

## Defining dependencies (and constituents)

Kim Gerdes, Sylvain Kahane

► **To cite this version:**

Kim Gerdes, Sylvain Kahane. Defining dependencies (and constituents). Intrenational Conference on Dependency linguistics (Depling 2011), Sep 2011, Barcelona, Spain. pp.17-27, 2011. <halshs-00634109>

**HAL Id: halshs-00634109**

**<https://halshs.archives-ouvertes.fr/halshs-00634109>**

Submitted on 20 Oct 2011

**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.

# Defining dependencies (and constituents)

**Kim Gerdes**

LPP

Sorbonne Nouvelle, Paris

kim@gerdes.fr

**Sylvain Kahane**

Modyco

University Paris Ouest

sylvain@kahane.fr

## Abstract

The paper proposes a mathematical method of defining dependency and constituency provided linguistic criteria to characterize the acceptable fragments of an utterance have been put forward. The method can be used to define syntactic structures of sentences, as well as discourse structures for texts or morphematic structures for words.

**Keywords:** connection graph, dependency tree, phrase structure.

## 1 Introduction

Syntacticians generally agree on the hierarchical structure of syntactic representations. Two types of structures are commonly considered: Constituent structures and dependency structures (or mixed forms of both, like headed constituent structures, sometimes even with functional labeling). However, these structures are rarely clearly defined and often purely intuition-based as we will illustrate with some examples. Even the basic assumptions concerning the underlying mathematical structure of the considered objects (ordered constituent tree, unordered dependency tree) are rarely motivated (why syntactic structures should be trees?).

In this paper, we propose a definition of syntactic structures that supersedes constituency and dependency, based on a minimal axiom: *If an utterance can be separated into two fragments, we suppose the existence of a connection between these two parts.* We will show that this assumption is sufficient for the construction of rich syntactic structures.

The notion of *connection* stems from Tesnière who says in the very beginning of his *Éléments de syntaxe structurale* that “Any word that is part of a sentence ceases to be isolated as in the dictionary. Between it and its neighbors the mind perceives **connections**, which together form the structure of the sentence.” Our axiom is less strong than Tesnière's, because we do not presuppose that the connections are formed between words only.

We will investigate the linguistic characteristics defining the notion of “fragment” and how this notion leads us to a well-defined graph-based structure, to which we can apply further conditions leading to dependency or constituent trees. We will start with a critical analysis of some definitions in the field of phrase structure and dependency based approaches (Section 2). Connection structures are defined in Section 3. They are applied to discourse, morphology, and deep syntax in Section 4. The case of surface syntax is explored in Section 5. Dependency structures are defined in Section 6 and constituent structures in Section 7.

## 2 Previous definitions

### 2.1 Defining dependency

Tesnière (1959) does not go any further in his definition of dependency and remains on a mentalist level (“the mind perceives connections”). The first formal definition of dependency stems from Lecerf (1960) and Gladkij (1966) (see also Kahane 1997) who showed that it is possible to infer a dependency tree from a constituent tree with heads (what is commonly called *phrase structure*). Further authors have tried to overcome these first definitions of constituency. Mel'čuk (1988: 130-132) proposes a definition of fragments of two words connected together. But

it is not always possible to restrict the definition to two-word fragments. Consider:

(1) *The dog slept.*

Neither *the slept* nor *dog slept* are acceptable syntactic fragments. Mel'čuk resolves the problem by connecting *slept* with the head of *the dog*, which means that his definitions of fragments and heads are mingled. Moreover Mel'čuk's definition of the head is slightly circular: "In a sentence, wordform w1 directly depends syntactically on wordform w2 if the passive [surface] valency of the phrase w1+w2 is (at least largely) determined by the passive [surface] valency of wordform w2." However, the concept of passive valency presupposes the recognition of a hierarchy, because the passive valency of a word or a fragment designates the valency towards its governor (see Section 6.1).

Garde (1977) does not restrict his definition of dependency to two-words fragments but considers more generally "significant elements" which allows him to construct the dependency between *slept* and *the dog*. However, he does not show how to reduce such a dependency between arbitrary "significant elements" to links between words. The goal of this article is to formalize and complete Garde's and Mel'čuk's definitions.

Schubert (1987:29) attempts to define dependency as "directed co-occurrence" while explicitly including co-occurrence relations between "distant words". He explains the directedness of the co-occurrence by saying that the "occurrence of certain words [the dependent] is made possible by the presence of other words," the governor. However, "form determination should not be the criterion for establishing co-occurrence lines." This adds up to lexical co-occurrences rather than syntactic dependencies. Hudson (1994) precisely proposes to keep this type of dependencies. For our part, we want to restrict connection and dependency to couples of elements which can form an acceptable text fragment in isolation (which is not the case of the radio *playing*). We do not disagree that some sort of dependency exists between radio and *playing*, but we consider this link as a lexical or semantic dependency (Mel'čuk 1988, 2011) rather than as a surface syntactic one.

## 2.2 Defining constituency

In order to evaluate the cogency of a definition of dependency based on a pre-existing definition of constituency, we have to explore how constituents are defined.

Bloomfield (1933) does not give a complete definition of syntactic constituents. His definition of the notion of *constituent* is first given in the chapter *Morphology* where he defines the morpheme. In the chapter on syntax it is said that "Syntactic constructions are constructions in which none of the immediate constituents is a bound form. [...] The actor-action construction appears in phrases like: *John ran, John fell, Bill ran, Bill fell, Our horses ran away.* [...] The one constituent (*John, Bill, our horses*) is a form of a large class, which we call *nominative expressions*; a form like *ran* or *very good* could not be used in this way. The other constituent (*ran, fell, ran away*) is a form of another large class, which we call *finite verb expressions*." Bloomfield does not give a general definition of constituents: They are only defined by the previous examples as instances of distributional classes. The largest part of the chapter is dedicated to the definition of the head of a construction. We think that in some sense Bloomfield should rather be seen as a precursor of the notions of connection (called *construction*) and dependency than as the father of constituency.

