

Chapter 21

FROM TREEBANK RESOURCES TO LFG F-STRUCTURES

Automatic F-Structure Annotation of Treebank Trees and CFGs Extracted from Treebanks

Anette Frank*

German Research Center for Artificial Intelligence (DFKI), Germany
frank@dfki.de

Louisa Sadler

University of Essex
louisa@essex.ac.uk

Josef van Genabith, Andy Way

Dublin City University
{josef;away}@computing.dcu.ie

Abstract We present two companion methods for automatically enriching phrase-structure oriented treebank resources with functional structures. Both methods define systematic patterns of correspondence between partial PS configurations and functional structures. These are applied to PS rules extracted from treebanks, or to flat term representations of treebank trees.

Keywords: Automatic Annotation, Higher-level Syntax, LFG f-structures, Corpus Linguistics, Robustness, SUSANNE Corpus, AP Treebank

*The work presented here was performed while the first author was at Xerox Research Centre Europe (XRCE), Grenoble

1. INTRODUCTION

In this contribution we address two important concerns: automatic annotation of treebanks and CFGs extracted from such treebanks with LFG f(eature)-structures (Kaplan and Bresnan 1982; Bresnan 2001; Dalrymple 2001).

Treebanks which encode higher-level functional or basic predicate–argument structure, in addition to pure phrase structure information, are required as training resources for probabilistic unification grammars and data-driven parsing approaches, (e.g. Bod and Kaplan 1998). Manual construction of treebanks with feature structure annotations is very labour and cost intensive. So is the development of new or the scaling-up of existing unification grammars which can be used to analyse large text corpora. What is more, even if a large-coverage unification grammar is available, typically, for each input string it would generate hundreds or thousands of candidate (constituent and feature structure) analyses from which a highly trained expert has to select. Although proposals have been made for filtering and ranking parsing ambiguities (e.g. Charniak 1993; Abney 1997; Frank et al. 2000), to date none is guaranteed to uniquely determine the best analysis. In order not to compromise the quality of the corpus under construction, a linguistic expert is required to find the best among a large number of candidate analyses.

Given this situation, is there a way to automate, or bootstrap, the construction of grammars and treebanks with feature structure annotations reusing existing resources?

In a number of papers van Genabith et al. (1999a,b,c) presented a new corpus based method. Their basic idea is the following: take an existing treebank, read off the CFG following (Charniak 1996), *manually* annotate the extracted CFG rules with f-structure annotations and provide macros for the lexical entries. Then deterministically “rematch” the structure of the original treebank trees (not the strings) with the annotated rules. During this rematching process, the f-structure annotations are resolved, and an f-structure is produced. The entire process is deterministic if the feature structure annotations are, and to a considerable extent, costly manual inspection of candidate analyses is avoided. The method is an improvement but still involves a large labour intensive component, namely *manual* annotation of the extracted grammar rules.

Treebank grammars (CFGs extracted from treebanks) are large and grow with the size of the treebank (Charniak 1996; Krotov et al. 1998). They feature rather flat rules, many of which share and/or repeat significant portions of their right-hand sides (RHS). This causes problems for manual rule annotation approaches such as (van Genabith et al. 1999a,b,c). Manual rule annotation is labour intensive, error prone, repetitive and risks missing generalisations.

In this paper we show how f-structure annotation of both grammar rules and treebank trees can (to a large extent) be *automated*.

The basic idea is simple: functional annotations define systematic correspondences between constituent and higher level feature structure representations. These can be captured in general annotation principles, which are applied either to grammar rules extracted from a treebank or directly to treebank trees.

The observation that constituent and higher-level feature structure representations stand in a systematic relationship informs theoretical work in LFG (Kaplan and Bresnan 1982) and HPSG (Pollard and Sag 1994). In LFG c(onstituent)-structure and f-structure are independent levels of representation which are related in terms of a correspondence function ϕ . The correspondence follows linguistically determined principles which are partly universal, and partly language specific (Bresnan 2001; Dalrymple 2001).

What is new in our approach is that (i) we employ *partial* and *underspecified* annotation principles in a principle-based c- to f-structure interface for the LFG architecture; (ii) we use these to automate functional structure assignment to flat and “noisy” treebank trees and CFGs extracted from them; and (iii) we reuse existing linguistic resources. In contrast to more theoretically informed work in LFG and HPSG, treebanks do not tend to follow highly abstract and general X-bar architectural design principles. The challenge in our approach is to develop grammars and annotation principles for real text.

The potential benefits of automation are considerable: substantial reduction in development effort, hence savings in time and cost for treebank annotation and grammar development; the ability to tackle larger fragments in a shorter time, a considerable amount of flexibility for switching between different treebank annotation schemes, and a natural approach to robustness. Our methods can also be viewed as a new corpus- and data-driven approach to grammar development, an approach that as much as possible recycles existing resources.

In our work to date we have developed two related but interestingly different methods. Both methods define annotation principles as correspondences between *partial* and *underspecified* c- and f-structure configurations. In one approach (Sadler et al. 2000) we read off a CFG treebank grammar following the method of Charniak (1996) and then compile regular expression based annotation principles over the extracted grammar. In the companion approach (Frank 2000) we operate on treebank trees encoded as flat term representations and annotate them with f-structures.

Both methods are partial and robust in the following further sense: they yield partial, unconnected f-structures in the case of missing annotation principles. In the case of conflicting feature assignments (Frank 2000) admits partially unresolved f-structures to achieve further robustness.

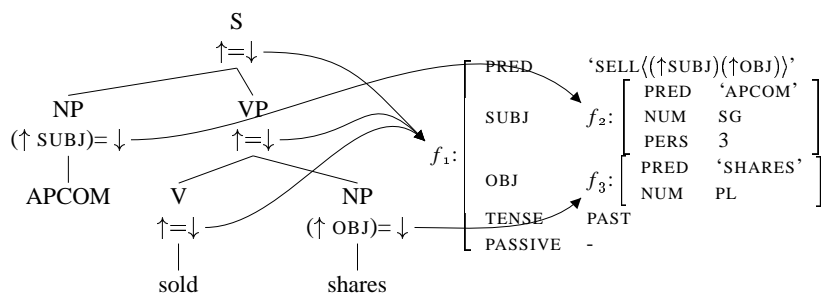
We describe two experiments, one for each method. For the first experiment we developed a regular expression based annotation principle interpreter which operates on grammar rules with *order independent* and monotonic inter-

pretation of annotation principles. For the second experiment we employed an existing term rewriting system (Kay 1999; Frank 1999), which we use to apply annotation principles to flat, term-based representations of treebank trees. The term rewriting system allows us to exploit both *order dependent*, cascaded and *order independent* formulations of annotation principles. In our first experiment we used the first 100 trees of the AP treebank (Leech and Garside 1991), in the second, 166 trees of the Susanne treebank (Sampson 1993).

The paper is structured as follows: in Section 2 we motivate and describe our annotation methods in more detail. In Section 3 we report on our two experiments. For each experiment we explain the design, describe the data and evaluate the results. In Section 4 we compare the two methods and outline ongoing research. Section 5 concludes.

2. METHODS FOR AUTOMATIC F-STRUCTURE ANNOTATION

In LFG the correspondence between c- (constituent) and f- (functional) structure is defined in terms of functional annotations of the RHS categories in CFG rules and lexical information.



