

---

# Grammatical Framework and Multiple Context-Free Grammars

PETER LJUNGLÖF

Grammatical Framework (GF) (Ranta, 2004) is a grammar formalism originating from logical frameworks for dependent type theory. It is already known that the parsing problem for GF is undecidable, which has to do with the possibility to formulate undecidable propositions. But for subclasses of GF, in particular GF with a context-free backbone, parsing is decidable. Until now the parsing complexity of context-free GF has been unknown, which we aim to change with this article.

We show that there is a simple one-to-one correspondence between Grammatical Framework with context-free backbone and Multiple Context-Free Grammars (MCFG) (Seki et al., 1991). Since the parsing complexity for MCFGs is known to be polynomial in the length of the input, we get the same result for context-free GF.

GF with dependent types is undecidable in general, but the separation of abstract and concrete syntax makes it possible to use a two-step process in parsing. First a MCFG parser is used to create a chart, then each item in the chart is converted to a Horn clause. The resulting logic program then solves the parsing problem.

## 6.1 Grammatical Framework

Grammatical Framework (Ranta, 2004) is a grammar formalism built upon a Logical Framework. What GF adds to the logical framework is a possibility to define concrete syntax, that is, notations express-

ing formal concepts in user-readable ways. Although GF grammars are bidirectional, the perspective of GF is on generation rather than parsing. A difference from usual grammar formalisms is the support for multilinguality; it is possible to define several concrete syntaxes upon one abstract syntax. The abstract syntax then works as an interlingua between the concrete syntaxes. The development of GF as an authoring system started as a plug-in to the proof editor ALF, to permit natural-language rendering of formal proofs (Hallgren and Ranta, 2000). The extension of the scope outside mathematics was made in the Multilingual Document Authoring project at Xerox (Dymetman et al., 2000). In continued work, GF has been used in areas like software specifications (Hähnle et al., 2002) and dialogue systems (Ranta and Cooper, 2004). In this article we have changed the notation somewhat, compared to other GF articles. The reason for this is to make the comparison with Generalized Context-Free Grammars simpler. There are also some features of GF that we do not mention, mainly syntactic sugar and some default notations.

### 6.1.1 GF with a context-free backbone

Grammatical Framework is a strongly typed functional language; functional in the sense that a grammar defines a set of functions, and strongly typed in the sense that each function has a given typing, known at compile-time. Basic types in GF are called *categories*, and are specified by the grammar. For now we assume that there are only a finite number of categories – in section 6.1.3 we generalize the concept to dependent and higher-order categories.

A function is specified by a typing, an (optional) definition and a linearization. To retain compositionality, all terms of a given category  $A$  must linearize to the same *linearization type*, written  $\llbracket A \rrbracket$ .

**Function typings** The typing of a function tells us how many arguments it takes, what their categories are and what category the result is.

$$f : A_1 \rightarrow \dots \rightarrow A_\delta \rightarrow A$$

The function may take no arguments; in which case it is called a *ground term*. Given a grammar, we define the abstract terms, or syntax trees, as follows. The term  $f(t_1, \dots, t_\delta)$  is of category  $A$  whenever each  $t_i$  is of category  $A_i$ .

**Abstract definitions** An abstract definition specifies a computation rule on terms, thereby defining a notion of equality between terms. This equality is a semantic equality, since it does not affect the concrete linearizations; even though two terms are equal by definition, they can

be linearized to different strings. Since they do not affect the concrete syntax, we do not mention them further in this article.

**Linearization terms and types** Linearizations are written as terms in a typed functional programming language, which is limited to ensure decidability in generation and in parsing. The language has records and finite-domain functions (called tables); and the basic types are terminal lists (called strings) and finite data types (called parameter types). There are also local definitions, lambda-abstractions and global function definitions. The parameters are declared by the grammar; they can be hierarcical but not recursive, to ensure finiteness.