For Chomsky, a constituent exists only inside the syntactic structure of a sentence, and he never gives precise criteria of what should be considered as a constituent. In Chomsky (1986), quarreling with the behaviorist claims of Quine (1986), he refutes it as equally absurd to consider the fragmentation of *John contemplated the problem* into *John contemplated – the problem* or into *John contemp – lated the problem* instead of the "correct" *John – contemplated the problem*. No further justification for this choice is provided.

Gleason (1961:129-130) proposes criteria to define constituents (like substitution by one word, possibility to be a prosodic unit) and to build a constituent structure bottom up: "We may, as a first hypothesis, consider that each of [the words of the considered utterance] has some statable relationships to each other word. If we can describe these interrelationships completely, we will have described the syntax of the utterance in its entirety. [...] We might start by marking those pairs of words which are felt to have the closest relationship. " But he makes the following assumption without any justification: "We will also lay down the rule that each word can be marked as a member of only one such pair." Gleason then declares the method of finding the best among all the possible pairings to be "the basic problem of syntax" and he notes him-

self that his method is “haphazard” as his “methodology has not as yet been completely worked out” and lacks precise criteria. We are not far from agreeing with Gleason but we do not think that we need to choose between various satisfactory pairings. For instance, he proposes the following analysis for the NP *the old man who lives there*:

the the the the	old	man	who	lives	there
	graybeard		who	survives	
	graybeard			surviving	
	survivor				
he					

We think that other analyses like

the he	old	man	who	lives	there there
	graybeard		living		
someone			surviving		
he					

are possible, that they are not in competition, but complementary, and that both (and others) can be exploited to find the structure of this NP.

Today, the definition of 'constituent' seems no longer be a significant subject in contemporary literature in syntax. Even pedagogical books in this framework tend to skip the definition of constituency, for example Haegeman (1991) who simply states that “the words of the sentence are organized hierarchically into bigger units called phrases.”

Commonly proposed tests for constituency include the “stand-alone test”, meaning that the segment can function as an “answer” to a question, the “movement test” including clefting and topicalization, and coordinability, the latter causing the “problems” of coordination of multiple constituents, gapping, and right-node raising.

In phrase structure frameworks, constituents are nothing but a global approach for the extraction of regularities, the only goal being the description of possible orders with few rules. However, it is never actually shown that the proposed phrase structure really is the most efficient way of representing the observed utterances.

We see that the notion of constituency is either not defined at all or in an unsatisfactory way, often based on the notion of one element, the *head*, being linked to another, its *dependent*, modifying it. It is clear that the notion of dependency cannot be defined as a derived notion of constituency, as the definition of the latter presupposes head-daughter relations, making such a definition of dependency circular.

## 2.3 Intersecting analyses

An interesting result of the vagueness of the definitions of constituency is the fact that different scholars invent different criteria that allow to choose among the possible constituent structures. For example, Jespersen's lexically driven criteria select particle verbs as well as idiomatic expressions. For instance, the sentence (2) is analyzed as “S W O” where W is called a “composite verbal expression” (Jespersen 1937:16)

(2) *She [waits on] us.*

Inversely, Van Valin & Lapolla 1997:26) oppose *core* and *periphery* of every sentence and obtain another unconventional segmentation of (3).

(3) [*John ate the sandwich*] [*in the library*]

Imposing one of these various fragmentations supposes to put forward additional statements (all legitimate) based on different types of information like head-daughter relations (for X-bar approaches), idiomaticity (for Jespersen) or argument structure or information packaging (for VanValin & Lapolla) and serve merely for the elimination of unwanted fragments.

We consider the fact that we find multiple decomposition of an utterance not to be a problem. There is no reason to restrict ourselves to one particular fragmentation as it is done in phrase-structure based approaches. On the contrary, we think that the best way to compute the syntactic structure of an utterance is to consider all its possible fragmentations and this is the idea we want to explore now. Steedman (1985) was certainly one of the first linguists to develop a formal grammar that allows various groupings of words. This work and later articles by Steedman corroborated the multi-fragment approach to syntactic structure.

## 3 Fragmentation and connection

### 3.1 Fragments

We will relax the notion of syntactic constituent. We call *fragment* of an utterance any of its subparts which is a linguistically acceptable phrase with the same semantic contribution as in the initial utterance. Let us take an example :

(4) *Peter wants to read the book.*

We consider the acceptable fragments of (4) to be: *Peter*, *wants*, *to*, *read*, *the*, *book*, *Peter wants*, *wants to*, *to read*, *the book*, *Peter wants to*, *wants to read*, *read the book*, *Peter wants to read*, *to read the book*, *wants to read the book*.

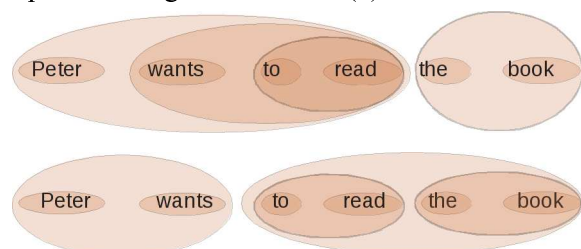
We will not justify this list of fragments at this point. We just say for the moment that *wants to*



*read*, just like *waits on*, fulfills all the commonly considered criteria of a constituent: It is a “significant element”, “functions as a unit” and can be replaced by a single word (*reads*). In the same way, *Peter wants* could be a perfect utterance. Probably the most unnatural fragment of (4) is the VP *wants to read the book*, traditionally considered as a major constituent in a phrase structure analysis.

### 3.2 Fragmentations

A *fragmentation (tree)* of an utterance *U* is a recursive partition of *U* into acceptable fragments. The following figure shows two of the various possible fragmentations of (4):



More formally, if *X* is set of minimal units (for instance the words of (4)), *fragments* are subsets of *X* and a *fragmentation* *F* is a subset of the powerset of *X* ( $F \subset P(X)$ ) such that:

1. for every  $f_1, f_2 \in F$ , either  $f_1 \subseteq f_2$ ,  $f_2 \subseteq f_1$ , or  $f_1 \cap f_2 = \emptyset$ ;
2. Each fragment is partitioned by its immediate sub-fragments.

