



HAL
open science

Complex discourse units and their semantics

Nicholas Asher, Antoine Venant, Philippe Muller, Stergos Afantenos

► **To cite this version:**

Nicholas Asher, Antoine Venant, Philippe Muller, Stergos Afantenos. Complex discourse units and their semantics. CID 2011 - Constraints In Discourse, Sep 2011, Agay Roches Rouges, Var, France. hal-03690403

HAL Id: hal-03690403

<https://hal.science/hal-03690403>

Submitted on 8 Jun 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Complex discourse units and their semantics

Nicholas Asher¹, Antoine Venant², Philippe Muller^{2,3} and Stergos Afantenos²

IRIT

(1) CNRS, (2) Univ. Toulouse, (3) Alpage-INRIA

1 Introduction

A natural and intuitive principle concerning the organization of content in discourse is that discourse structure and rhetorical function operate at several levels of granularity at once. There are low level discourse connections between elementary discourse units (EDUs), even within a single sentence; but there are also discourse connections between larger constituents, complex discourse units or CDUs, which may include only two or three EDUs or may correspond to several paragraphs. CDUs and the constraints they impose on the discourse structure have not been an object of study in computational or formal work on discourse, as they are generally the by-product of processes either focused on elementary units, e.g. in RST, [17, 14] or on thematic cohesion, e.g. in text tiling—[10, 9].

The purpose of this paper is to fill this lacuna. First, we demonstrate the centrality of CDUs in an account of discourse structure. We then provide formal definitions of equivalences involving discourse graphs, which enables us to prove some results about how CDUs relate to EDUs and to each other. This in turn leads us to provide separation axioms or existence principles for CDUs. We work within the framework of SDRT [2, 4], which is a more general framework than that of RST.

2 Empirical evidence: the Annodis corpus

The present work is supported by the gathering of human annotations on discursive phenomena within the Annodis project, and the resulting Annodis corpus.¹ The Annodis corpus [18] is a French language corpus consisting of 86 documents annotated by experts in the field of discourse and about 100 documents doubly annotated by naive subjects. The documents come from the French wikipedia and the French news journal *Est Republicain*².

During the naive phase, a manual was provided to the subjects which contained information on segmentation of EDUs as well as an intuitive account of the semantics of the 18 discourse relations used. These relations are common to most theories of discourse structure, including RST and SDRT. The manual contained a brief mention of complex discourse units. The annotators were instructed to group together as many EDUs as they wanted whenever they felt that the group was semantically coherent and that there was a discourse relation between this group and another EDU (or CDU). No other structural postulates were provided to the subjects. In particular, no distinction was introduced about the difference between hierarchical relations (e.g. elaboration), called subordinating relations in SDRT, or “dominance” in the theory of [11], and other relations (called coordinating relations in SDRT, corresponding to satisfaction-precedence in [11]). Many authors postulate restrictions on discourse attachment with respect to these two types of rhetorical relations, such as what is sometimes called the Right Frontier Constraint (RFC), restricting attachment of new segments to the last introduced segment at each level of the already built discourse hierarchy.³

On the other hand, the annotators that produced the expert corpus followed these constraints. They explicitly respected SDRT’s RFC and allowed for no overlapping CDUs. They also allowed for discourse asides,

¹ Cf. <http://w3.erss.univ-tlse2.fr/textes/pagespersos/annodis/ANNODISen.html>.

² Available at <http://www.cnrtl.fr/corpus/estrepublikain/>.

³ See also Figure 1.

which are constituents attached to a constituent within a CDU but not themselves part of the CDU. In Table 1 we provide some statistics on the EDUs, CDUs, relations as well as on the sizes of the CDUs. As we can see, about a fourth of all the discourse units are CDUs and about 68% of them are small, that is they contain 5 or less EDUs.

# EDUS	3188
# CDUS	1086
# relations	3173
# CDUS with ≤ 5 EDUS	738
# CDUS with > 5	348

Table 1. Statistics on the number of EDUS, CDUS, relations and small and big CDUS

3 Some observations about uses of CDUs

CDUs provide thematic and rhetorical coherence. They may include only two or three EDUs or may correspond to one or more paragraphs. As Table 1 shows, a quarter of the discourse units participating in rhetorical relations are CDUs and 68% are small, containing 5 EDUs or less. A CDU has a semantic content, and can have an internal structure, as a Segmented Discourse Representation Structure (SDRS) in SDRT⁴, see below. But it is also a complex speech act and as such enters as a term into other discourse relations, making it part of a larger SDRS [3].

