic ones.

The supports and generates relations that must sometimes be inferred to determine domination are also more basic than rhetorical relations; they are general relations that hold between propositions and actions. Hence, the inferring of relationships such as supports and generates is simpler than that of rhetorical relationships. The determination of whether a supports or generates relationship exists depends only on facts of how the world is, not on facts of the discourse. In contrast, the recognition of rhetorical relations requires the combined use of discourse and domain information.

For several reasons, rhetorical relationships do not have a privileged status in the account given here. Although they appear to provide a metalevel description of the discourse, their role in discourse interpretation remains unclear. As regards discourse processing, it seems obvious that the ICP and OCP have essentially different access to them. In particular, the ICP may well have such rhetorical relationships "in mind" as he produces utterances (as in McKeown's (1985) system), whereas it is much less clear when (if at all) the OCP infers them. A claim of the theory being developed in this paper is that a discourse can be understood at a basic

level even if the OCP never does or can construct, let alone name, such rhetorical relationships. Furthermore, it appears that these relationships could be recast as a combination of domain-specific information, general relations between propositions and actions (e.g., supports and generates), and general relations between intentions (e.g., domination and satisfaction-precedence).<sup>20</sup> Even so, rhetorical relationships are, in all likelihood, useful to the theoretician as an analytical tool for certain aspects of discourse analysis.

## 8 CONCLUSIONS AND FUTURE RESEARCH

The theory of discourse structure presented in this paper is a generalization of theories of task-oriented dialogues. It differs from previous generalizations in that it carefully distinguishes three components of discourse structure: one linguistic, one intentional, and one attentional. This distinction provides an essential basis for explaining interruptions, cue phrases, and referring expressions.

The particular intentional structure used also differs from the analogous aspect of previous generalizations. Although, like those generalizations, it supplies the principal framework for discourse segmentation and determines structural relationships for the focusing structure (part of the attentional state), unlike its predecessors it does not depend on the special details of any single domain or type of discourse.

Although admittedly still incomplete, the theory does provide a solid basis for investigating both the structure and meaning of discourse, as well as for constructing discourse-processing systems. Several difficult research problems remain to be explored. Of these, we take the following to be of primary importance:

1. Specification of the relationship between discourse-level (DP/DSP) and utterance-level intentions;
2. Identification of the information that discourse participants use to recognize these intentions, and the ways in which they utilize it;
3. Development of an adequate treatment of the interaction among intentions of multiple participants;
4. Investigation of the effect of multiple DSPs on the theory;
5. Investigation of alternative models of attentional state.

Finally, the theory suggests several important conjectures. First, that a discourse is coherent only when its discourse purpose is shared by all the participants and when each utterance of the discourse contributes to achieving this purpose, either directly or indirectly, by contributing to the satisfaction of a discourse segment purpose. Second, general intuitions about "topic" correspond most closely to DP/DSPs, rather than to syntactic or attentional concepts. Finally, the theory suggests that the same intentional structure can give rise to different attentional structures through different discourses. The different attentional structures will be manifest in part

because different referring expressions will be valid, and in part because different cue phrases and other indicators will be necessary, optional, or redundant.

### AKNOWLEDGMENTS

We have benefited greatly from discussions with Martha Pollack, Ray Perrault, and Scott Weinstein. The paper has benefited from the comments of Jon Barwise, Marcia Derr, Brad Goodman, David Israel, Amichai Kronfeld, Mitch Marcus, Martha Pollack, Ray Perrault, John Perry, Jane Robinson, Stuart Shieber, Ralph Weischedel, Scott Weinstein, and the anonymous reviewers for **Computational Linguistics**. Whatever errors remain are, of course, all ours.

This paper was made possible by a gift from the System Development Foundation. Support was also provided for the second author by the Advanced Research Projects Agency of the Department of Defense and was monitored by ONR under Contract No. N00014-85-C-0079. The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the U.S. Government.