A fragmentation whose fragments are constituents is nothing else than a constituency tree.

A fragmentation is *binary* if every fragment is partitioned into 0 or 2 fragments.

### 3.3 Connection structure and fragmentation hypergraph

We consider that each segmentation of a fragment in two pieces induces a *connection* between these two pieces.<sup>1</sup> This allows us to define graphs on the fragments of a set *X*. An *hypergraph* *H* on *X* is a triplet  $(X, F, \phi)$  where  $F \subset$

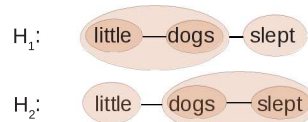
<sup>1</sup> The restriction of the connections to binary partitions can be traced back all the way to Becker (1827:469), who claims that “every organic combination within language consists of no more than two members.” (*Jede organische Zusammensetzung in der Sprache besteht aus nicht mehr als zwei Gliedern*). Although we have not encountered irreducible fragments of three or more elements in any linguistic phenomena we looked into, this cannot be *a priori* excluded. It would mean that we encountered a fragment *XYZ* where no combination of any two elements forms a fragment, i.e. is autonomizable in any without the third element. Our formal definition does not exclude this possibility at any point and a connection can in theory be, for example, ternary.

$P(X)$  and  $\phi$  is a graph on *F*. If *F* is only composed of singletons, *H* corresponds to an ordinary graph on *X*. For each binary fragmentation *F* on *X*, we will define a *fragmentation hypergraph*  $H = (X, F, \phi)$  by introducing a connection between every couple of fragments which partitions another fragment.

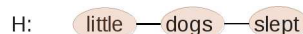
Let us illustrate this with an example:

(5) *Little dogs slept*.

There are two natural fragmentations of (5) whose corresponding hypergraphs are:<sup>2</sup>



As you can see, these two hypergraphs tell us that *little* is connected to *dogs* and *dogs* to *slept*. *H*<sub>2</sub> also show a connection between *little* and *dogs slept*, but in some sense, this is just a rough version of the connection between *little* and *dogs* in *H*<sub>1</sub>. The same observation holds for the connection between *little dogs* and *slept* in *H*<sub>1</sub>, which correspond to the connection between *dogs* and *slept* in *H*<sub>2</sub>. In other words, the two hypergraphs contains the same connections (in more or less precise versions). We can thus construct a finer-grained hypergraph *H* with the finest version of each connection:



We will call this hypergraph (which is equivalent to a graph on the words in this case) the *connection structure* of the utterance. We will now see how to define the connection structure in the general case.

### 3.4 A complete partial order on hypergraphs

We saw with our example that the connection structure is a finer-grained version of the different fragmentation hypergraphs of the utterance. So we propose to define the connection structure as the *infimum*<sup>3</sup> of the fragmentation hypergraphs for a natural order of fineness.

A connection  $f - g$  is *finer* than a connection  $f' - g'$  if  $f \subseteq f'$  and  $g \subseteq g'$ . For instance the con-

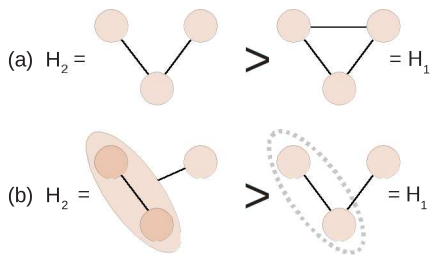
<sup>2</sup> It is possible that, for most readers, *H*<sub>1</sub> seems to be more natural than *H*<sub>2</sub>. From our point of view, it is not the case: *dogs slept* is a fragment as valid as *little dogs*. See nevertheless footnote 6.

<sup>3</sup> If  $\leq$  is a partial order on *X* and *A* is a subset of *X*, a *lower bound* of *A* is an element *b* in *X* such that  $b \leq x$  for each *x* in *A*. The *infimum* of *A*, noted  $\wedge A$ , is the greatest lower bound of *A*. A partial order for which every subset has an infimum is said to be *complete*. (As a classical example, consider the infimum for the divisibility on natural integers, which is the greatest common divisor:  $9 \wedge 12 = 3$ ).

nection  $[dogs]-[slept]$  is finer than the connection  $[little\ dogs]-[slept]$ . A connection is *minimal* when it cannot refine.

Intuitively, the *fineness order*, henceforth noted  $\leq$ , represents the precision of the hypergraph, ie.  $H_1 \leq H_2$  if  $H_1$  is a finer-grained analysis than  $H_2$ . A hypergraph  $H_1$  is *finer* than a hypergraph  $H_2$  (that is  $H_1 \leq H_2$ ) if every connection in  $H_2$  has a finer connection in  $H_1$ .

In other words,  $H_1$  must have more connections than  $H_2$ , but  $H_1$  can have some connections pointing to a smaller fragment than in  $H_2$ , and in this case the bigger fragment can be suppressed in  $H_1$  (if it carries no other connections) and  $H_1$  can have less fragments than  $H_2$ . This can be resumed by the following schemata:



In case (a),  $H_1$  is finer because it has one connection more. In case (b),  $H_1$  is finer because it has a finer-grained connection and the dotted fragment can be suppressed. It is suppressed when it carries no further connection.

We think that this partial order on hypergraphs is *complete* (see note 3). We have not proven this claim but it appears to be true on all the configurations we have investigated.

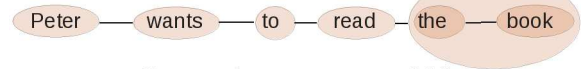
If we have an utterance  $U$  and linguistic criteria characterizing the acceptable fragments of  $U$ , we define the *connection structure* of  $U$  as the infimum of its all fragmentation hypergraphs.