One purpose of CDUs is to exhibit the scope of a discourse relation. Consider *[John said that]_a [Bill claimed that]_b [Pat was sick]_c*.⁵ Its discourse structure is intuitively, *Attribution(a, Attribution(b, c))*, which involves a CDU as the right argument for an Attribution. This is not equivalent to the “distributed” formula, *Attribution(a, b) \wedge Attribution(b, c)*; the latter entails that John and Bill said something whereas the former only entails that John said something.

The clear distinction in SDRT between attaching to a CDU and to an EDU that might be considered its “head” distinguishes SDRT from most computational approaches to discourse parsing in the literature, which adopt tree structures for discourse. The latter deliver structures of the form *R'(c, R(a, b))*, where it’s not *a priori* clear what it means to have as a term something of the form *R(a, b)*, where *R* is a discourse relation and *a* and *b* are DUs—is it the span or the nucleus of *R* if it has only one?⁶ In any case, the argument cannot be determined just based on the type of *R* alone; it can neither always be the nucleus of *R* nor can it always be the span. Consider:

- (1) [For the last two decades,]₁ [the German central bank had a restrictive monetary policy,]₂ [because it viewed inflation as the number one problem.]₃ (discourse structure: Elaboration(1, [2,3]), Explanation(2,3).)⁷
- (2) [John worked at U.T. for two decades]₁ [He worked in the library]₂ [because he wanted to be in charge of large collections.]₃ (discourse structure: Elaboration(1,2), Explanation(2,3).)

The fact that Elaboration in (1) has scope over both subsequent segments implies that the Germans viewed inflation as the number one problem for the last two decades. The structure in (2) where the Elaboration

⁴ For a full definition, see [4].

⁵ The square brackets in the text delimit the EDUs.

⁶ See [8]. According to [16]’s Nuclearity Principle, complex segments or spans are given as the argument of a discourse relation, whenever the complex segment is made up of constituents linked by multi-nuclear discourse relations; if the span of one of the arguments consists of two DUs linked by a nuclear satellite type relation *R*, then the argument is the nucleus of *R*.

⁷ In the Annodis corpus the relation between 1 and [2,3] is called Frame, which SDRT analyzes as a species of Elaboration [20].

does not have scope over EDU 3, and hence (2) does not imply that John wanted to be in charge of large collections for two decades. In the Annodis corpus 89% of the attachments between CDUs and DUs are cases where distinguishing between the CDU and its head makes a semantic difference. As we will see below, discourse attachments to a CDU that are equivalent to an attachment involving an EDU head of the CDU are very few—only about 11% of all attachments involving a CDU fall under such equivalences, roughly 5% of all attachments in the corpus.

CDUs are needed sometimes to express a content that can't be constructed in theories that have a "right frontier" constraint on attachment. Consider a case where there is a Narration and a Result involving one event and then a CDU containing two events that happen simultaneously—e.g., [John kissed Mary]₁, [She then slapped him]₂ [and his wife did too, at the same time]₃. This example would have the Annodis style annotation— $\text{Narration}(1, [2, 3]) \wedge \text{Result}(1, [2, 3]) \wedge \text{Parallel}(2, 3)$. In terms of propositional content, we have $\text{Narration}(1, \text{Parallel}(2, 3)) \wedge \text{Result}(1, \text{Parallel}(2, 3))$. A desirable property, which we will call "Right Distributivity", (see the corresponding theorem in section 3) should entail $\text{Narration}(1, 2)$, $\text{Narration}(1, 3)$ as well as $\text{Result}(1, 2)$ and $\text{Result}(1, 3)$. But we can't build this discourse structure either in RST or SDRT; the attachment of EDU 3 to EDU 1 given an attachment via a coordinating relation of 2 to 1 is not possible without violating SDRT's right frontier constraint (or RST's). The only way to express this content is to use a CDU. CDUs are thus not eliminable from a discourse graph without a loss of expressive power or the introduction of important ambiguities.