PS rules define f-structure via functional descriptions

$$S \rightarrow \begin{array}{cc} \text{NP} & \text{VP} \\ (\uparrow \text{SUBJ}) = \downarrow & \uparrow = \downarrow \end{array} \quad \text{VP} \rightarrow \begin{array}{cc} \text{V} & \text{NP} \\ \uparrow = \downarrow & (\uparrow \text{OBJ}) = \downarrow \end{array}$$

$$\begin{array}{l} \text{APCOM: N} \quad (\uparrow \text{PRED}) = \text{'APCOM'} \quad \text{sold: V} \quad (\uparrow \text{PRED}) = \text{'SELL'((↑SUBJ)(↑OBJ))'} \\ (\uparrow \text{NUM}) = \text{SG} \quad (\uparrow \text{TENSE}) = \text{PAST} \\ (\uparrow \text{PERS}) = 3 \quad (\uparrow \text{PASSIVE}) = - \end{array}$$

The c-structure/f-structure correspondence follows universal and language specific principles. In our work, we define annotation principles as involving *partial* and *underspecified* phrase structure configurations and apply them to CFG rules or tree fragments that meet the relevant *partial* configuration. To illustrate the idea: a head principle assigns $\uparrow = \downarrow$ to the X daughter in all $XP \rightarrow \dots X \dots$ configurations, irrespective of the surrounding categorial context. For the example at hand, the challenge in our approach is to provide annotation principles that identify heads in the flat treebank tree and rule configurations

which deviate significantly from X-bar design principles. Annotation principles capture generalisations and can be used to *automatically* annotate PS configurations with functional structures in a highly general and economical way. Both our annotation methods are built on this insight: in the first, annotation principles are applied to CFG rules extracted from treebanks while in the second annotation principles are applied directly to flat term representations of treebank trees and tree fragments.

2.1 Regular expression based f-structure annotation of extracted CFGs

In this method, described in (Sadler et al. 2000), we extract a CFG from the treebank following (Charniak 1996) and develop a set of regular expression based annotation principles. The principles are applied to the extracted CFG to produce an annotated CFG. Annotated rules are then rematched against the original treebank trees and f-structures are produced from the annotations.

Annotation Principle Interpreter. Our CFG rule annotation principles are of the form $L \rightarrow R @ A$. A is a set of attribute-value structure annotations (rule decorations). L and R are regular expressions (under)specifying LHSs and RHSs of CFG rules in terms of categorial and configurational constraints. The regular expressions provided include Kleene and positive Kleene “*”, “+”, optionality “()”, disjunction “|” and a limited form of complement “~”. Operators are prefix and “{ }” is used to indicate grouping. “*” without argument denotes any string.

Given a grammar rule of the form $M \rightarrow Ds$ (expanding a mother category M into a sequence of daughter categories Ds) and a regular expression based annotation principle $L \rightarrow R @ A$, if the LHS L of the principle matches M and the RHS R matches Ds , then $M \rightarrow Ds$ is annotated with A . A single grammar rule can match multiple principles and a single principle may match a given grammar rule in more than one way. The annotations resulting from all possible matches are collected and the grammar rule is annotated accordingly.

More formally, let the denotation $[[E]]$ of a regular expression E be the set of strings denoted by E . Given a CFG rule $M \rightarrow Ds$ and a set of annotation principles AP of the form $L \rightarrow R @ A$, $M \rightarrow Ds$ is annotated with the set of feature structure annotations F :

$$M \rightarrow Ds @ F \text{ iff } F = \{A \mid \exists P \in AP \text{ with } P \equiv L \rightarrow R @ A \text{ and } M \in [[L]] \text{ and } Ds \in [[R]]\}$$

Annotation is monotonic and order independent.

Partial and Underspecified Annotation Principles. In our Prolog implementation, CFG grammar rules extracted from the treebank are represented as

$C:F \rightarrow C1:F1, \dots, Cn:Fn.$

where syntactic categories C and (optional) logical variables F representing feature-structure information are paired $C:F$. Regular expression based annotation principles can underspecify the LHS and RHS of grammar rules. To give a simple example, the following annotation principle¹ states that infinitival phrases $infp$ following the final $v0$ in vp rules are open complements ($xcomp$) controlled by the subject of the final $v0$:

$vp > * v0:V0 * \{\sim v0\} infp:I *$
 $@ [V0:xcomp = I, V0:subj = I:subj].$

The next principle states that in non-conjunctive contexts² $v0$ sequences, possibly separated by adverbials adv , form open complement sequences where the subject of the preceding $v0$ controls that of the following:

$vp > * \{\sim conj\} v0:V1 (adv) v0:V2 * \{\sim conj\}$
 $@ [V1:xcomp = V2, V1:subj = V2:subj].$

Note that the principle applies twice to a $\dots v0:V1, v0:V2, v0:V3 \dots$ RHS rule configuration with $[V1:xcomp = V2, V1:subj = V2:subj, V2:xcomp = V3, V2:subj = V3:subj]$ as the resulting annotation. Finally observe that the formalism supports the statement of generalisations over LHSs of CFG rules:

$\{fn:X|infp:X|tgp:X|si:X|vp:X\}$
 $> * \{\sim \{v0|conj\}\} v0:V0 * \{\sim conj\}$
 $@ [X = V0].$

This principle states that for a variety of constructions including verbal (vp) and infinitival ($infp$) phrases in non-conjunctive contexts the initial $v0$ is the head of the clause.

Example output (automatically annotated grammar rules from the AP fragment) is shown below:³

$vp:A \rightarrow v0:B, v0:C, v0:D, np:E, fa:F$
 $@ [A=B, D:obj=E, C:xcomp=D, C:subj=D:subj,$
 $B:xcomp=C, B:subj=C:subj, A:vp_adjunct:1=F].$

$vp:A \rightarrow v0:B, v0:C, v0:D, rp:E, pp:F$
 $@ [(D:obl=F; D:vp_adjunct:1=F), A=B, D:part=E,$
 $C:xcomp=D, C:subj=D:subj, B:xcomp=C,$
 $B:subj=C:subj].$

$vp:A \rightarrow vp:B, pnct:_, vp:C, pnct:_, conj:D, vp:E$
 $@ [A:conj:3=C, A=D, A:conj:2=B, A:conj:1=E].$

$vp:A \rightarrow vp:B, conj:C, vp:D, pp:E, fa:F$

@ [(D:obl=E;D:vp_adjunct:1=E),A=C,A:conj:2=B,
A:conj:1=D,A:vp_adjunct:1=F]).

In the first and in the second rule the leftmost v0 is identified as the head of the construction. In v0, v0 sequences the second v0 provides an open complement xcomp to the first with the subject of the second controlled by the subject of the first. The np in the first rule is analysed as the object of the rightmost v0, while the pp in the second rule is either an adjunct or an oblique argument to the vp. The last two example rules show coordinate structures. Note that in the final rule the pp is analysed as oblique or as an adjunct to the rightmost vp. Here our current annotation principles miss a possible attachment of the pp to the mother vp.