### 3.5 Constructing the connection structure

Our definition could appear as being slightly complicated. In practice, it is very easy to build the connection graph of an utterance as soon as you have decided what the acceptable fragments of an utterance are. Indeed, because the fineness order on hypergraphs is complete, you can begin with any fragmentation and refine its connections until you cannot refine them any further. The connection structure is obtained when all the connections are minimal. The completeness ensures, due to the uniqueness of the greatest lower bound, that you obtain always the same structure. The only problem stems from cycles and sometimes connections must be added (see 3.7). Let us see what happens with example (4). Suppose the first step of your fragmentation is :

$f_1 = \text{Peter wants to}$   
 $f_2 = \text{read the book}$

This means that you have a connection between  $f_1$  and  $f_2$  that will correspond in the final connection structure to a link between two minimal fragments, possibly words. Now, you want to discover these minimal fragments. For that you are looking for the minimal fragment  $g$  overlapping both  $f_1$  and  $f_2$ :  $g = \text{to read}$ . It is fragmentable into *to* and *read*. Therefore the connection between  $f_1$  and  $f_2$  is finally a connection between *to* and *read*. It now remains to calculate the connection structures of  $f_1$  and  $f_2$  in order to obtain the complete connection structure of the whole sentence :



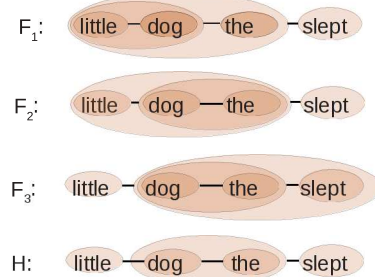
Connection structure of (4)

### 3.6 Irreducible fragment

The connection structure of (4) is not equivalent to a graph on its words because some fragments are irreducible. An *irreducible fragment* is a fragment bearing connections which cannot be attributed to one of its parts. For instance, *the book* in (4) is irreducible because there is no fragment overlapping *the book* and including only *the* or only *book* (neither *read the* nor *read book* are acceptable).

(6) *The little dog slept.*

Example (6) poses the same problem, because *little* can be connected to *dog* (*little dog* is acceptable), but *slept* must be connected to *the dog* and cannot be refined (neither *dog slept* or *the slept* is acceptable). One easily verifies that (6) has the fragmentation hypergraphs  $F_1$ ,  $F_2$ , and  $F_3$  and the connection graph  $H$  (which is their infimum). Note that the fragmentation *the dog* persists in the final connection graph  $H$  because it carries the link with *slept* but *little* is connected directly to *dog* and not to the whole fragmentation *the dog*.



Connection structure of (6):  $H = F_1 \wedge F_2 \wedge F_3$

Irreducible fragments are quite common with grammatical words. We have seen the case of determiners but conjunctions, prepositions, or relat-

ive pronouns can also cause irreducible fragments:

(7) *I think [ that [ Peter slept ] ]*

(8) *Pierre parle [ à Marie ]*

Peter speaks [to Mary]

(9) *[ the (old) man ] [ who lives ] there*

### 3.7 Cycles

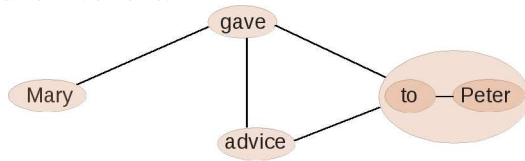
Usually the connection graph is acyclic (and could be transformed into a tree by choosing a node as the root, as we have shown for example (5). But we can have a *cycle* when a fragment XYZ can be fragmented into XY+Z, YZ+X, and XZ+Y. This can happen in examples like :

(10) *Mary gave advice to Peter.*

(11) *I saw him yesterday at school.*

(12) *the rise of nationalism in Catalonia*

In (10), *gave advice*, *gave to Peter*, and *advice to Peter* are all acceptable. We encounter a similar configuration in (11) with *saw yesterday*, *saw at school*, and *yesterday at school* (*It was yesterday at school that I saw him*). In (12), *in Catalonia* can be connected both with *nationalism* and *the rise* and there is no perceptible change of meaning. We can suppose that the hearer of these sentences constructs both connections and does not need to favor one.<sup>4</sup>



Cyclic connection graph for (10)<sup>5</sup>

### 3.8 Connection structures and fragments

We have seen that the connection structure is entirely defined from the set of fragments. Conversely the set of fragments can be reconstructed from the connection graph. Every initial fragment can be obtained by cutting connections in

<sup>4</sup> The fact that we cannot always obtain a tree structure due to irreducible fragment and cycle suggests that we could add weights on fragments indicating that a fragment (or a fragmentation) is more likely than another. We do not pursue this idea here, but we think that *weighted connection graphs* are certainly cognitively motivated linguistic representations.

Note also that the preferred fragmentation is not necessary the constituent structure. For instance, the most natural segmentation of (i) is just before the relative clause, which functions as a second assertion in this example and can be preceded by a major prosodic break (Deulofeu *et al.* 2010).

(i) *He ran into a girl, who just after entered in the shop.*

<sup>5</sup> The irreducibility of *to Peter* is conditioned by the given definition of fragments. If we considered relativization as a criteria for fragments, the possibilities of preposition stranding in English may induce the possibility to affirm that *gave* and *advice* are directly linked to the preposition *to*.

the structure and keeping the segment of the utterance corresponding to continuous pieces of the connection structure.

For instance in the connection structure of (4) cutting the connections between *to* and *read*, gives the segment *read the book*. But the segment *read the* cannot be obtained because even when cutting the connection between *the* and *book*, *read* remains connected to the entire group *the book*.

## 4 Discourse, morphology, semantics

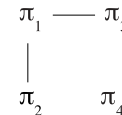
Dependency structures are usually known to describe the syntactic structures of sentences, i.e the organization of the sentence's words. In the next sections, we will give a precise definition of fragments for surface syntax in order to obtain a linguistically motivated connection structure and to transform it into a dependency tree. Let us now at first apply our methodology to construct connection structures for discourse, morphology, and the syntax-semantics interface.

### 4.1 Discourse

Nothing in our definition of connection graphs is specific to syntax. We obtain syntactic structures if we limit our maximal fragment to be sentences and our minimal fragments to be words. But if we change these constraints and begin with a whole text and take “discourse units” as minimal fragments, we obtain a discourse connection graph. This strategy can be applied to define discourse relations and discourse structures such as RST or SDRT. Of course, to obtain linguistically motivated structures, we need to define what is an acceptable sub-text of a text (generally it means to preserve coherency and cohesion).