### REFERENCES

- Allen, J.F. 1979 A Plan-based Approach to Speech Act Recognition. Technical Report 131, Department of Computer Science, University of Toronto, Toronto, Canada.
- Allen, J.F. 1983 Recognizing Intentions from Natural Language Utterances. In Brady, M. and Berwick, R.C., Eds., *Computational Models of Discourse*. MIT Press: 107-166.
- Allen, J.F. and Perrault, C.R. 1980 Analyzing intention in dialogues. *Artificial Intelligence* 15(3): 143-178.
- Appelt, D. 1985 Planning English Referring Expressions. *Artificial Intelligence* 26: 1-33.
- Butterworth, B. 1975 Hesitation and semantic planning in speech. *Journal of Psycholinguistic Research* (4): 75-87.
- Chafe, Wallace L. 1979 The Flow of Thought and the Flow of Language. In Givon, T., Ed., *Syntax and Semantics, Vol. 12, Discourse and Syntax*. Academic Press, New York, New York: 159-182.
- Chafe, W.L. 1980 The Deployment of Consciousness in the Production of a Narrative. In Chafe, W.L., Ed., *The Pear Stories: Cognitive, Cultural and Linguistic Aspects of Narrative Production. Vol. 3. Advances in Discourse Processes*. Ablex Publishing Corp, Norwood, New Jersey: 9-50.
- Cohen, P.R. and Levesque, H.L. 1980 Speech Acts and the Recognition of Shared Plans. *Proceedings of the Third Biennial Conference of the Canadian Society for Computational Studies of Intelligence*. Victoria, British Columbia: 263-271.
- Cohen, R. 1983 A Computational Model for the Analysis of Arguments. Technical Report CSRG-151, Computer Systems Research Group, University of Toronto, Toronto, Canada.
- Cohen, P.R. and Levesque, H.J. 1985 Speech Acts and Rationality. *Proceedings of 23rd Annual Meeting of the Association for Computational Linguistics*. Chicago, Illinois: 49-60.
- Firbas, J. 1971 *On the Concept of Communicative Dynamism in the Theory of Functional Sentence Perspective*. Brno Studies in English, Vol. 7. Brno University, Brno, Czechoslovakia: 12-47.
- Goldman, A.I. 1970 *A Theory of Human Action*. Princeton University Press, Princeton, New Jersey.
- Grice, H.P. 1969 Utterer's Meaning and Intentions. *Philosophical Review* 68(2): 147-177.
- Grimes, J.E. 1975 *The Thread of Discourse*. Mouton Press, The Hague, Netherlands.
- Grosz, Barbara [Deutsch] 1974 The Structure of Task Oriented Dialogs. In *IEEE Symposium on Speech Recognition: Contributed Papers*. Carnegie Mellon University Computer Science Dept., Pittsburgh, Pennsylvania: 250-253.
- Grosz, B.J. 1977 The representation and use of focus in dialogue understanding. Technical Report 151, Artificial Intelligence Center, sri International, Menlo Park, California.
- Grosz, B.J. 1978a Discourse Analysis. In Walker, D., Ed.: 235-268.
- Grosz, B.J. 1978b Focusing in Dialog. *Theoretical Issues in Natural Language Processing-2*. University of Illinois at Urbana-Champaign, Champaign, Illinois: 96-103.
- Grosz, B.J. 1981 Focusing and Description in Natural Language Dialogues. In Joshi, A.; Webber, B.; and Sag, I., Eds., *Elements of Discourse Understanding*. Cambridge University Press, New York, New York: 84-105.