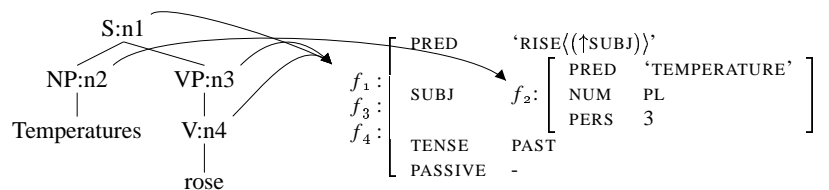
Automatic annotation is completed with macros for the preterminal tags contributing lexical information, e.g.:

nn1(Word):A @ [A:pred=Word, A:num=sg, A:pers=3rd].

The annotation principles together with the lexical macros constitute a principle-based c-structure/f-structure interface architecture for LFG.

2.2 F-structure annotation of treebank trees using flat tree descriptions

This method, described in (Frank 2000), builds on a pure correspondence view of the LFG architecture, where the mapping from c- to f-structure is encoded by the projection function ϕ . Annotation principles define ϕ -projection constraints which associate partial c-structures with their corresponding partial f-structures. Application of annotation principles to flat set-based encodings of treebank trees directly induces the f-structure, allowing us to skip the (re)matching process for f-structure composition. What is more, the principles can apply to non-local tree fragments, as opposed to local CFG rules.



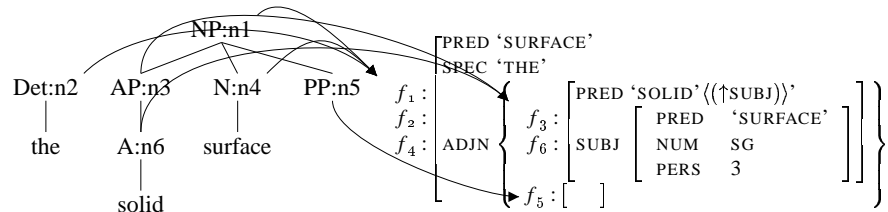
ϕ -correspondence: $\phi(n1) = f_1$, $\phi(n2) = f_2$, $\phi(n3) = f_3$, $\phi(n4) = f_4$, $\phi(n1) = \phi(n3) = \phi(n4)$

f-structure: $(f_1 \text{ SUBJ}) = f_2$, $(f_2 \text{ PRED}) = \text{'temperature'}$, $(f_4 \text{ PRED}) = \text{'rise'}$...

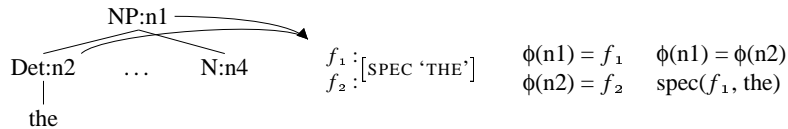
Modular projection principles for f-structure annotation of tree fragments.

To illustrate the key idea of partial f-structure annotation principles, below we display the representation of a complex NP. This complex configuration

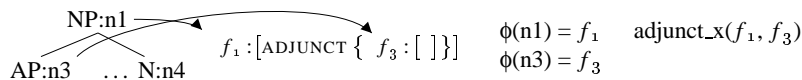
can be broken down into modular, piece-wise correspondences of *partial* c- and f-structures, abstracting away from irrelevant material in the surrounding context.



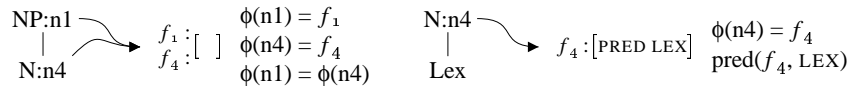
The functional contribution of the pronominal determiner *the* is independent of the presence of AP or PP, and is captured by the partial correspondence constraints stated on the right hand side.



An AP daughter of NP is analysed as an ADJUNCT of the nominal head, unless the N head is omitted. This generalisation is captured below.

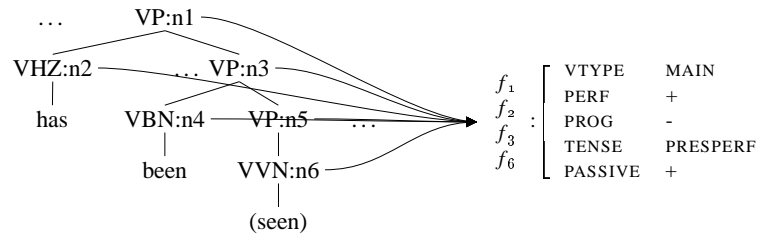


Projection principles for head categories and lexical nodes (here for nominal categories) are straightforward:



Similar correspondences are defined for the remaining c-structure fragments. These correspondences all apply to the complex NP structure above, conspiring to define the ϕ -projection and f-structure in a modular, declarative way. By abstracting away from immaterial c-structure context, the principles generalise over specific tree configurations, and therefore apply to fragments of unseen trees.

In the correspondence-based approach annotation principles can apply to *non-local* tree fragments. This allows us to associate partial f-structures with complex c-structure fragments. For example, by specifying non-local c-structure fragments in binary branching VPs, we capture tense and active/passive distinctions of the verbal complex in a natural way. This is illustrated for the characteristic construction indicative of present perfect tense.



The idea of modular annotation principles is much in the spirit of projection principles as proposed by (Dalrymple 2001) and (Bresnan 2001), and provides a principle-based c- to f-structure interface in the LFG architecture.⁴

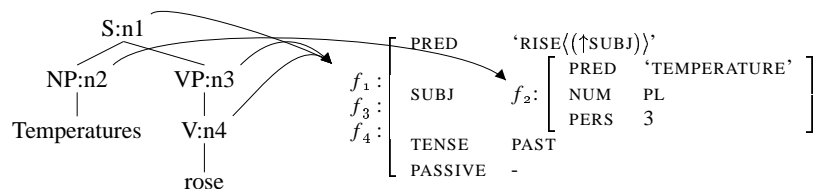
A term rewriting system for f-structure annotation. To define and process annotation principles we make use of an existing term rewriting system, originally designed for transfer-based Machine Translation (Kay 1999; Frank 1999).

The system takes as input an unordered set of n-ary terms $p, q, r \dots$, and an ordered set of rewrite rules $p_i \dots p_j \Rightarrow q_k \dots q_l$.⁵ If the LHS terms $p_i \dots p_j$ match the input, the matching terms are eliminated from the input set, and the terms $q_k \dots q_l$ are added to the output set. A rule applies to *each instantiation* of the LHS terms in the input. Besides terms p that are to be eliminated from the input, the LHS may state positive $+p$ and negative $-p$ terms. A rule with positive term $+p$ only applies if p matches some term in the input but positive terms are not eliminated from the input set. A rule with negative term $-p$ only applies if p does not match any term in the input. The order in which the rules are stated is crucial: Each rule applies to the *current* input set, and yields an output set. The output set of a rule constitutes the input set for the next rule.