(13) ( $\pi_1$ ) *A man walked in.* ( $\pi_2$ ) *He sported a hat.*  
( $\pi_3$ ) *Then a woman walked in.* ( $\pi_4$ ) *She wore a coat.* (Asher & Pogodalla 2010)

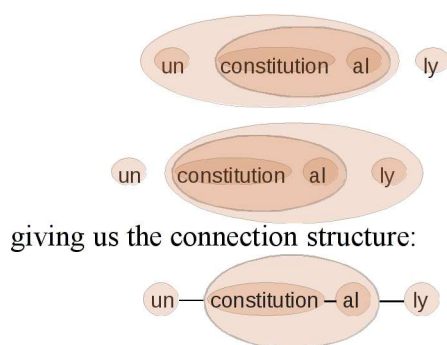
We have the fragments  $\pi_1\pi_2$ ,  $\pi_1\pi_3$ ,  $\pi_3\pi_4$  but we don't have  $\pi_2\pi_3$  nor  $\pi_1\pi_4$ . This gives us the following connection graph:



### 4.2 Morphology

On the other side, we can fragment words into morphemes. To define the acceptable fragmentations of a word, we need linguistic criteria like the commutation test. As an example for constructional morphology consider the word “*unconstitutionally*”. The two possible fragmentations are:



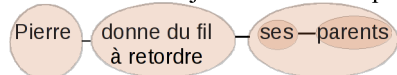


giving us the connection structure:

### 4.3 Deep Syntax

The *deep syntactic representation* is the central structure of the semantics-syntax interface (Mel'čuk 1988, Kahane 2009). If we take compositionality as a condition for fragmentation, we obtain a structure that resembles Mel'čuk's deep syntactic structure. In other words, idioms must not be fragmented and semantically empty grammatical words are not considered as fragments.

(14) *Pierre donne du fil à retordre à ses parents.*  
lit. Peter gives thread to twist to his parents  
'Peter has become a major irritant to his parents'



## 5 Fragmentations for surface syntax

### 5.1 Criteria for syntactic fragments

The connection structure we obtain completely depends on the definition of acceptable fragments. We are now interested in the linguistic criteria we need in order to obtain a connection structure corresponding to a usual surface syntactic structure. As a matter of fact, these criteria are more or less the criteria usually proposed for defining constituents. A *surface syntactic fragment* of an utterance U:

- is a subpart of U (in its original order),
- is a linguistic sign and its meaning is the same when it is taken in isolation and when it is part of U,<sup>6</sup>
- can stand alone (for example as an answer of a question),<sup>7</sup>

<sup>6</sup> This condition has to be relaxed for the analysis of idiomatic expressions as they are precisely characterized by their semantic non-compositionality. The fragments are in this case the elements that appear autonomizable in the paradigm of parallel non-idiomatic sentences.

<sup>7</sup> Mel'čuk (1988, 2011:130-132) proposes a definition of two-word fragments. Rather than the stand alone criterion, he propose that a fragment must be a prosodic unit. This is a less restrictive criterion, because the possibility to stand alone supposes to be a speech turn and therefore to be a prosodic unit. For instance *little dog* can never be a prosodic unit in *the little dog* but it is a prosodic unit when it stands

- belongs to a distributional class (and can for instance be replaced by a single word).

Mel'čuk (2006) proposes, in his definition of wordforms, to weaken the stand-alone property (or autonomizability). For instance in (6), *the* or *slept* are not autonomizable, but they can be captured by subtraction of two autonomizable fragments: *slept* = *Peter slept* \ *Peter*, *the* = *the dog* \ *dog*.<sup>8</sup> We call such fragments *weakly autonomizable*.<sup>9</sup>

Of course, even if our approach resolves most of the problems arising when trying to directly define constituents, some problems remain. For instance, if you consider the French noun phrase *le petit chien* 'the little dog', the three fragments *le chien*, *petit chien*, and *le petit* 'the little one' are acceptable. Eliminating the last fragment *le petit* supposes to put forward non trivial arguments: *le petit*, when it stands alone, is an NP (it commutes with NPs) but it cannot commute with NPs like for example *la fille* 'the girl' in *le petit chien* as *\*la fille chien* 'the girl dog' is ungrammatical. Many exciting questions posed by other phenomena like coordination or extraction cannot be investigated here for lack of space.

### 5.2 Granularity of the fragmentation

Syntactic structures can differ in the minimal units. Most of the authors consider that the wordforms are the basic units of dependency structure, but some authors propose to consider dependencies only between chunks and others between lexemes and grammatical morphemes. The following figure shows representations of various granularity for the same sentence (15).

(15) *A guy has talked to him.*

Tree A is depicting an analysis in chunks (Vergne 1990), Tree B in words, Tree D in lexemes and inflectional morphemes (and can be

alone. We think that this criterion is interesting, but not easy to use because the delimitation of prosodic units can be very controversial and seems to be a gradual notion. Note also that clitics can form prosodic units which are unacceptable fragments in our sense, like in:

(i) *the king* | *of England's* | *grandmother*

(ii) *Je crois* | *qu'hier* | *il n'est pas venu*

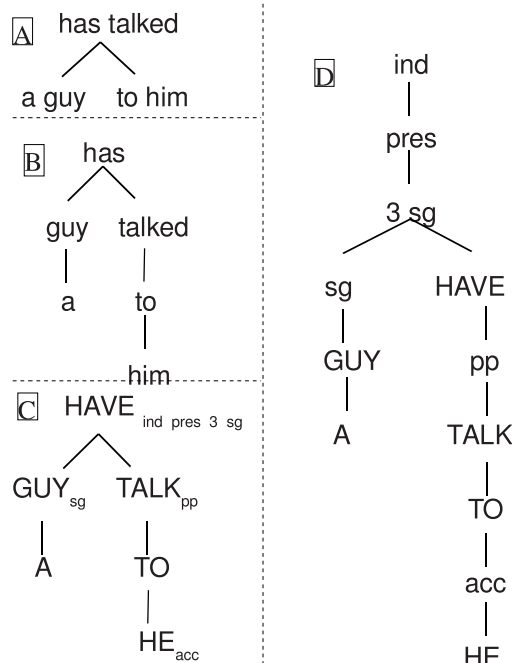
'I think | that yesterday | he didn't come'

<sup>8</sup> Note that singular bare noun like *dog* are not easily autonomizable in English, but they can for instance appear in titles.