- Grosz, B.J.; Joshi, A.K.; and Weinstein, S. 1983 Providing a Unified Account of Definite Noun Phrases in Discourse. *Proceedings of the 21st Annual Meeting of the Association for Computational Linguistics*. Cambridge, Massachusetts: 44-50.
- Hajičová, E. 1983 Topic and Focus. *Theoretical Linguistics* 10(2/3): 268-276.
- Hendrix, G.G. 1979 Encoding Knowledge in Partitioned Networks. In Findler, N. V., Ed., *The Representation and Use of Knowledge in Computers*. Academic Press, New York, New York: 51-92.
- Hirschberg, J. and Pierrehumbert, J. 1986 The Intonational Structuring of Discourse. *Proceedings of the 24th Annual Meeting of the Association for Computational Linguistics*. New York, New York: 136-144.
- Hobbs, J. 1979 Coherence and Co-reference. *Cognitive Science* 3(1): 67-82.
- Holmes, H.W. and Gallagher, O. 1917 *Composition and Rhetoric*. D. Appleton and Co., New York, New York.
- Linde, C. and Goguen, J. 1978 Structure of Planning Discourse. *J. Social Biol. Struct.* 1: 219-251.
- Linde, C. 1979 Focus of Attention and the Choice of Pronouns in Discourse. In Givon, T., Ed., *Syntax and Semantics, Vol. 12, Discourse and Syntax*. Academic Press, New York, New York: 337-354.
- Litman, Diane. 1985 Plan Recognition and Discourse Analysis: An Integrated Approach for Understanding Dialogues. PhD dissertation, University of Rochester, Rochester, New York.
- Mann, W.C.; Moore, M.A.; Levin, J.A.; and Carlisle, J.H. 1975 Observation Methods for Human Dialogue. Technical Report RR/75/33, Information Sciences Institute, Marina del Rey, CA.
- Mann, W.C. and Thompson, S.A. 1983 Relational Propositions in Discourse. Technical Report RR-83-115, Information Sciences Institute, Marina del Rey, CA.
- Marcus, M.P.; Hindle, D.; and Fleck, M.M. 1983 D-Theory: Talking about Talking about Trees. *Proceedings of the 21st Annual Meeting of the Association for Computational Linguistics*. Cambridge, MA: 129-136.
- McKeown, Kathleen R. 1985 *Text Generation*. Cambridge University Press, New York, New York.
- Polanyi, L. and Scha, R. 1983 On the Recursive Structure of Discourse. In Ehlich, K. and van Riemsdijk, H., Eds., *Connectedness in Sentence, Discourse and Text*. Tilburg University, Tilburg: 141-178.
- Polanyi, L. and Scha, R. 1984 A Syntactic Approach to Discourse Semantics. *Proceedings of International Conference on Computational Linguistics*. Stanford University, Stanford, CA: 413-419.
- Polanyi, L. and Scha, R.J.H. forthcoming Discourse Syntax and Semantics. In Polanyi, L., Ed., *The Structure of Discourse*. Ablex Publishing Co., Norwood, New Jersey.
- Pollack, Martha E. 1986 Inferring Domain Plans in Question-Answering. PhD dissertation, University of Pennsylvania.
- Reichman, R. 1981 Plain-speaking: A theory and grammar of spontaneous discourse. PhD dissertation, Department of Computer Science, Harvard University. Also, BBN Report No. 4681, Bolt Beranek and Newman Inc., Cambridge, Massachusetts.