A flat, term-based representation of the LFG architecture We encode the LFG projection architecture in a term representation language as follows:

- immediate dominance: `arc(MNode, MLabel, DNode, DLabel)`
- immediate precedence: `prec(CsNode_x, CsNode_y)`
- lexical insertion: `lex(TerminalNode, Lex)`
- ϕ -correspondence: `phi(CsNode, FsNode), equal(FsNode_x, FsNode_y)`
- f-structure attributes: `attr(FsNode_x, FsNode_y), attr(FsNode, Value)`

With this, the traditional representation

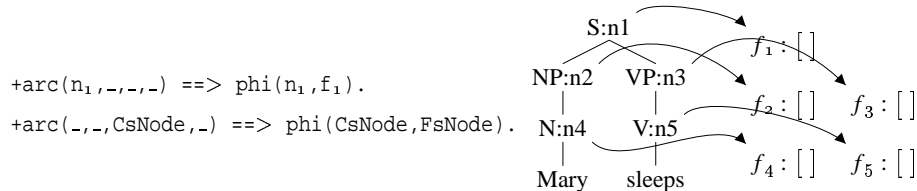


is translated into the following set of terms:

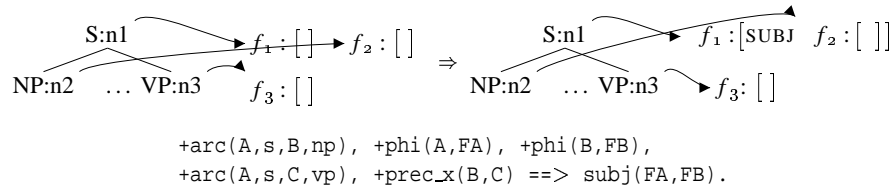
```
arc(n1,s,n2,np), arc(n1,s,n3,vp), arc(n3,vp,n4,v), prec(n2,n3),
lex(n2,Temperatures), lex(n4,rose),
phi(n1,f1), phi(n2,f2), phi(n3,f3), phi(n4,f4),
equal(f1,f3), equal(f3,f4),
pred(f1,raise), subj(f1,f2), pred(f2,temp.), num(f2,pl), tense(f1,past),..
```

Automatic annotation of flat tree descriptions with f-structures.

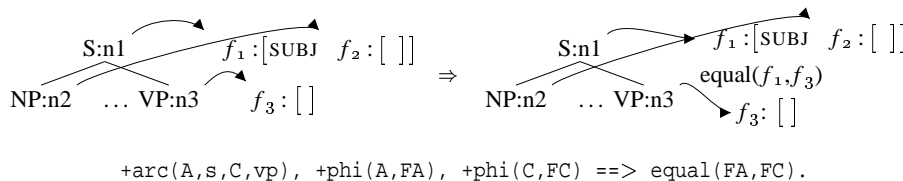
Initialisation Starting from the c-structure term representation, we induce a 1-1 ϕ -correspondence from c-structure nodes to empty f-structure nodes.⁶



Partial and underspecified annotation principles associate partial c-structure configurations with their corresponding partial f-structures, and further restrict the trivial 1-1 ϕ -correspondence via the predicate `equal(Fx,Fy)`. The rule below defines the VP-external NP as the SUBJ of f_1 , the f-structure projected from the S node. The predicate `prec_x(B,C)` is defined (by use of macros) as a finitely constrained transitive closure over the precedence relation `prec`. It can be used to underspecify precedence constraints holding between nodes n_x and n_y , allowing for an arbitrary or else a restricted sequence of intervening categories.



The following rule applies to the output resulting from the previous rule application. The predicate `equal(Fx,Fy)` restricts the ϕ -function to map the VP and S nodes to identical nodes in f-structure.



Formal restrictions Apart from initialisation we restrict ϕ predicates to only occur in LHSs of rules as *positive constraints*. Given the input specification of a 1-1 ϕ -projection, this guarantees that the functional property of the ϕ -correspondence is preserved. *equal* predicates only *restrict* the ϕ -correspondence, while preserving its functional property.

Order independence in a cascaded rewrite system Although annotation rules operate in a cascaded, order dependent way, order independence can be obtained by requiring that no annotation rule refers to f-structure information introduced by other rules, and no rule consumes (or adds) any c-structure information referred to by other rules. These constraints ensure that annotation rules have access to the full initial input structure, and no more than this, and thereby guarantee order independence of annotation, irrespective of the order in which the rules are stated and applied. The effect of order independence can be observed by inverting the application order of the subject and head-projection rules above: while the intermediate term set will be different, the final output set will be identical.

There is a trade-off between order dependence and independence. Constraining rules to c-structure information only can require complex rule constraints to prevent application of conflicting annotation rules to the same tree fragment, thereby avoiding inconsistencies. Moreover, reference to f-structure information can be used to generalise annotation rules. If several PS configurations are indicative of e.g. a subject function, or passive voice, such diverse configurations can be captured by referring to the more abstract f-structure information to further guide f-structure construction. The order of annotation rules must then ensure that the required f-structure information is introduced by previous annotation rules.

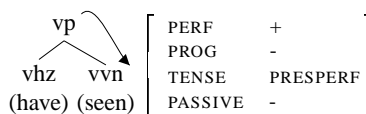
An annotation grammar consists, just like an ordinary LFG grammar, of different types of annotation rules: lexical, morphosyntactic, and phrasal.

Lexical and morphosyntactic rules Morphosyntactic rules introduce morphological (and some semantic) information encoded in lexical category labels into the f-structure space. The example given below illustrates how highly specific category distinctions in treebank encodings can be neutralised: once NUMBER is encoded in f-structure, based on the *nn1* vs. *nn2* distinction, this categorial distinction can be neutralised by mapping both lexical category labels to the generalised label *nn* (see van Genabith et al. 1999b for a similar approach). Such generalisations are essential for compact definition of annotation principles. For example, below the instantiation of the PRED-value of nouns is captured in a single lexical rule which applies to all “generalised” *nn*-daughters.

```
arc(A,ML,B,nn1) ==> num(B,sg), ntype(B,common), arc(A,ML,B,nn).
arc(A,ML,B,nn2) ==> num(B,pl), ntype(B,common), arc(A,ML,B,nn).
+arc(A,n,B,nn), +lex(B,Lex) ==> equal(A,B), pred(B,Lex), pers(B,'3').
```

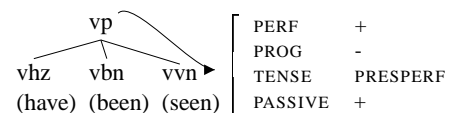
Tense information as well as the active/passive distinction can be captured by stating constraints on the partial c-structure context of verbs, as illustrated below for present perfect tense in a flat VP, as it is assigned in the Susanne corpus. For binary branching VPs (as assigned in the Penn-II Treebank, cf. Marcus et al 1994), we can define complex tense information in similar ways, by extending annotation rules to *non-local* tree fragments (see above and Frank 2000).

```
+arc(A,vp,B,vhz) % have-aux
-arc(A,vp,D,vbn) % no been-aux !
+arc(A,vp,C,vvn) % main verb participle
==> perf(A,+), prog(A,-),
    tense(A,presperf), passive(A,-).
```



PERF	+
PROG	-
TENSE	PRES PERF
PASSIVE	-

```
+arc(A,vp,B,vhz), % have-aux
+arc(A,vp,C,vbn), % been-aux
+arc(A,vp,D,vvn), % main verb part.
==> perf(A,+), prog(A,-),
    tense(A,presperf), passive(A,+).
```



PERF	+
PROG	-
TENSE	PRES PERF
PASSIVE	+

Partial phrasal rules and underspecification Annotation principles are designed to apply to modular, partial c-structure configurations, to define their corresponding functional projections. Even though treebanks do not tend to follow classical X-bar syntax, specific types of tree branches correspond to functional dependencies in f-structure. Annotation principles apply, in the general case, to single tree branches, with some contextual constraints, and generalise to unseen tree configurations. Below, *that*-clauses (category f) are associated with a function COMP in f-structure by referring to a single branch (arc) in c-structure, abstracting away from irrelevant co-occurrences in the c-structure context.