**Linearization definitions** The linearization of a function term  $f$  is a function (also written  $f$ ) from the linearization types of the argument categories to the linearization type of the result category; or as a functional typing,  $f : \llbracket A_1 \rrbracket \rightarrow \dots \rightarrow \llbracket A_\delta \rrbracket \rightarrow \llbracket A \rrbracket$ .

$$f(x_1, \dots, x_\delta) = \phi$$

The linearization  $\phi$  must of course be of type  $\llbracket A \rrbracket$  (written  $\phi : \llbracket A \rrbracket$ ), given that each  $x_i$  is of type  $\llbracket A_i \rrbracket$ . Note that a limitation on the linearization is that it is not possible to examine the structure of variables by e.g. case analysis; which is one way of defining compositionality for grammars.

### 6.1.2 Canonical GF

The concrete syntax of any GF grammar can be partially evaluated to a grammar in canonical form, as shown in Ranta (2004). In canonical form, all local and global definitions disappear, as well as function applications; and all tables are fully expanded. Hierarcical parameters can be flattened; thus we can assume that the parameters are declared by giving a finite set  $\text{Par}$  of parameter types, each  $P \in \text{Par}$  being a set of parameters  $p_1, \dots, p_n$ . A linearization term in canonical GF is of the following form:

- A string constant is of type  $\text{Str}$ ; and a concatenation  $s_1 \cdot s_2 : \text{Str}$ , whenever  $s_1, s_2 : \text{Str}$ .
- A constant parameter  $p : P$ , whenever  $p \in P$ .
- A record  $\{r_1 = \phi_1; \dots; r_n = \phi_n\}$  is of type  $T = \{r_1 : T_1; \dots; r_n : T_n\}$ , whenever each  $\phi_i : T_i$ .
- A record projection  $\phi.r_i : T_i$ , whenever  $\phi$  is of the record type  $T$  above.
- A table  $[p_1 \Rightarrow \phi_1; \dots; p_n \Rightarrow \phi_n]$  is of type  $P \Rightarrow T$ , whenever  $P = \{p_1, \dots, p_n\}$ , and each  $\phi_i : T$ .
- A table selection  $\phi! \psi : T$ , whenever  $\phi : P \Rightarrow T$  and  $\psi : P$ .

$$\begin{aligned}
\llbracket S \rrbracket &= \{ s : \text{Str} \} \\
\llbracket N \rrbracket &= \{ s : \text{Str}; n : \text{Num} \} \\
\llbracket V \rrbracket &= \{ s : \text{Num} \Rightarrow \text{Str} \} \\
p &: N \rightarrow V \rightarrow S \\
p(x, y) &= \{ s = x.s \cdot y.s!(x.n) \} \\
q &: V \rightarrow N \rightarrow V \\
q(x, y) &= \{ s = [n \Rightarrow x.s!n \cdot y.s] \} \\
a : N &= \{ s = \text{“animals”}; n = \text{Pl} \} \\
b : N &= \{ s = \text{“Bernie”}; n = \text{Sg} \} \\
l : V &= \{ s = [\text{Sg} \Rightarrow \text{“loves”}; \text{Pl} \Rightarrow \text{“love”}] \}
\end{aligned}$$

FIGURE 1 Example GF grammar

- An argument variable  $x_i : \llbracket A_i \rrbracket$ .

Together with this there are computation rules for string concatenation, record projection and table selection.

**An example grammar** In figure 1, there is a simple example GF grammar, which is not entirely in canonical form. The linearization of  $q$  contains a non-expanded table  $[n \Rightarrow x.s!n \cdot y.s]$ , whose canonical form is  $[\text{Sg} \Rightarrow x.s!\text{Sg} \cdot y.s; \text{Pl} \Rightarrow x.s!\text{Pl} \cdot y.s]$ .

### 6.1.3 GF with dependent categories

Full GF have a much more expressive abstract syntax than above; categories can depend on other categories. E.g. the grammar could specify that category  $A$  depends of category  $B$ , meaning that  $A(x)$  is a category whenever  $x$  is a term of category  $B$ .

Another extension is that arguments to function typings (and dependent categories) can be functions and not just of a basic category.