- Reichman-Adar, R. 1984 Extended Person-Machine Interface. *Artificial Intelligence* 22(2): 157-218.
- Reinhart, T. 1981 Pragmatics and Linguistics: An Analysis of Sentence Topics. *Philosophica* 27(1):53-94.
- Robinson, A. 1981 Determining Verb Phrase Referents in Dialogs. *American Journal of Computational Linguistics* 7(1): 1-16.
- Schank, R.C.; Collins, G.C.; Davis, E.; Johnson, P.N.; Lytinen, S.; and Reiser, B.J. 1982 What's the Point? *Cognitive Science* 6(3): 255-275.
- Sgall, P.; Hajičová, E.; and Benesová, E. 1973 *Topic, Focus and Generative Semantics*. Scriptor Verlag, GmbH, and Co, Kronberg, Taunus, East Germany.
- Sidner, C.L. 1979 Towards a Computational Theory of Definite Anaphora Comprehension in English Discourse. Technical Report 537, Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts.
- Sidner, C.L. 1982 Protocols of Users Manipulating Visually Presented Information with Natural Language. Technical Report 5128, Bolt Beranek and Newman Inc., Cambridge, Massachusetts.
- Sidner, C.L. 1983 What the Speaker Means: The Recognition of Speakers' Plans in Discourse. *International Journal of Computers and Mathematics*, Special Issue in Computational Linguistics 9(1): 71-82.
- Sidner, C.L. 1985 Plan Parsing for Intended Response Recognition in Discourse. *Computational Intelligence* 1(1): 1-10.
- Sidner, C.L. and Israel, D.J. 1981 Recognizing Intended Meaning and Speaker's Plans. Proceedings of the Seventh International Joint Conference in Artificial Intelligence. University of British Columbia, British Columbia, Canada: 203-208.
- Walker, D. 1978 *Understanding Spoken Language*. Elsevier North-Holland, New York, New York.
- Wittgenstein, L. 1953 *Philosophical Investigations*. Oxford University Press, London, England.
- 1978) did not use a prebuilt tree, but constructed the tree – based on a partially-ordered model – only as a given discourse evolved.
6. The observant reader will note that this was written in the early days of the cinema, before the advent of sound; hence the quotation marks around “movies.” Note also that utterance (7) contains a somewhat odd preposition, and utterance (16) somewhat odd definite noun phrases. We have quoted the text exactly as it was printed.
  7. The segmentation omits some levels of detail. For example, utterances 19-24 are a segment within DS5. Rather than present this detail, we concentrate on the larger segments here so as to focus on the major issues with which this paper is concerned.
  8. This modification “folds in” an informing action with the request. Such combining of two types of speech acts is similar to the action subsumption that Appelt (1985) discusses in regard to referring expressions.
  9. Hirschberg and Pierrehumbert (1986) have shown recently that intonational features, most notably pitch range, can also be used to indicate discourse segment boundaries.
  10. We assume here that the OCP must recognize intentions rather than actions. The argument that such is the case is beyond the scope of this paper. At a very general level, it centers on the possibility that the very same sequence of utterance actions will correspond to two different discourse structures with the difference statable only in terms of the ICP's intentions. The possibility of such sequences was suggested to us by Michael Bratman [personal communication]. The irony contained in such a clause as *you're a real sweetheart* illustrates the need to consider intentions.
  11. This knowledge may be available prior to the discourse or from information supplied by previous utterances in the discourse.
  12. This boundary is clearly atypical of stacks. It suggests that ultimately the stack model is not quite what is needed. What structure should replace the stack remains unclear to us.
  13. Because this is so clearly the case on other grounds, the segment boundary is obvious even to a reader after the fact.
  14. From just the fragment presented, all that can be determined is that the two dominates relationships are domination but not direct domination.
  15. *OK* is many ways ambiguous. It may also mean (at least) *I heard what you said, I heard and intend to do what you intend me to intend, I am done what I undertook to do, or I approve what you are about to do*.
  16. This portion is taken from Redefinition IVB: a further redefinition deals with abstracting about audience and would unnecessarily complicate our initial view of intentions and discourse.
  17. Grice (1969) mentions iconic, conventional, and associative modes, giving examples of each.
  18. This analogy is meant to help clarify and motivate the discussion. Although it also suggests some important problems in common between research on discourse and research on theories of action and intention, those issues are the subject of another paper.
  19. Here again we use a notational shorthand rather than a formal language to make some of the relationships clearer.
  20. This claim reflects a move analogous to the one made by Cohen and Levesque (1985) in showing that the definitions of various speech acts can be derived as lemmas within a general theory of rational behavior.

## NOTES

1. The use of the phrase “linguistic structure” to refer to the structure of sequences of utterances is a natural extension of its use in traditional linguistic theories to refer to the syntactic structure of individual sentences. To avoid confusion the phrase “linguistic structure” will be used in this paper only to refer to the structure of a sequence of utterances composing a discourse or discourse segment.
2. Mann has also reported that the subjects did not label segments nearly so consistently. We believe this fact is related to the kinds of relations the labels were dependent upon. As discussed in Section 7.4, there is a difference between the intentional structure we describe and the relations that others use.
3. Referring expressions can also be used to mark a discourse boundary. For example, novelists sometimes use pronouns to indicate a new scene in a story.
4. These two relations are similar to ones that play a role in parsing at the sentence level: immediate dominance and linear precedence. However, the dominance relation, like the one in Marcus and Hindle's D-theory (Marcus et al. 1983), is partial (i.e., nonimmediate).
5. Even in the task case the orderings may be partial. In fact, the systems built for task-oriented dialogues (Robinson 1981, Walker