The CDUs in the Annodis corpus have several interesting features. First of all, though the definition of SDRSs in [4] allows them, we found very few overlapping CDUs amongst the annotations done by experts in the ANNODIS corpus (those that existed involved notation errors), and, more surprisingly few amongst those done by the naive annotators. Finally, while DUs within CDUs do not attach outside of their CDU, some DUs can see into a CDU and link to their elements. These DUs that attach inside a CDU C but are not elements of C are discourse "asides", digressions from the main story line or the rhetorical purpose of the CDU. The presence of asides has a very interesting consequence: CDUs may be discontinuous in the sense that a CDU need not span a continuous sequence of words[1]. The ANNODIS corpus has 377 discourse asides, of which this is a translated example:

- (3) [Historically,]₁₁ [the black keys were covered with ebony]₁₂ [and the white keys with ivory,]₁₃ [Obviously, [since elephants are now protected,]₁₅ synthetic materials have replaced the ivory,]₁₄ [Nevertheless, ivory is still available]₁₆

This text talks about the construction of pianos; annotators annotated the text with $\text{Frame}(11, [12,13])$ with $\text{Continuation}(12,13)$ and $\text{Contrast}(11,[14,16,...])$. These segments all contribute to the piano's description. But, the aside in 15 explains why the contrast took place and is clearly not part of the description.

4 A formal definition of complex discourse units

We want to precisify the semantics of CDUs. What are they exactly? Intuitively, they are elements of the descriptive structure of the text, which include other smaller constituents. This is rather vague, and we will need a slightly more formal background in the following sections, which we introduce here.

As we use the formal framework of SDRT theory, the descriptive structures we mentioned are called SDRSs (Segmented Discourse Representation Structures).

Definition 1 (SDRS, CDU) *A SDRS can be thought of as a labeled graph with two kind of edges:*

- Each node (or label) stands for a constituent (complex or not) of the text.
- Directed labeled edges express rhetorical relations between the constituents. The label identifies the relation, e.g. Explanation, Narration.

- Directed unlabeled edges connect a complex constituent to its sub-constituents, introducing recursivity in the structure.

Hence, a **Complex Discourse Unit** or Complex Constituent is a node of the graph that has some sub-constituents identified by the second kind of edge. We may write $\alpha \in \pi$ as shortcut for α is a sub-constituent of π .

Such a structure has its semantics given in term of truth conditional content. The labels corresponding to *elementary* discourse unit (*i.e.* those that are not complex) are assigned a low level logical formula (using for instance Kamp’s DRT [12]). The semantics of a complex constituent, as one could expect, depends recursively on its sub-constituents and the relations that hold between them.

To each relation type correspond specific semantic constraints that are expressed using various logical tools. The causal relations and Conditional have their dynamic update content (which is the standard notion of dynamic semantics but generalizes standard truth conditional content) defined relative to the evaluation points singled out by the discourse context D using counterfactuals and normality conditionals [7, 5, 13]: $(w, f) \parallel \text{goal}(a, b) \parallel (w', g)$ iff $(w, f) \parallel^{\forall} a \parallel (w', g)$ and for all the best a -worlds w' relative to w' , $\exists w'', h (w', g) \parallel^{\forall} b \parallel (w'', h)$. We express similar equivalences giving update conditions with the abbreviation \equiv : $\text{Result}(a, b) \equiv^{\forall} a$ and $^{\forall} b$ and $^{\forall} a >^{\forall} b$ and $\neg^{\forall} a \square \rightarrow \neg^{\forall} b$.⁸ $\text{Explanation}(a, b) \equiv^{\forall} a$ and $^{\forall} b$ and $^{\forall} b >^{\forall} a$ and $\neg^{\forall} b \square \rightarrow \neg^{\forall} a$. Thus, if there are no context sensitive expressions in a and b , $\text{Explanation}(a, b) \equiv \text{Result}(b, a)$. The thematic relations are difficult to define precisely even at the static level [19, 20]; for *Elab* and *Frame*, we will make use of the constraint that if $\text{Elab}(a, b)$ or $\text{Frame}(a, b)$, then the main eventuality in a temporally includes the main event in b and the content of b is part of a [2].

To conclude this section, we also remind the **right frontier constraint** which, in SDRT, restricts the available attachment points of a SDRS graph to the following:

Definition 2 (Right frontier) Let S be an SDRS.

We first need to define the relation $<_S$ on the nodes of S . Let α and δ be two nodes of S ; $\alpha <_S \delta$ holds iff:

- $\alpha \in \delta$
- There is a directed edge from δ to α in S labeled by a **subordinating** relation R .