These extensions turn the abstract syntax into a logical framework, where e.g. undecidable propositions can be formulated, thus giving a very expressive formalism. A natural question is then how to parse and linearize terms with these extensions, which will be addressed in section 6.4.1.

## 6.2 Generalized CFG

Generalized Context-Free Grammars (GCFG) were introduced by Pollard in the 80's as a way of formally describing Head Grammars (Pollard, 1984). Later people have used GCFG as a framework for describing

many other formalisms, such as Linear Context-Free Rewriting Systems (Vijay-Shanker et al., 1987) and Multiple Context-Free Grammars (Seki et al., 1991); and here we will use it to describe GF with a context-free backbone.

There are several definitions of GCFG in the literature, and we introduce yet another one, in which a Generalized Context-Free Grammar consists of an abstract grammar together with a concrete interpretation.

**Abstract grammar** The abstract grammar is a tuple  $(C, S, F, R)$ , where  $C$  and  $F$  are finite sets of categories and function symbols respectively,  $S \in C$  is the starting category, and  $R \subseteq C \times F \times C^*$  is a finite set of abstract syntax rules. For each function symbol  $f \in F$  there is an associated context-free syntax rule.

$$A \rightarrow f(A_1, \dots, A_\delta)$$

The tree rewriting relation  $A \Rightarrow t$  is defined as  $A \Rightarrow f(t_1, \dots, t_\delta)$  whenever  $A_1 \Rightarrow t_1, \dots, A_\delta \Rightarrow t_\delta$ .

**Concrete interpretation** To each category  $A$  is associated a *linearization type*  $\llbracket A \rrbracket$ , which is not further specified. To each function symbol  $f$  is associated a partial *linearization function* (also written  $f$ ), taking as many arguments as the abstract syntax rule specifies.

$$f \in \llbracket A_1 \rrbracket \rightarrow \dots \rightarrow \llbracket A_\delta \rrbracket \rightarrow \llbracket A \rrbracket$$

The linearization of a syntax tree is defined as  $\llbracket f(t_1, \dots, t_\delta) \rrbracket = f(\llbracket t_1 \rrbracket, \dots, \llbracket t_\delta \rrbracket)$ , if the application is defined. Note that we impose no restrictions on the linearization types or the linearization functions; this is left to the actual grammar formalism.

### 6.2.1 GF with a context-free backbone

Grammatical Framework with a context-free backbone is obviously an instance of GCFG, where the abstract GF rule  $f : A_1 \rightarrow \dots \rightarrow A_\delta \rightarrow A$  is just another way of writing the abstract GCFG rule  $A \rightarrow f(A_1, \dots, A_\delta)$ .

### 6.2.2 Multiple Context-Free Grammars

Multiple Context-Free Grammars (Seki et al., 1991) were introduced in the late 80's as a very expressive formalism, incorporating Linear Context-Free Rewriting Systems and other mildly context-sensitive formalisms, but still with a polynomial parsing algorithm. MCFG is an instance of GCFG, with the following restrictions on linearizations:

- Linearization types are restricted to tuples of strings.

- The only allowed operations in linearization functions are tuple projections and string concatenations.

Since records can be seen as syntactic sugar for tuples, we can use records in this article without changing the definition of MCFG.

**Comparison with GF** Obviously MCFG is an instance of context-free GF, but without tables and table selections. The fact that GF can have nested records does not change anything – all nestings can be flattened. Also, an expanded table

$$[\mathbf{p}_1 \Rightarrow \phi_1; \dots; \mathbf{p}_n \Rightarrow \phi_n] \quad : \quad P \Rightarrow T$$

is equivalent to a record

$$\{\mathbf{p}_1 = \phi_1; \dots; \mathbf{p}_n = \phi_n\} \quad : \quad \{\mathbf{p}_1 : T; \dots; \mathbf{p}_n : T\}$$

and an instantiated selection  $\phi! \mathbf{p}_i$  is equivalent to a projection  $\phi. \mathbf{p}_i$ .

There are two fundamental differences:

- MCFG cannot have parameters as linearization values.
- A table selection in GF does not necessarily have to be an instantiated parameter; it can be any term of the correct linearization type.