The example also illustrates the effect of underspecification. *That*-clauses can appear in different syntactic contexts. By referring to an *underspecified* (variable) mother node label ML, we generalise over various possible mother labels (e.g. (in)finite, modal, nominal or adjective phrases).

```
+arc(A,ML,B,f), +comp_form(B,that) => comp(A,B).
```

Finer categorial restrictions can be captured by defining classes of category labels in disjunctive templates.⁷ Below, the disjunctive template `np_cat(XL)` defines a class of category labels (n, d, m). The template is called (by logical

“and” &&) in the annotation rule for PPs (p) to define this restricted class of alternative NP-types as complements (i.e., OBJ) of prepositions.

```
template definition:  np_cat(XL) :: { XL == n } ==> 0; % n: nominal phrase
                       { XL == d } ==> 0; % d: determiner phrase
                       { XL == m } ==> 0. % m: number phrase

annotation rule:      +arc(A,p,B,NPL) ==> obj(A,B) && np_cat(NPL).
```

Grammatical function assignment In languages like English, grammatical function assignment relies heavily on c-structure configurations, while still not being fully deterministic. In case marking languages, morphological marking will be used to constrain grammatical function assignment. Below we give an example for the assignment of OBJ vs. OBJ2 functions for transitive and ditransitive verbs in English, which is determined by surface order. Long-distance phenomena are captured by path expressions (see Frank 2000 for further details).

```
+arc(A,vp,C,np), +arc(A,vp,D,np), +prec_x(C,D) ==> obj2(A,D).% OBJ2 ditrans
+arc(A,vp,C,np), +arc(A,vp,D,np), +prec_x(C,D) ==> obj(A,C). % OBJ ditrans
+arc(A,vp,C,np), -arc(A,vp,D,np), {D \== C} ==> obj(A,C).8 % OBJ trans
```

Subcategorisation assignment We induce subcategorisation frames (the semantic forms) by collecting grammatical functions assigned by annotation rules into the predicate’s semantic form, following the method of (van Genabith et al. 1999a).

Obviously, pure c-structure information does not allow us to distinguish between NP, PP, and infinitival arguments vs. adjuncts. Similarly, lacking lexical information, raising and control constructions can only be represented as involving anaphoric control. In (Frank 2000) we show how to extend this model by integration of lexical subcategorisation information, combined with strategies for OT-based ambiguity ranking and filtering (cf. Frank et al. (2000)).

Partial annotation and robustness Our f-structure annotation method embodies an important aspect of robustness. In the case of missing annotation principles the system does not fail, but partial trees are left without f-structure annotation. We obtain (typically large) partial, unconnected f-structures.

Moving treebanks Our framework can also be used to adjust particular treebank encodings, by “moving” treebanks to a different structural encoding, thereby facilitating principle-based f-structure induction. In our treatment of the Susanne corpus, we defined a set of c-structure rewriting rules to transform the encoding of coordination and flat modal VP structures into more standard PS analyses, which lend themselves to principle-driven f-structure annotation.

3. TWO EXPERIMENTS

3.1 Experiment I

Experiment Design. Our first experiment involves the first 100 trees of the AP treebank (Leech and Garside 1991). We refer to this subsection as AP01. We preprocess the treebank using the structure preserving grammar compaction method reported in (van Genabith et al. 1999b) preserving as much categorial fine-grainedness as is required to guide annotation. From this we extract a treebank grammar following (Charniak 1996). We develop a set of feature structure annotation principles. The regular expression based interpreter described in Section 2.1 compiles the principles over the rules extracted from the AP01 treebank fragment. The results obtained are compared against a manually annotated “gold standard” reference grammar and precision and recall measures are reported.⁹

Data. The AP treebank annotation schema employs 183 lexical tag types and 53 phrasal category types, with tree structure encoded in terms of labelled bracketing. The corpus is “skeletonally parsed”, that is, it contains some unlabelled brackets. We remove these in an automatic pre-editing step. The sentences in the AP01 fragment range from 4 to 50 leaf tokens (including punctuation symbols). The AP01 section of the corpus attests 94 of the 183 lexical tag types and 25 of the 53 phrasal tag types. The large number of highly discriminating lexical and phrasal categories results in a large number of flat and often very specific rules. To facilitate annotation we use the structure preserving grammar compaction method presented in (van Genabith et al. 1999b) to compact the grammar into a more general one that still preserves important categorial information to drive automatic annotation. Compaction works by generalising tags, i.e. collapsing tags (and categories) into supertags. This reduces the number of CFG rule types in the fragment from 511 to 330. AP01 and the compacted AP01^c are summarised in **Table 1** below:

T1	sentences	average length	phrasal types	lexical types	CFG rule types
AP01	100	20	25	94	511
AP01 ^c	100	20	12	28	330

Manually Annotated Reference Grammar. In order to evaluate Experiment I, we manually constructed a “gold standard” reference grammar following (van Genabith et al. 1999a,b,c). The grammar features 1143 annotations, on average 3.46 annotations per rule.

Automatic Annotation and Evaluation. For the experiment we constructed 119 annotation principles, this against 330 CFG rules resulting in a

template/rule ratio of 0.36. We expect the ratio to skew substantially in favour of templates as we proceed to larger fragments (see Section 4). Automatic annotation generates 1029 annotations, on average 3.12 annotations per rule. Experiment I is evaluated in terms of *precision* and *recall* measures:

$$\text{precision} = \frac{\# \text{ generated annotations also in reference}}{\# \text{ generated annotations}}$$

$$\text{recall} = \frac{\# \text{ reference annotations also generated}}{\# \text{ reference annotations}}$$

The results are summarised in **Table 2**:

T2	Experiment I
precision	93.38
recall	91.58

The numbers are conservative: *precision* and *recall* are computed automatically for a first pass encoding of annotation principles as regular expressions. The results are encouraging and indicate that automatic annotation is more often partial than incorrect.

3.2 Experiment II

Our method for f-structure annotation of trees in Section 2.2 is evaluated in Experiment II, this time based on the Susanne corpus (Sampson 1993).

Data The Susanne treebank encodes labelled bracketed structures with surface form and lemmatised lexical entries. Functional category labels (subj, obj) and traces indicating control or long-distance dependencies are eliminated in preprocessing, to guarantee a non-biased evaluation with conventional PS trees as input. In preprocessing we also collapse overspecific phrasal categories.

Some decisions on PS assignment in the Susanne corpus are debatable. We defined a set of c-structure rewriting rules that transform the encoding of coordination and flat modal VP structures to more standard PS analyses.