The right frontier $RF(S)$ of S contains last_S the last EDU that has been added to S , and any node δ such that $\text{last}_S <_S^* \delta$.⁹

Figure 1 presents an example text and the corresponding SDRS and its right frontier. Subordinating relations are drawn vertically, coordinating ones horizontally, and edges that link a constituent to its sub-constituents are dashed.

5 SDRSs and their entailments

To characterize CDUs, we introduce two kinds of SDRT graph entailment. Define a *permissible continuation* of G to be a graph G' in which one or more DUs are attached to accessible attachment points in G . Then:

Definition 3 (i) $G \models G'$ iff in every point of evaluation in which G is satisfied, G' is satisfied (using the usual SDRS truth definition). (ii) $G \equiv G'$ iff $G \models G'$ and $G' \models G$. (iii) $G \Vdash G'$ iff $G \equiv G'$ and to every

⁸ the *and* in our \equiv gloss of update conditions should be understood to be dynamic.

⁹ $<_S^*$ denotes the transitive closure of $<_S$

- a Max has a great evening last night.
- b He had a great meal.
- c He ate salmon,
- d He devoured lots of cheese.
- e He then won a dancing competition.

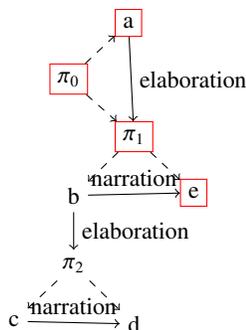


Fig. 1. SDRS example. Nodes of the right frontier are inside red rectangles

permissible continuation c of G we can assign a permissible continuation d of G' such that $G^c \equiv G'^d$, where X^c stands for the SDRS X updated with the continuation c .¹⁰ (iv) $G \cong G'$ iff $G \models G'$ and $G' \models G$.

The dynamic notion of equivalence \cong entails that if $G \cong G'$, then the right frontier of G is identical to G' . \cong is not sufficient to determine graph isomorphism. But the following theorem shows that the right frontier for all nodes in a graph does characterize many unlabelled graphs (none of the edges are labelled with discourse relations) up to isomorphism. Call an SDRS graph *proper* iff (a) if $R(a,b)$ and R' is subordinating and $R'(c,a)$ where R' is subordinating then we have $R'(c,[a,b])$ (a and b form a complex segment), (b) all attachments must be made to elements on the right frontier, (c) discourse asides always attach with subordinate relations immediately to their left unless sentential syntax dictates otherwise.¹¹ Then:

Theorem 1. For a proper, unlabelled SDRS graph G , suppose given the set of all its nodes (EDUs or CDUs if any) and for each node n in the graph, the set of its right frontier elements $rf(n)$. Then $rf(n)$ suffices to construct a proper unlabelled SDRS graph G^* that is an elementary substructure of G .¹²

A special case of Theorem 1 is that if each, constituent $n+1$ in G has only one attachment point in G_n (the unlabelled graph determined by the first n constituents), then G^* in theorem 1 is isomorphic to G . SDRT allows attachments with multiple discourse relations to one or more available nodes in the contextually given graph. Thus, not all unlabelled SDRS graphs can be recovered only from right frontier information, although all graphs that have an RST or DLTAG tree equivalent can be.

To make use of our two notions of entailment, we distinguish DRs into types. We divide relations into those that are left or right veridical using the definition of [4]. The discourse relations used in ANNODIS come in six further types: causal (explanation, goal, result), structural (contrast, parallel, continuation), point of view (attribution, comment), logical (consequence, alternation), thematic (elaboration, frame, e-elab), and narrative (narration, background, flashback). Note that $\text{explanation}(a,b) \neq \text{result}(b,a)$, even making the simplifying assumption that there are no anaphors in a or b ; \equiv does not yield definitions of dynamic update conditions. Dynamically, $\text{Explanation}(a,b)$ provides possible continuations that Result does not.

¹⁰ In the rest of the paper we make two approximations in using this definition: first we forget about the nonmonotonic character of SDRS updating. Therefore, continuation can be of three types: 1) introducing a new constituent inside an existing complex segment and attaching it to one or more elements of the right frontier 2) creating a new complex segment as second argument of a relation R by replacing $R(x,a)$ with the updated $R(x,C)$ where $a \in C$ and $b \in C$ (b is the new element to be introduced) 3) creating a new complex segment which was not explicit before, as first argument of a relation by attaching a new element to it. We do not consider here the case 3), focusing on the behaviour of complex segments as second arguments of relations (see section 6).