In the example grammar, the linearization of  $p$  contains a selection  $y.s!(x.n)$ , which is not instantiated.

### 6.3 Converting GF to MCFG

In this section we show that it is possible to convert a GF grammar into an equivalent grammar where all table selections are instantiated, and containing no parameters. By the argument above, this means that context-free GF and MCFG are equivalent. This conversion is done in two steps, described later in sections 6.3.2 and 6.3.3, and the following theorem is a consequence.

**Theorem 23** *Any GF grammar with a context-free backbone can be reduced to an equivalent MCFG grammar.*

#### 6.3.1 Preliminaries

A linearization term  $\phi$  is in  $\eta$ -normal form if the structure follows the structure of its linearization type; i.e.  $\phi$  is a record if the type is a record type, and  $\phi$  is an expanded table if the type is a table type. The subterms of  $\phi$  which are of the basic linearization types,  $\text{Str}$  or  $\text{P} \in \text{Par}$ , are called the *leaves* of  $\phi$ . A *path* is a sequence of record projections and table selections; meaning that  $\epsilon$ ,  $\sigma.r$  and  $\sigma! \phi$  are paths if  $\sigma$  is a path. A path that does not contain any argument variables  $x_i$  is called instantiated; in which case the selections  $\phi$  can only be parameters. A non-instantiated path is called nested; this is because if

a path contains an argument variable  $x_i$ , then that variable is always followed by a (possibly empty) path.

A linearization type  $T$  as well as a linearization  $\phi$  can be partitioned into parameter paths and string paths:

$$\begin{aligned} T^{\text{Str}} &= \{ \sigma : \text{Str} \mid T.\sigma = \text{Str} \} \\ T^{\text{Par}} &= \{ \sigma : \text{P} \mid T.\sigma = \text{P} \in \text{Par} \} \\ \phi^{\text{Str}} &= \{ \sigma = \phi.\sigma \mid \phi.\sigma : \text{Str} \} \\ \phi^{\text{Par}} &= \{ \sigma = \phi.\sigma \mid \phi.\sigma : \text{P} \in \text{Par} \} \end{aligned}$$

Note that we equate nestings of tables/records and sets of path-value pairs, and that we extend paths to linearization types in the obvious way. Also note that there are only a finite number of instantiated parameter records  $\pi : T^{\text{Par}}$ , since there are only finitely many parameters.

**The example grammar** For the term  $a : \llbracket N \rrbracket$  in the example grammar we have that

$$\begin{aligned} \llbracket N \rrbracket^{\text{Str}} &= \{ s : \text{Str} \} \\ \llbracket N \rrbracket^{\text{Par}} &= \{ n : \text{Num} \} \\ a^{\text{Str}} &= \{ s = \text{“animals”} \} \\ a^{\text{Par}} &= \{ n = \text{Pl} \} \end{aligned}$$

### 6.3.2 A normal form for linearizations

**Definition 1** A GF linearization is in table normal form if it is of the form

$$f(x_1, \dots, x_\delta) = [\pi_1 \Rightarrow \phi_1; \dots; \pi_n \Rightarrow \phi_n]! \xi$$

and the following hold:

- $\xi$  contains all parameter paths of the arguments  $x_1, \dots, x_\delta$ ; in other words  $\xi = (x_1^{\text{Par}}, \dots, x_\delta^{\text{Par}})$ .
- Each  $\pi_k$  is a possible parameter instantiation of  $\xi$ ; in other words  $\pi_k : \llbracket A_1 \rrbracket^{\text{Par}} \times \dots \times \llbracket A_\delta \rrbracket^{\text{Par}}$ .
- Each  $\phi_k$  is in  $\eta$ -normal form where the leaves are either parameters or concatenations of constant strings and instantiated string paths.

The following algorithm converts any GF linearization in canonical form into normal form.