<sup>9</sup> Some complications arise with examples like Fr. *il dormait* 'he slept'. Neither *il* (a clitic whose strong form *lui* must be used in isolation), nor *dormait* are autonomizable. But if we consider the whole distributional class of the element which can commute with *il* in this position, containing for example *Peter*, we can consider *il* to be autonomizable by generalization over the distributional class.



compared to an X-bar structure with an IP, governed by agreement and tense). The tree C (corresponding to the surface syntactic structure of Mel'čuk 1988) can be understood as an underspecified representation of D.



These various representations can be captured by our methods. The only problem is to impose appropriate criteria to define what we accept as minimal fragments. For instance, trees C and D are obtained if we accept parts of words which commute freely to be “syntactic” fragments (Kahane 2009). Conversely, we obtain tree A if we only accept strongly autonomizable fragments.

## 6 Heads and dependencies

### 6.1 Defining head and dependency

Most of the syntactic theories (if not all) suppose that the syntactic structure is hierarchized. This means that connections are directed. A directed connection is called a *dependency*. For a dependency from A to B, A is called the *governor* of B, B, the *dependent* of A, and A, the *head* of the fragment AB.<sup>10</sup> The introduction of the term “head” into syntax is commonly attributed to Henry Sweet (1891-96, I:16, sections 40 and 41): “The most general relation between words in sentences from a logical point of view is that of

<sup>10</sup>Dependency relation are sometimes called head-daughter relations in phrase structure frameworks. Note the distinction between *head* and *governor*. For a fragment  $f$ , the governor of  $f$  is necessary outside  $f$ , while the head of  $f$  is inside  $f$ . The two notions are linked by the fact that the governor  $x$  of  $f$  is the head of the upper fragment composed of the union of  $f$  and  $x$ .

**adjunct-word** and **head-word**, or, as we may also express it, of **modifier** and **modified**. [...] The distinction between adjunct-word and head-word is only a relative one : the same word may be a head-word in one sentence or context, and an adjunct-word in another, and the same word may even be a head-word and an adjunct-word at the same time. Thus in *he is very strong*, *strong* is an adjunct-word to *he*, and at the same time head-word to the adjunct-word *very*, which, again, may itself be a head-word, as in *he is not very strong*.”

Criteria for the recognition of the direction of relations between words have been proposed by Bloomfield (1933), Zwicky (1985), Garde (1977), or Mel'čuk (1988). The most common criterion is that the head of a constituent is the word controlling its distribution, which is the word that is most sensitive to a change in its context. But for any fragment, its distribution does not depend only on its head (and, as we have said in the introduction, constituents cannot easily be defined without using the notion of head). As an example, consider the fragment *little dogs* in (16):

(16) *Very little dogs slept.*

As *little* is connected to *very* and *dogs* to *slept*, *little dogs* does not have the distribution of *dogs* nor of *little* in (16) as *very dogs slept* and *very little slept* are both unacceptable. Determining the head of the fragment *little dogs* (i.e. the direction of the relation between *little* and *dogs*) is equivalent to the identification of the governor of this fragment (between *very* and *slept*). But, as soon as we have identified the governor of the fragment, the head of the fragment is simply the word of the fragment which is connected to the governor, that is the main word outside the fragment. For example, in (16), the identification of *slept* as the governor of the fragment *little dogs* also chooses *dogs* as the head of *little dogs*.

Problems occur only if we are dealing with an irreducible fragment like the determiner-noun connection.<sup>11</sup> To sum up: In order to direct the

<sup>11</sup>Various criteria have been proposed in favor of considering either the noun or the determiner as the head of this connection, in particular in the generative framework (Principles and Parameters, Chomsky (1981), remains with NP, and, starting with Abney (1986), DP is preferred). It seems that the question is triggered by the assumption that there has to be one correct directionality of this relation, in other words that the syntactic analysis is a (phrase structure) tree. This overly simple assumption leads to a debate whose theoretical implications do not reach far as any DP analysis has an isomorphic NP analysis. The NP/DP debate was triggered by the observation of a parallelism in the relation between the lexical part of a verb and its inflection (reflec-

connections and to define a dependency structure for a sentence, it is central to define the head of the whole sentence (and to resolve the case of irreducible fragments if we want a dependency tree). We consider that the head of the sentence is the main finite verb, because it bears most of the illocutionary marks: Interrogation, negation, and mood morphemes are linked to the main finite verb. In English, interrogation changes the verbal form, and in French, interrogation, negation, or mood can be marked by adding clitics or inflectional morphemes on the finite verb even if it is an auxiliary verb.

(17) a. *Did very little dogs sleep?*

b. *Pierre a-t-il dormi?*

lit. Peter has-he slept? ‘Did Peter sleep?’

c. *Pierre n'a pas dormi.*

lit. Peter neg has neg slept ‘Peter didn't sleep’

d. *Pierre aurait dormi.*

lit. Peter have-COND slept?

‘Peter would have slept’

Once the head of the sentence has been determined, most of the connections can be directed by a top down strategy. Consequently the main criterion to determine the head of a fragment  $f$  is to search if one of the words of  $f$  can form a fragment with the possible governors of  $f$ , that is if one of the words of  $f$  can be connected with the possible governors of  $f$ . If not, we are confronted with an irreducible fragment, and other criteria must be used, which we cannot discuss here (see Mel'čuk 1988, 2011).<sup>12</sup> Nevertheless, it is well known that in many cases, the head is difficult to find (Bloomfield called such configurations *exocentric*). It could be advocated not to attempt to direct the connections and to settle with an only *partially directed connection structure*.<sup>13</sup>

## 6.2 Refining the dependency structure

Even when the connection structure is completely directed, the resulting dependency structure is not necessary a tree due to irreducible fragments and cycles. We can use two principles to refine the dependency structure and to get closer to a dependency tree. The fineness order

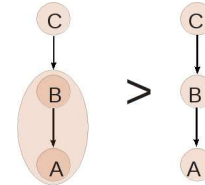
ted by the opposition between IP and VP in the generative framework). This carries over to dependency syntax: The analysis D of sentence (15) captures the intuition that the inflection steers the passive valency of a verb form.