¹¹ (a) is a consequence of SDRT's CDP [4]; (b) and (d) are SDRT assumptions for coherent texts; (c) is an observed fact of the Annodis corpus.

¹² An elementary substructure corresponds to a subgraph, where G' is a subgraph of G iff the nodes of $G' \subseteq$ the nodes of G ; the edges of $G' \subseteq$ the edges in G , and the labelling function from edges to discourse relations for G' is a subfunction of that for G .

For the narrative relations, we have assuming no anaphors in a or b : $\text{narration}(a, b) \equiv \text{flashback}(b, a)$, but $\text{narration}(a, b) \not\equiv \text{flashback}(b, a)$.

Given the SDRT semantics for discourse relations, \cong entails only one type of strong equivalence between attachments with CDUs and attachments with their heads (5% of the attachments in our corpus).

Fact 1 For left veridical R and $R' \in \{\text{Elaboration}, \text{Frame}\}$, (i) $R'(a, b) \wedge R(a, d) \cong R(C, d)$, where C is the CDU consisting of a and b , together with $R'(a, b)$. This is represented graphically on figure 2.

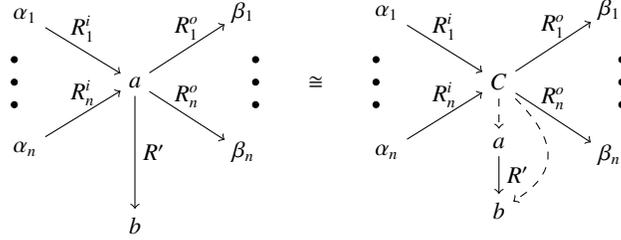


Fig. 2. Dynamic equivalence for $R' = \text{Elaboration}, \text{Frame}, \text{Attribution}$

We now characterize DRs in terms of their distributivity properties, which corrects and generalizes the Continuing Discourse Patterns (CDP) constraint of [4]:

Fact 2 Right Distribution: For left right veridical R' and for $R := \text{Frame}, \text{Elab}, \text{Conditional}, \text{Result}, \text{Goal}, \text{Attribution}, \text{or Commentary}$ ¹³, $R(a, [b, c]) \wedge R'(b, c) \models R(a, c) \wedge R(a, b)$
The result does not hold for $R := \text{Explanation}, \text{Flashback}, \text{Narration}, \text{Parallel}, \text{Alternation}$ and $\text{Contrast}, \text{E-elab}$.

Fact 3 Left Distribution: For left right veridical R' and where $R := \text{Explanation}, \text{E-Elab}$ and Attribution , $R([a, b], c) \wedge R'(a, b) \models R(a, c) \wedge R(b, c)$,
The result does not hold for $R := \text{Result}, \text{Conditional}, \text{Result}, \text{Goal}, \text{Narration}$.

We have illustrated static and dynamic equivalences with small structures. To make use of them in real-world examples, we need a way to apply them locally in more complex structures. This is the purpose of the replacement theorem we propose here, which allows under certain hypotheses (most of them simply guarantee that *replacing* is a well defined operation) to replace a subgraph by an equivalent one while keeping the whole graph equivalent to what it was at first.

First we need to define what replacing a sub-structure means. We define a sub-structure of a SDRS g to be a connected subgraph of g closed under the dominance relation (which means that if it contains a complex segment C , it contains all sub-constituents of C). A subgraph g of G may be replaced by g' only if you know what to do with the relations holding between an element that is external to g (i.e. in $G \setminus g$) and another one that is in g . Intuitively, those relations shall be identified the corresponding elements of g' . We ensure the existence of these corresponding elements with the following hypothesis:

$$\text{if } R(x, a) \text{ is a relation occurring in } G \text{ with } x \notin g \wedge a \in g \text{ then } a \in g'$$

If this hypothesis holds, the graph $G[g'/g]$ where g' replaces g is obtained following these steps:

¹³ All of these have a semantics such that their right term is closed under logical consequence.

1. Remove all nodes of g , and all relations holding between two nodes of g or a node of g and a node of $G \setminus g$.
2. Add all nodes of g' , and all relations between them.
3. For every relation $R(x, a)$ or $R(a, x)$ where $x \in G \setminus g$ and $a \in G$ that was in G before step 2), re-add this relation between x and the re-introduced a of g' (which exists by hypothesis).