Experiment Design We chose two sections of the Susanne corpus, J01 and J02 (text type J: learned writing). On these, we ran an experiment in 3 steps:

First, we develop f-structure annotation principles for the first 66 sentences of J01. These generate fully connected f-structures for 50 out of the 66 sentences. In step 2 we apply the resulting annotation grammar AG1 to the first 50 (unseen) sentences of J02 (J02-1), and measure the annotation results. Grammar AG1 is then upgraded to AG2, to fully cover these additional 50 sentences. We record the number of principles added or modified. In step 3, the annotation grammar AG2 is applied to the remaining 46 (unseen) sentences of the

second part of J02 (J02-2). Again, we measure the results. In this experiment we applied an order dependent annotation scheme that consumes c-structure terms while building up the f-structure (cf. Frank 2000). We established a natural order for the different types of annotation principles discussed in Section 2.2.

Evaluation and Results Table 3 provides basic data of the treebank subsections: the number of sentences and average sentence length; the number of phrasal and lexical categories and the number of distinct PS rules and PS branches encoded by the corpus trees. Note that the percentage of new (unseen) PS rules in J02-1 and J02-2 is considerably higher than for new (unseen) tree branches. This is not surprising, and supports our annotation scheme, where annotation involves underspecified, partial trees (often single branches; cf. discussion in Section 4).

	sent.	av. length	phrasal cat	lexical cat	PS rules	tree branches
J01	66	34.27	32	73	430	281
J02-1	50	21.68	25 (3 new)	64 (8 new)	249 (60.34% new)	172 (20.93% new)
J02-2	46	24.8	24 (4 new)	57 (3 new)	212 (45.28% new)	163 (15.95% new)

The results of automatic f-structure annotation are summarised in Table 4. We measured correctness of f-structure assignment modulo the argument/ad-junct distinction for PPs and infinitival VPs, and the missing assignment of control/raising equations. Also, attachment or labelling mistakes in the treebank are not counted as annotation mistakes if the resulting f-structure is predicted from the given tree.

AG1 features 118 non-lexical (phrasal) annotation principles and assigns correct f-structures to 48% of the unseen section J02-1. As expected, the upgrade from AG1 to AG2 required little effort: it involves 28 new and 5 modified rules and required approximately one person-day of work. AG2 applied to the unseen section J02-2 yields 76.09% of correct f-structures.

	correct fs		partial fs		tag rules	lexical rules	phrasal rules	all rules
	#	%	#	%				
J01 w/ AG1	50	75.76	16	24.24	41	132	118	291
J02-1 w/ AG1	24	48	26	52	41	132	118	291
J02-1 w/ AG2	49	98	1	2%	41+4	132+4 (2 mod)	118+20 (3 mod)	291+28
J02-2 w/ AG2	35	76.09	11	23.91	45	136	138	319

Although small scale, we consider these results promising. Our experiment yields 76% correctly assigned complete and fully connected f-structures when applied to unseen trees, on the basis of a stepwise extended annotation grammar, developed for about 100 sentences. The increase of coverage when moving from AG1 to AG2 is considerable. Upgrading to larger fragments takes

little effort due to the generalisation capacity of partial annotation principles. The latter is confirmed by the increasing percentage of correct f-structure assignments to unseen trees, and the fact that partial f-structure assignments generally consist of large pieces of partial f-structures.

4. DISCUSSION AND CURRENT RESEARCH

We have presented two companion automatic f-structure annotation methods (Sadler et al. 2000; Frank 2000) for treebanks and grammars. Both methods and the experiments show considerable overlap and several interesting differences.

Annotation principles can apply to extracted PS rules or to PS tree fragments encoded as flat term representations. Our second method can be specialised to PS rules by restricting trees to depth one. The first method generates an annotated grammar, which can be used to rematch treebank trees to induce f-structures or serve as a basis for developing a stand-alone LFG resource. In the second approach an f-structure is built during the annotation process. In order to parse free text, this method can be applied to the output of (P)CFG parsing. The same architecture can be implemented using the principles designed in the first approach. Our second approach can be modified to annotate (non-local) tree fragments with f-descriptions for the rematching scenario applied in the first method. Both methods use compaction techniques for generalising over-specific categorisation. In the first experiment the structure of treebank entries remains unchanged, while in the second certain structures are transformed to conventional PS analyses to support principle-based annotation. For our first method, we implemented an order independent and monotonic annotation principle interpreter. For the second, a more general term rewriting system was used. The term rewriting system allows us to define order dependent, cascaded processing of annotation principles. Alternatively, the term rewriting system can implement order independent annotation. Order independence can sometimes ease maintenance of annotation principles, but requires more complex and verbose constraints in order to avoid inconsistent annotations. By contrast, order dependent cascaded rewriting allows for a compact representation of annotation principles. The extra power of an order dependent system can be useful in category generalisation and subcategorisation induction during the annotation process. Experiment I uses a manually constructed “gold standard” reference grammar for evaluation, experiment II is evaluated with respect to how it performs on unseen, extended treebank fragments.

Robustness is an inherent property of the approaches presented here. It resides in a number of levels: First, our principles are partial and underspecified and will match new, as yet unseen configurations. Second, the principles are conditional. If a certain context (a regular expression or a constraint set) is

met, a principle applies. Even if only few principles apply, the system will not fail but deliver partial annotations. Third, the constraint solver employed in our second method can cope with conflicting information. A constraint solver of this type can also be imported into the processing of rules annotated by our first method.

Both approaches factor out information spread over many CFG grammar rules into a smaller number of modular and general principles. To a first approximation, the reason why our principles allow a compact representation of grammatical knowledge is the following: by and large the annotation principles capture statements about single mother-daughter relationships in CFG rules or local trees of depth one. This means that the principles are essentially about single branches in local configurations. Given a treebank (grammar) with n distinct categories the worst case number of distinct branches is n^2 . Contrast this with the worst case number of possible grammar rules:

$$\begin{array}{rcl} \#(x \rightarrow y_1) & \hat{=} & n^2 \\ \#(x \rightarrow y_1 y_2) & \hat{=} & n^3 \\ \dots & \dots & \dots \\ \#(x \rightarrow y_1 \dots y_m) & \hat{=} & n^{m+1} \end{array}$$

Clearly, given a grammar with n categories and a RHS rule length of at most m , the worst case number of different grammar rules

$$\sum_{i=1}^m n^{i+1} \gg n^2$$

for $m > 2$ is much higher than the worst case number n^2 of distinct branches.

In recent research we have scaled an automatic f-structure annotation approach evolved from the methods presented here to the complete Penn-II treebank resource (Cahill et al. 2002a, 2002b) to generate f-structures for 49,000 trees and 1 million words.

In order to develop stand-alone LFG grammars we need semantic forms (subcategorisation lists) to enforce subcategorisation requirements. We are currently exploring a number of ways of semi-automatically compiling these from machine readable dictionaries and the f-structure annotated corpus resources produced.

We expect that our approach can also feed into grammar development efforts. To be sure, because treebank grammars are large and flat, automatically annotated treebank grammars are less maintainable than the more compact, linguistically designed grammars which follow X-bar design principles. However, as pointed out above, our approaches allow for a novel grammar design and processing architecture: given a treebank, a probabilistic context-free

grammar compiled from the treebank parses new text. For each input string, the (possibly n -) best parse trees are passed on to the annotation interpreters which annotate or rewrite the parse trees and induce f-structures. This and other probabilistic parsing architectures are developed in (Cahill et al. 2002b) and applied to parse the WSJ section of the Penn-II treebank into proto-f-structures. We consider this a promising new approach to partially automate large-coverage, corpus-based unification grammar development.