<sup>12</sup>Conversely, whenever the fragmentation tests do not give clear results on whether or not a connection must be established, criteria used to determine the head can be helpful to confirm the validity of the connection.

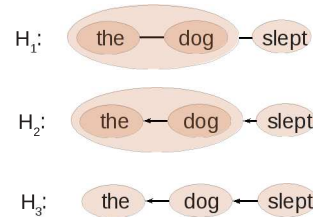
<sup>13</sup>Equally, the problem of PP attachment in parsing is certainly partially based on true ambiguities, but in many cases, it is an artificial problem of finding a tree structure where the human mind sees multiple connections, like for instance

on hypergraphs will be prolonged for directed hypergraph in accordance with these principles.

The first principle consists of avoiding double government: if  $C$  governs  $AB$  and  $B$  is the head of  $AB$ , then the dependency from  $C$  to  $AB$  can be replaced by a dependency from  $C$  to  $B$  (if  $[A \leftarrow B] \leftarrow C$ , then  $A \leftarrow B \leftarrow C$ ). In other words, the directed hypergraph with the connection  $B \leftarrow C$  is finer than the hypergraph with the connection  $[AB] \leftarrow C$ .



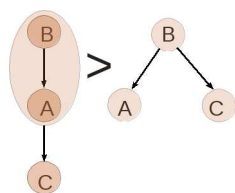
Suppose, for instance, that for the sentence (1) *The dog slept*, we obtained the connection graph  $H_1$  below. We can then add directions: The head principle easily gives the link from *slept* to the rest of the sentence, and some additional criteria may direct the connection between *dog* and *the* to give us  $H_2$ . We can now carry over this directionality to a complete dependency graph  $H_3$ .



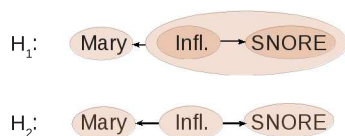
Inversely, the second principle consist of avoiding the creation of unacceptable projections: if  $C$  depends on  $AB$  and  $B$  is the head of  $AB$ , then the dependency from  $AB$  to  $C$  can be replaced by a dependency from  $B$  to  $C$  (if  $[A \leftarrow B] \rightarrow C$ , then  $A \leftarrow B \rightarrow C$ ). In other words, the directed hypergraph with the connection  $B \rightarrow C$  is finer than the hypergraph with the connection  $[AB] \rightarrow C$ .<sup>14</sup>

in *He reads a book about syntax* or in the examples (10) to (12). We can assume that a statistical parser will give better results when trained on a corpus that uses the (circular) graph structure, reserving the simple tree structures for the semantically relevant PP attachments.

<sup>14</sup> The two principles could be generalized into only one: if  $C$  is connected to  $AB$  and  $B$  is the head of  $AB$ , then the connection between  $AB$  and  $C$  can be replaced by a connection between  $B$  and  $C$  (if  $[A \leftarrow B] - C$ , then  $A \leftarrow B - C$ ). Nevertheless we think that the two principles are different and that the second one is less motivated. For instance, *the most famous of the world* can be analyzed in  $[[the\ most] \leftarrow [famous]] \rightarrow [of\ the\ world]$  and neither *famous of the world* or *the most of the world* are acceptable, but we think that  $[of\ the\ world]$  is rather selected by the superlative marker *the most* rather than by the adjective *famous* (because for any adjective  $X$  we have *the most X of the world*). The problem can be also solved by declaring *the most of the world* acceptable based on previous more general arguments.



For example, in the sentence *Mary snored*, based on the observation that the distribution of the sentence depends on the inflection of the verb, we decide to direct the relation between the inflection and the lexical part of the verb *snored* as  $\text{Infl.} \rightarrow \text{SNORE}$ . This implies, following Principle 2, that the subject depends on the inflection, and not on the lexical part of the verb. This corresponds to the observation that other, non-finite forms of the verb cannot fill the subject slot of the verbal valency.



## 7 Constituency

We saw in section 3.8 that any fragmentation can be recovered from the connection structure. As soon as the connections have been directed, some fragmentations can be favored and constituent structures can be defined.

Let us consider nodes A and B in a dependency structure. A *dominates* B if  $A = B$  or if there is a path from A to B starting with a dependency whose governor is A. The fragment of elements dominated by A is called the *maximal projection* of A. Maximal projections are major constituents (XPs in X-bar syntax). The maximal projection of A can be fragmented into  $\{A\}$  and the maximal projections of its dependents. This fragmentation gives us a flat constituent structure (with possibly discontinuous constituents).

*Partial projections* of A are obtained by considering only a part of the dependencies governed by A. By defining an order on the dependency of each node (for instance by deciding that the subject is more “external” than the object), we can privilege some partial projections and obtain our favorite binary fragmentation equivalent to the phrase structure trees we prefer. In other words, a phrase structure for a given utterance is just one of the possible fragmentations and this fragmentation can only be identified if the notion of *head* is considered.

We can thus say that phrase structure contains a definition of dependency at its very base, a fact that already appears in Bloomfield's work, who

spends much more time on defining head-daughter relations than on the notion of constituency. Jackendoff's X-bar theory is based on a head-centered definition of constituency, as each XP contains an X being the (direct or indirect) governor of the other elements of XP.