This yields the following Replacement theorem:

Theorem 2. *If G is a SDRS graph, g a substructure of G , $g' \cong g$ (resp. $g' \equiv g$) and there is no edge $R(a, b)$ which is a discourse aside in G but not in g then $G[g'/g] \cong G$ (resp. $G[g'/g] \equiv G$).*

6 Separation Axiom and non-overlapping CDUs

The existence of a CDU depends on the nature of the relations that it has to the surrounding discourse context; facts 2 and 3 characterize certain relations as being distributive and others not, which specify different separation principles. We concentrate here on cases where a CDU C is a second argument of a relation R , written as $\rightarrow_R C$. Right distributive (RD) R relations for $\rightarrow_R C$ define something like SDRT's CDP as a separation principle for CDUs. Non-right distributive relations (\neg RD), however, require that DUs in a CDU contribute “equally” to the establishment of the relation on the CDU in the following sense: if one says N was ill because he drank too much and smoked too much, both the smoking and drinking contribute equally to the explanation—eliminating either will weaken or destroy the explanation relation. $C - x$ provides the largest proper subgraph of C without x ; note in addition that elaborations of constituents are treated together as one large constituent in keeping with fact 1.

Definition 4 $b \in C \leftrightarrow \exists y \in C (\text{Attaches}(b, y) \wedge \forall R \forall x R(x, C) \rightarrow \text{Role}_{R,x}(b))$.¹⁴ *For each relation R and each constituent x , $\text{Role}_{R,x}$ defined the common rhetorical role of the elements of the complex segment. Using our distinction between right distributive and non-right distributive relation, we propose the following formalisation of this role:*

– *If R is right distributive then*

$$\text{Role}_{R,x}(b) \leftrightarrow (\text{Attaches}(b, x) \rightarrow R(x, b))$$

– *Otherwise*

$$\text{Role}_{R,x}(b) \leftrightarrow \forall y \in C (R(x, C - b) \leftrightarrow R(x, C - y))$$

A final coherence axiom says: in a CDU any two DUs have to have the same rhetorical functions; that is $\forall C \forall b_1, b_2 \in C$ and for any R such that $\rightarrow_R C$, b_1 will have the appropriate property for R iff b_2 does. These license discourse asides and immediately yield our observed fact, which imposes an important constraint on any computation of CDUs.

Fact 4 *There are no complex segments C and D such that $\exists b \in C \wedge b \in D$ but $C \not\subseteq D$ and $D \not\subseteq C$*

7 Related work on CDUs

Most discourse approaches ignore distant connections between larger constituents. In the Rhetorical Structure Theory (RST, [14]), for example, rhetorical relations can hold between elementary discourse units or

¹⁴ *Attaches*(b, x) says that b attaches to x . We ignore here SDRT's non-monotonic GL, and suppose that $R(b, x)$ holds—how this is determined can be done either by machine learning or symbolic algorithms.

larger textual chunks, given that: 1) those larger textual chunks are spanned by one and the same rhetorical relation, and 2) the candidate for attachment textual spans are contiguous [15].

Another often cited discourse annotation project is the Penn Discourse Treebank (PDTB) (Miltsakaki et al., 2004) which is based on an earlier theoretical framework developed by Webber and Joshi (1998). In this project everything evolves around discourse connectors, essentially. They are considered as predicates whose arguments are abstract objects—such as events, states and propositions—realized in terms of textual phrases. Connectors can either *explicitly* appear in the text or they can be *implicit* in which case the annotators should provide the missing discourse connector. In any case, the notion of complex discourse units is completely missing from PDTB.

A final approach, which allows for complex segments is the Graphbank corpus developed by [21]. In this project complex discourse units, called "groups of relations", are encoded using the pseudo-relation *Same*.¹⁵ As in the case of RST, in Graphbank as well groups of EDUs should be *contiguous* [21, p 31], something which is in contrast with the approach that has been taken in SDRT and the Annodis corpus which has been encoded using this theory (see section 2). Another interesting point of difference between the ANNODIS corpus and the Graphbank corpus is that the former approach does not allow for overlapping complex segments while the later does allow. Interestingly, in the Graphbank annotated corpus no instance whatsoever was found of overlapping complex segments [21, p 25], despite the explicit permission to the annotators to do so.

8 Conclusions

In this paper, we set up principles with which to explore CDUs. We provided dynamic equivalences for discourse structures with and without CDUs, and we've provided principles of CDU construction. We also showed that given certain assumptions there are no overlapping CDUs, which accords with the findings of our annotators. In future work, we hope to attack the problem of automatic CDU detection, which will greatly help efforts to build complete and detailed discourse structures for texts.