Current research also investigates further applications of flat, term-based tree structure conversion to induce grammars for alternative formalisms from existing treebanks. (Frank 2001) describes a treebank conversion method, applied to the German NEGRA corpus (Brants et al. 1997) to extract an LTAG grammar of German. The same method and corpus was used in (Becker and Frank 2002) to extract a stochastic topological grammar of German, to be used for integrated shallow and deep parsing.¹⁰ (Liakata and Pulman 2002) present a method based on flat, term-based tree representations that closely resembles the original approach in (Frank 2000), in order to annotate Penn-II treebank trees with Quasi-Logical Forms information while (Cahill et al. 2003) show how simple Quasi-Logical Forms can be generated from f-structures produced for the Penn-II trees in (Cahill et al. 2002a, 2002b).

5. SUMMARY

We have presented two companion automatic f-structure annotation methods (Sadler et al. 2000; Frank 2000) for treebanks and grammars. The approaches make use of a corpus-based strategy that takes disambiguated tree structures as input, and annotate using (linguistically motivated) annotation principles. The principles are used to automatically enrich treebanks or extracted treebank grammars with higher-level functional information not present in the original corpora. Automatic annotation holds considerable potential in curtailing f-structure bank development costs and opens up the possibility of tackling large fragments. The work reported here is proof of concept. (Cahill et al. 2002a, 2002b) have further developed automatic f-structure annotation technology based on the methods described here and successfully scaled it to the Penn-II treebank resource. Here, we have presented a grammar development and treebank annotation methodology which is data-driven, semi-automatic, reuses existing resources and covers real text. We found the LFG framework very conducive to our experiments. We do believe, however, that the methods can be generalised, and we intend to apply them in an HPSG scenario and we and other researchers have applied similar technology to automatic semantic representation based annotation (Liakata and Pulman 2002; Cahill et al. 2003).

Our approach encourages work in the best linguistic tradition as (i) it is concerned with real language and (ii) enforces generalisations in the form of

annotation principles. Our methods factor out information spread over many CFG rules into a small number of modular and general principles. What is new in our approach is that (i) the principles state partial and underspecified correspondences between c- and f-structure configurations and (ii) they are applied to flat and noisy treebank representations that do not follow general X-bar design principles. Our experiments show how theoretical work and ideas on principles can translate into grammar development for real texts. In this sense our approach may contribute to bridge the often-perceived gap between theoretically motivated views of grammar as a set of principles vs. grammars for “real” text.

Acknowledgements

The authors wish to thank Tracy H. King, the members of the Pargram group, in particular Ron Kaplan, Mary Dalrymple and John Maxwell as well as Joan Bresnan, for helpful discussions and feedback.

Notes

1. For expository purposes, these are slightly simplified principles from our annotation grammar.
2. The annotation principles have to take into consideration that, in many cases, the representation of coordination in treebank rules is overly flat.
3. The annotation process itself is fast: in our experiments the interpreter annotates about 40 treebank CFG rules per second (Sparc 400Mhz).
4. It is also closely related to the principle-based grammar architecture of HPSG, cf. related work by (Neumann and Flickinger 1999) and (Neumann, this volume).
5. There are obligatory (\Rightarrow) and optional ($?\Rightarrow$) rewrite rules.
6. n_1 refers to the tree’s root node.
7. Disjunctive templates encode alternative rewrite rules, and can be unioned (by logical “and” $\&\&$) with annotation rules. While this does still involve disjunctive processing, the rules can be stated in a generalised, compact way.
8. We require B and C to be distinct variables through inequality constraints (in curly brackets).
9. Templates, grammars and f-structures generated are available at: <http://www.compapp.dcu.ie/~away/Treebank/treebank.html>.
10. In this work, we developed a simple rewriting system modeled after the term rewriting system of Kay(1999).

References

- S. Abney. (1997). Stochastic Attribute-Value Grammars. In: *Computational Linguistics*, 23(4), p. 597–618.
- M. Becker, A. Frank. (2002). A Stochastic Topological Parser of German. *Proceedings of COLING 2002*, Taipei, Taiwan.
- R. Bod, R. Kaplan. (1998). A Probabilistic Corpus-driven Model for Lexical-Functional Analysis. *Proceedings of COLING/ACL’98*, p. 145–151.

- T. Brants, W. Skut, B. Krenn. (1997). Tagging Grammatical Functions. *Proceedings of EMNLP*, Providence, RI, USA.
- J. Bresnan. (2001). *Lexical-Functional Syntax*. Blackwells Publishers, Oxford.
- A. Cahill, M. McCarthy, J. van Genabith, A. Way. (2002a). Automatic Annotation of the Penn-Treebank with LFG F-Structure Information. A. Lenci, S. Montemagni and V. Pirelli, editors, In: *LREC 2002 Workshop on Linguistic Knowledge Acquisition and Representation - Bootstrapping Annotated Language Data*, LREC 2002, Third International Conference on Language Resources and Evaluation, p. 8-15.
- A. Cahill, M. McCarthy, J. van Genabith, A. Way. (2002b). Parsing with PCFGs and Automatic F-Structure Annotation. In: M. Butt and T.-H. King, editors: *Proceedings of the LFG 2002 Conference*, Technical University Athens, Greece.
- A. Cahill, M. McCarthy, J. van Genabith, A. Way. (2003). Quasi-Logical Forms from F-Structures for the Penn Treebank. Fifth International Workshop on Computational Semantics (IWCS-5). *Proceedings of the Workshop*, Tilburg, The Netherlands.
- E. Charniak. (1993). *Statistical Language Learning*. MIT Press, Cambridge MA.
- E. Charniak. (1996). Tree-bank Grammars. *AAAI-96. Proceedings of the Thirteenth National Conference on Artificial Intelligence*, p. 1031–1036. MIT Press.
- M. Dalrymple, R.M Kaplan, J.T. Maxwell III, and A. Zaenen, editors. (1995). *Formal Issues in Lexical-Functional Grammar*. CSLI Lecture Notes, No. 47. CSLI Publications.
- M. Dalrymple. (2001). *Lexical-Functional Grammar*. Syntax and Semantics 34, Academic Press.
- A. Frank. (1999). From Parallel Grammar Development towards Machine Translation. A Project Overview. *Proceedings of Machine Translation Summit VII "MT in the Great Translation Era"*, p. 134–142.
- A. Frank. (2000). Automatic F-Structure Annotation of Treebank Trees. In: M. Butt and T.H. King editors, *Proceedings of the LFG00 Conference*, University of California at Berkeley, CSLI Online Publications, Stanford, CA, <http://www-csli.stanford.edu/publications/>.
- A. Frank, T. King, J. Kuhn, J. Maxwell. (2000). Optimality Theory Style Constraint Ranking in Large-scale LFG Grammars. In: Sells, P., editor, *Optimality Theoretic Syntax*. CSLI Publications, Stanford, CA.
- A. Frank. (2001). Treebank Conversion. Converting the NEGRA Treebank to an LTAG Grammar. *Proceedings of the Workshop on Multi-layer Corpus-based Analysis*, Workshop of the EUROLAN 2001 Summer Institute on Creation and Exploitation of Annotated Language Resources, Iasi, Romania.