**Algorithm 1** First, add the outer table as in the definition of table normal form:

$$\begin{aligned} f(x_1, \dots, x_\delta) &= [\pi_1 \Rightarrow \phi; \dots; \pi_n \Rightarrow \phi]! \xi \\ \xi &= (x_1^{\text{Par}}, \dots, x_\delta^{\text{Par}}) \\ \pi_k &: \llbracket A_1 \rrbracket^{\text{Par}} \times \dots \times \llbracket A_\delta \rrbracket^{\text{Par}} \end{aligned}$$

Second, for each instantiation  $\pi_k$ , convert  $\phi$  to  $\phi_k$ , by repeating the following substitution until there are no parameter paths left:

- Substitute each instantiated parameter path  $x_i.\sigma$  for its  $\pi_k$ -instantiation  $(\pi_k)_i.\sigma$ .

**Lemma 24** *The algorithm, together with the standard computation rules, yields an equivalent linearization in table normal form.*

**The example grammar** There are two linearizations in the example that are not in table normal form, and this is how they look after conversion:

$$\begin{aligned}
 p(x, y) &= [\text{Sg} \Rightarrow \{s = x.s \cdot y.s ! \text{Sg}\}; \\
 &\quad \text{PI} \Rightarrow \{s = x.s \cdot y.s ! \text{PI}\}] ! x.n \\
 q(x, y) &= [\text{Sg} \Rightarrow \{s = [\text{Sg} \Rightarrow x.s ! \text{Sg} \cdot y.s ; \\
 &\quad \text{PI} \Rightarrow x.s ! \text{PI} \cdot y.s]\}; \\
 &\quad \text{PI} \Rightarrow \{s = [\text{Sg} \Rightarrow x.s ! \text{Sg} \cdot y.s ; \\
 &\quad \text{PI} \Rightarrow x.s ! \text{PI} \cdot y.s]\}] ! y.n
 \end{aligned}$$

### 6.3.3 Refining the abstract syntax

To get an MCF grammar, we have to get rid of the parameters in some way; and this we do by moving them to the abstract syntax. Each table row  $\pi_k \Rightarrow \phi_k$  above will then give rise to a unique function symbol with linearization  $\phi_k$ .

**Algorithm 2** *Given a GF grammar where all linearizations are in table normal form, create a grammar with the following categories, function symbols and linearizations:*

- For each category  $A$  and each instantiated parameter record  $\pi : \llbracket A \rrbracket^{\text{Par}}$ , create a new category  $\hat{A} = A[\pi]$ . The linearization type is the same as the string paths of the original linearization type,  $\llbracket \hat{A} \rrbracket = \llbracket A \rrbracket^{\text{Str}}$
- For each syntax rule  $A \rightarrow f(A_1, \dots, A_\delta)$ , and all new categories  $\hat{A}, \hat{A}_1, \dots, \hat{A}_\delta$ , create a new syntax rule  $\hat{A} \rightarrow \hat{f}(\hat{A}_1, \dots, \hat{A}_\delta)$ ; where  $\hat{f}$  is a unique function symbol,  $\hat{f} = f[\hat{A} \rightarrow \hat{A}_1 \dots \hat{A}_\delta]$ .
- For each linearization function

$$f(x_1, \dots, x_\delta) = [\pi_1 \Rightarrow \phi_1 ; \dots ; \pi_n \Rightarrow \phi_n] ! \xi$$