References

1. Stergos Afantenos and Nicholas Asher. Testing sdr't's right frontier. In *Proceedings of COLING 2010*, pages 1–8, 2010.
2. N. Asher. *Reference to Abstract Objects in Discourse*. Kluwer Academic Publishers, 1993.
3. N. Asher. Dynamic discourse semantics for embedded speech acts. In S. Tsohatzidis, editor, *John Searle's Philosophy of Language*, pages 211–244. Cambridge University Press, 2007.
4. N. Asher and A. Lascarides. *Logics of Conversation*. Cambridge University Press, 2003.
5. Nicholas Asher and Michael Morreau. Commonsense entailment: A modal theory of nonmonotonic reasoning. In *Proceedings of the 12th International Joint Conference on Artificial Intelligence*, Sydney, Australia, August 1991.
6. I. Borisova and G. Redeker. Same and elaboration relations in the discourse graphbank. In *Proceedings of SIGDIAL 2010*, pages 63–66. ACL, 2010.
7. Myriam Bras, Anne Le Draoulec, and Nicholas Asher. Evidence for a scalar analysis of result in sdr't from a study of the french temporal connective *alors*. In *Proceedings of the SPRIK Conference, Oslo, Norway, 2006*, pages 75–79. University of Oslo, Norway, 2007.
8. Laurence Danlos. Strong generative capacity of RST, SDRT and discourse dependency DAGs. In A. Benz and P. Kohnlein, editors, *Constraints in Discourse*. Benjamins, 2007.
9. J. Eisenstein. Hierarchical text segmentation from multi-scale lexical cohesion. In *Proceedings of HLT-NAACL*, pages 353–361. ACL, 2009.
10. M. Galley, K. R. McKeown, E. Fosler-Lussier, and H. Jing. Discourse segmentation of multi-party conversation. In *Proceedings of ACL*, pages 562–569. ACL, 2003.
11. B. Grosz and C. Sidner. Attention, intentions and the structure of discourse. *Computational Linguistics*, 12:175–204, 1986.

¹⁵ The approach of using the pseudo-relation *Same* to group elementary discourse units has been criticized by [6] for inciting confusion to annotators between this relation and the relation of *Elaboration*.

12. H. Kamp. A theory of truth and semantic representation. In J. A. G. Groenendijk, T. M. V. Janssen, and M. B. J. Stokhof, editors, *Formal Methods in Study of Languages*. Mathematical Centre Tracts, Amsterdam, 1981.
13. David Lewis. *Counterfactuals*. Blackwell, 1973.
14. W. C. Mann and S. A. Thompson. Rhetorical structure theory: A framework for the analysis of texts. *International Pragmatics Association Papers in Pragmatics*, 1:79–105, 1987.
15. D. Marcu. Building up rhetorical structure trees. In *Proceedings of AAAI*, pages 1069–1074, 1996.
16. D. Marcu. The rhetorical parsing of unrestricted natural language texts. In *Proceedings of ACL and EACL*, pages 96–103, 1997.
17. D. Marcu. A formal and computational synthesis of grosz and sidner’s and mann and thompson’s theories. In *Workshop on the Levels of Representation of Discourse*, pages 101–107, Edinburgh, 1999.
18. Marie-Paule Pery-Woodley, Nicholas Asher, P. Enjalbert, Farah Benamara, Myriam Bras, Cécile Fabre, Stéphane Ferrari, Lydia-Mai Ho-Dac, Anne Le Draoulec, Yann Mathet, Philippe Muller, Laurent Prvot, Josette Rebeyrolle, Ludovic Tanguy, Marianne Vergez Couret, Laure Vieu, and Antoine Widlcher. Annodis : une approche outille de l’annotation de structures discursive (poster). In *TALN*, 2009.
19. Laurent Prévot, Laure Vieu, and Nicholas Asher. Une formalisation plus précise pour une annotation moins confuse: la relation d’élaboration d’entité. *Journal of French Language Studies*, 19(2):207–228, juillet 2009.
20. Laure Vieu, Myriam Bras, Nicholas Asher, and Michel Aurnague. Locating adverbials in discourse. *Journal of French Language Studies*, 15(2):173–193, 2005.
21. F. Wolf and E. Gibson. *Coherence in Natural Language: Data Structures and Applications*. The MIT Press, 2006.