- R.M. Kaplan, J. Bresnan. (1982). *Lexical Functional Grammar*, p. 173–281. MIT Press, Cambridge, Mass.
- M. Kay. (1999). Chart Translation. *Proceedings of Machine Translation Summit VII "MT in the Great Translation Era"*, p. 9–14.
- A. Krotov, M. Hepple, R. Gaizauskas, Y. Wilks. (1998). Compacting the Penn Treebank Grammar. *Proceedings of COLING/ACL'98*, p. 699–703.
- G. Leech, R. Garside, (1991). *Running a Grammar Factory: On the Compilation of Parsed Corpora, or 'Treebanks'* in S. Johansson et al (eds) *English Computer Corpora: selected papers*, p. 15–32. Mouton de Gruyter, Berlin.
- M. Liakata, S. Pulman. (2002). From Trees to Predicate-Argument Structures. *Proceedings of COLING 2002*, Taipei, Taiwan.
- M. Marcus, G. Kim, M. Marcinkiewicz, R. MacIntyre, A. Bies, M. Ferguson, K. Katz, B. Schasberger. (1994). The Penn treebank: Annotating predicate argument structure. In: *ARPA Human Language Technology Workshop*.
- G. Neumann, D. Flickinger. (1999). HPSG-DOP: Data-oriented Parsing with HPSG. Learning Stochastic Lexicalized Tree Grammars from HPSG. DFKI Technical Report, Saarbrücken, 1999.
- G. Neumann. (2003). A Uniform Method for Automatically Extracting Stochastic Lexicalized Tree Grammars from Treebanks and HPSG. In this volume.
- C. Pollard, I. Sag. (1994). *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago, Illinois.
- L. Sadler, J. van Genabith, A. Way. (2000). Automatic F-Structure Annotation from the AP Treebank *Proceedings of the LFG 2000 Conference*, The University of California at Berkeley, CSLI Publications, Stanford, CA, <http://www-csli.stanford.edu/publications/>
- G. Sampson, (1993). *The Susanne Corpus*. Release 2.
- J. van Genabith, L. Sadler, A. Way. (1999a). Data-driven Compilation of LFG Semantic Forms. In: *EACL'99 Workshop on Linguistically Interpreted Corpora (LINC-99)*, Bergen, Norway, June 12th, p. 69–76.
- J. van Genabith, L. Sadler, A. Way. (1999b). Structure Preserving CF-PSG Compaction, LFG and Treebanks. *Proceedings ATALA Workshop - Treebanks*, Journées ATALA, Université Paris 7 p. 107–114.
- J. van Genabith, A. Way, L. Sadler. (1999c). Semi-Automatic Generation of f-Structures from Tree Banks. In: M. Butt and T.H. King, editors, *Proceedings of the LFG99 Conference*, Manchester University, CSLI Online Publications, Stanford, CA. <http://www-csli.stanford.edu/publications/>.

Appendix: Example of an Automatically Generated F-Structure (Susanne Corpus)

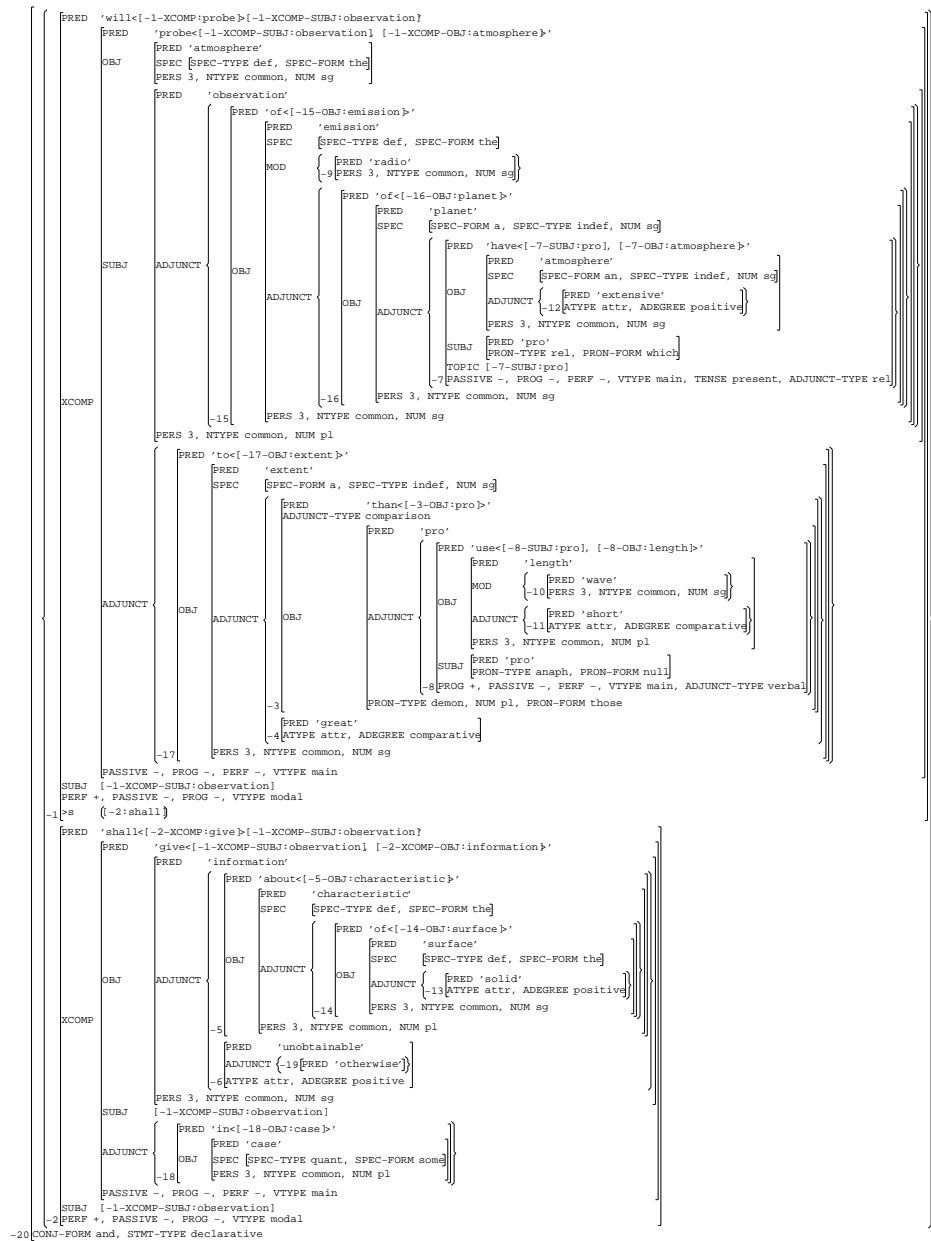


Figure 21.A.1. F-structure for: "Observations of the radio emission of a planet which has an extensive atmosphere will probe the atmosphere to a greater extent than those using shorter wave lengths and should in some cases give otherwise unobtainable information about the characteristics of the solid surface."