$$\begin{aligned}
 \llbracket \hat{N}_1 \rrbracket &= \llbracket \hat{N}_2 \rrbracket = \llbracket \hat{S} \rrbracket = \{ s : \text{Str} \} \\
 \llbracket \hat{V} \rrbracket &= \{ s ! \text{Sg} : \text{Str}; s ! \text{Pl} : \text{Str} \} \\
 \hat{p}_1 &: \hat{N}_1 \rightarrow \hat{V} \rightarrow \hat{S} \\
 \hat{p}_1(x, y) &= \{ s = x.s \cdot y.s ! \text{Sg} \} \\
 \hat{p}_2 &: \hat{N}_2 \rightarrow \hat{V} \rightarrow \hat{S} \\
 \hat{p}_2(x, y) &= \{ s = x.s \cdot y.s ! \text{Pl} \} \\
 \hat{q}_1 &: \hat{V} \rightarrow \hat{N}_1 \rightarrow \hat{V} \\
 \hat{q}_1(x, y) &= \{ s ! \text{Sg} = x.s ! \text{Sg} \cdot y.s; \\
 &\quad s ! \text{Pl} = x.s ! \text{Pl} \cdot y.s \} \\
 \hat{q}_2 &: \hat{V} \rightarrow \hat{N}_2 \rightarrow \hat{V} \\
 \hat{q}_2(x, y) &= \hat{q}_1(x, y) \\
 \hat{a} : \hat{N}_2 &= \{ s = \text{"animals"} \} \\
 \hat{b} : \hat{N}_1 &= \{ s = \text{"Bernie"} \} \\
 \hat{l} : \hat{V} &= \{ s ! \text{Sg} = \text{"loves"}; s ! \text{Pl} = \text{"love"} \}
 \end{aligned}$$

FIGURE 2 Grammar after conversion to MCFG

and each table row  $\pi_k \Rightarrow \phi_k$ , create a new linearization function for  $\hat{f} = f[\hat{A} \rightarrow \hat{A}_1 \dots \hat{A}_\delta]$ :

$$\begin{aligned}
 \hat{f}(x_1, \dots, x_\delta) &= \phi_k^{\text{Str}} \\
 \hat{A} &= A[\phi_k^{\text{Par}}] \\
 \hat{A}_i &= A_i[(\pi_k)_i] \quad (1 \leq i \leq \delta)
 \end{aligned}$$

where we by  $(\pi_k)_i$  mean the  $i$ th component of  $\pi_k$ .

Obviously the resulting grammar is an MCF grammar, since all linearizations are records of strings.

**Lemma 25** *The resulting grammar is equivalent to the original.*

**The example grammar** Figure 2 shows how the example grammar looks like after conversion to MCFG.

### 6.3.4 Non-deterministic reduction

There is a more direct conversion, using a non-deterministic substitution algorithm. This can also reduce the size of the resulting grammar, when argument parameters are not mentioned in linearizations.

**Algorithm 3** Assume the following abstract syntax rule, together with its linearization function:

$$\begin{aligned} A &\rightarrow f(A_1, \dots, A_\delta) \\ f(x_1, \dots, x_\delta) &= \phi \end{aligned}$$

Repeat the following non-deterministic substitution until there are no instantiated parameter paths left, accumulating the parameter records  $\pi_1, \dots, \pi_\delta$ :

- Substitute each instantiated parameter path  $x_i.\sigma : P$  with any  $p \in P$ , such that adding the row  $\sigma = p$  to  $\pi_i$  is consistent.

Supposing that the final substituted linearization is  $\psi$ , we can add the following rule for the new function symbol  $\hat{f}$ :

$$\begin{aligned} \hat{A} &\rightarrow \hat{f}(\hat{A}_1, \dots, \hat{A}_\delta) \\ \hat{f}(x_1, \dots, x_\delta) &= \psi^{\text{Str}} \\ \hat{A} &= A[\psi^{\text{Par}}] \\ \hat{A}_i &= A_i[\pi_i] \quad (1 \leq i \leq \delta) \end{aligned}$$

The algorithm is non-deterministic, and we get the final grammar by finding all solutions for each function symbol  $f$ ; which can be done by a standard all-solutions predicate, such as `findall` in Prolog.

**Coercions between categories** There is a difference between algorithm 3 and the previous algorithms; if an argument parameter  $x_i.\sigma$  is not mentioned in  $\phi$ , then there will be no  $\sigma$ -row in the constraint record  $\pi_i$ . This means that the new category  $\hat{A}_i = A_i[\pi_i]$  will only contain a subrecord of  $A_i[\phi_i^{\text{Par}}]$ , where  $\phi_i$  is a linearization of type  $\llbracket A_i \rrbracket$ .