If we accept to mix criteria for identifying fragments and heads, it is possible to directly define a constituent structure without considering all the fragmentations. The strategy is recursive and top-down (beginning with the whole sentence at first constituent); each step consists of first identifying the head of the constituent we want to analyze and then looking at the biggest fragments of the utterance without its head: These biggest fragments are constituents.<sup>15</sup>

## 8 Conclusion

We have shown that it is possible to formally define a syntactic structure solely on the basis of fragmentations of an utterance. The definition of fragments does not have to keep the resulting constituent structure in mind, but can be based on simple observable criteria like different forms of autonomizability. Even (and especially) if we obtain intersecting fragmentations, we can obtain a connection graph. This operation can be applied to any type of utterance, yielding connections from the morphological to the discourse level. This delegates the search for the head of a fragment to a secondary optional operation. It is again possible to apply the known criteria for heads only when they provide clear-cut answers, leaving us with partially unresolved connections, and thus with a hypergraph, and not necessarily a tree structure. It is possible, and even frequent, that the syntactic structure is a tree, but our definition does not presuppose that it must be one. This two step definition (connection and directionality) allows for a more coherent definition of dependency as well as constituency avoiding the commonly encountered circularities. It finds *connection* as a primary notion, preliminary to constituency and dependency.

<sup>15</sup> If the head of the constituent is a finite verb, clefting can be a useful test for characterizing sub-constituents. But clefting can only capture some constituents and only if the head of the constituent has been identified and is a finite verb. As noted by Croft (2001), such constructions can only be used to characterize the constituents once we have defined them. We know that constructions like clefting select constituents because we were able to independently define constituents with other techniques. We cannot inversely define constituents by use of such language-specific constructions.

Another interesting feature of our approach is not to presuppose a segmentation of a sentence into words and even not suppose the existence of words as an indispensable notion.

In this paper we could explore neither the concrete applicability of our approach to other languages nor the interesting interaction of this new definition of dependency with recent advances in the analysis of coordination in a dependency based approach, like the notion of pile put forward in Gerdes & Kahane (2009). It also remains to be shown that the order on hypergraphs is really complete, i.e. that we can actually always compute a greatest connection graph refining any set of fragmentation hypergraphs. We also leave it to further research to explore the inclusion of weights on the connection which could replace the binary choice of presence or absence of a connection.

## Acknowledgments

We would like to thank Igor Mel'čuk, Federico Sangati, and our three anonymous reviewers.

## References

- Steven Abney. 1986. *The English Noun Phrase in its Sentential Aspect*. Unpublished Ph.D., MIT.
- Nicholas Asher, Sylvain Pogodalla. 2010. "SDRT and Continuation Semantics", *Logic and Engineering of Natural Language Semantics* 7 (LENLS VII).
- Karl Ferdinand Becker. 1841 [1827]. *Organismus der Sprache*. 2<sup>nd</sup> edition. Verlag von G.F. Kettembeil, Frankfurt am Main.
- Leonard Bloomfield. 1933. *Language*. Allen & Unwin, New York.
- Rens Bod. 1998. *Beyond grammar: an experience-based theory of language*. Stanford, CA: CSLI Publications.
- Andrew Carnie. 2011. *Modern Syntax: A Coursebook*. Cambridge University Press.
- Noam Chomsky. 1981. *Lectures On Government and Binding*. Foris, Dordrecht.
- Noam Chomsky. 1986. *New horizons in the study of language and mind*, Cambridge University Press.
- William Croft. 2001. *Radical construction grammar: syntactic theory in typological perspective*. Oxford University Press.
- José Deulofeu, Lucie Dufort, Kim Gerdes, Sylvain Kahane, Paola Pietrandrea. 2010. "Depends on what the French say", *The Fourth Linguistic Annotation Workshop (LAW IV)*.
- Paul Garde. 1977. "Ordre linéaire et dépendance syntaxique : contribution à une typologie", *Bull. Soc. Ling. Paris*, 72:1, 1-26.
- Aleksej V. Gladkij. 1966. *Leckii po matematicheskoj lingvistike dlja studentov NGU*, Novosibirsk (French translation: *Leçons de linguistique mathématique*, fasc. 1, 1970, Paris, Dunod)
- Henry A. Gleason. 1955. *An Introduction to Descriptive Linguistics*. New York: Holt, Rinehart & Winston, 503 pp. Revised edition 1961.
- Kim Gerdes, Sylvain Kahane. 2009. "Speaking in piles: Paradigmatic annotation of a French spoken corpus", *Corpus Linguistics 2009*, Liverpool.
- Otto Jespersen. 1937. *Analytic syntax*. Copenhagen.
- Yves Lecerf. 1960. "Programme des conflits, module des conflits", *Bulletin bimestriel de l'ATALA*, 4,5.
- Liliane M. V. Haegeman. 1991. *Introduction to Government and Binding Theory*. Blackwell Publishers
- Richard Hudson. 1994. "Discontinuous phrases in dependency grammars", *UCL Working Papers in Linguistics*, 6.
- Sylvain Kahane. 1997. "Bubble trees and syntactic representations", *MOL5*, Saarbrücken, 70-76.
- Sylvain Kahane. 2009. "Defining the Deep Syntactic Structure: How the signifying units combine", *MTT 2009*, Montreal.
- Igor Mel'čuk. 1988. *Dependency Syntax: Theory and Practice*. The SUNY Press, Albany, N.Y.
- Igor Mel'čuk. 2006. *Aspects of the Theory of Morphology*. de Gruyter, Berlin, New York.
- Igor Mel'čuk. 2011. "Dependency in language", *Proceedings of Dependency Linguistics 2011*, Barcelona.
- Klaus Schubert. 1987. *Metataxis: Contrastive dependency syntax for machine translation*. <http://www.mt-archive.info/Schubert-1987.pdf>
- Henry Sweet. 1891-1896. *A New English Grammar*, 2 vols. Clarendon Press. Oxford.
- Willard Quine. 1986. "Reply to Gilbert H. Harman." In E. Hahn and P.A. Schilpp, eds., *The Philosophy of W.V. Quine*. La Salle, Open Court.
- Mark Steedman. 1985. "Dependency and coordination in the grammar of Dutch and English", *Language*, 61:3, 525-568.
- Lucien Tesnière. 1959. *Éléments de syntaxe structurale*. Klincksieck, Paris.
- Arnold M. Zwicky. 1985. "Heads", *Journal of Linguistics*, 21: 1-29.