**Algorithm 4** Given two reduced syntax rules,

$$\begin{aligned} \hat{A} &\rightarrow \hat{f}(\dots \hat{B}_1 \dots) \\ \hat{B}_2 &\rightarrow \hat{g}(\dots) \end{aligned}$$

where  $\hat{B}_1 = B[\pi_1]$  and  $\hat{B}_2 = B[\pi_2]$ . If  $\pi_1$  is a subrecord of  $\pi_2$ , add the coercion function  $\hat{c} = c[\pi_1 \pi_2]$ :

$$\begin{aligned} \hat{B}_1 &\rightarrow \hat{c}(\hat{B}_2) \\ \hat{c}(x) &= x \end{aligned}$$

**The example grammar** One function symbol gets a different linearization from algorithm 3 than in figure 2; the functions  $\hat{q}_1$  and  $\hat{q}_2$  get merged into one function  $\hat{q}$ .

$$\begin{aligned} \hat{q} &: \hat{V} \rightarrow \hat{N} \rightarrow \hat{V} \\ \hat{q}(x, y) &= \{ s! \text{Sg} = x.s! \text{Sg} \cdot y.s; \\ &\quad s! \text{Pl} = x.s! \text{Pl} \cdot y.s \} \end{aligned}$$

where  $\hat{N} = N[\ ]$ . This yields coercions for the more specific types  $\hat{N}_1 = N[n = \text{Sg}]$  and  $\hat{N}_2 = N[n = \text{Pl}]$ :

$$\begin{aligned} \hat{N} &\rightarrow \hat{c}_i(\hat{N}_i) \quad (i = 1, 2) \\ \hat{c}_i(x) &= x \end{aligned}$$

## 6.4 Implications to Parsing

**Definition 2** A chart for a GCFG is a finite set of tuples  $(f, \phi, \phi_1, \dots, \phi_\delta)$ , where  $\phi = f(\phi_1, \dots, \phi_\delta)$ .

A tree  $t = f(t_1, \dots, t_\delta)$  is represented by the chart if it contains  $(f, \llbracket t \rrbracket, \llbracket t_1 \rrbracket, \dots, \llbracket t_\delta \rrbracket)$ , and each subtree  $t_i$  is represented by the chart.

The following lemma is just another way of saying that GCFG grammars are compositional:

**Lemma 26** The set of GCFG trees  $\{ t \mid \llbracket t \rrbracket = \phi \}$ , for a given  $\phi$ , can be represented by a single chart.

In other words, a correct parsing algorithm for GCFGs does not have to return anything more than a chart.

If we have translated a context-free GF grammar into MCFG using algorithm 1+2, it is straight-forward to translate back a chart for the MCFG into a chart for the original grammar. Each item  $(\hat{f}, \phi, \phi_1, \dots, \phi_\delta)$ , where  $\hat{f} = f[\hat{A} \rightarrow \hat{A}_1 \dots \hat{A}_\delta]$  and  $\hat{A}_i = A_i[\pi_i]$ , can be converted to the item  $(f, \phi \cup \pi, \phi_1 \cup \pi_1, \dots, \phi_\delta \cup \pi_\delta)$ . Back-translation of trees are even simpler; just strip off the extra information from the nodes – each tree node  $\hat{f} = f[\dots]$  is converted to  $f$ .

If we have converted using algorithm 3+4, back-translation is only slightly more complicated; if there is a coercion  $\hat{A}_k \rightarrow \hat{c}(\hat{A}_k)$ , where  $\hat{A}'_k = A_k[\pi'_k]$ , use the linearization  $\phi_k \cup \pi'_k$  instead of  $\phi_k \cup \pi_k$ .

### 6.4.1 Dependent categories

If we have a GF grammar with dependent categories, there is a straight-forward two-step parsing process for that grammar. First we simply remove all dependencies from the abstract syntax, thereby getting a grammar with a context-free backbone. This grammar is over-generating, so when parsing we get a chart containing all parse trees we want, but perhaps also some unwanted trees.

The second step is to convert the chart into Horn clauses, which can be solved by any proof search, e.g. standard Prolog. This conversion is done one item at the time; suppose the following chart item:

$$(f, \phi, \phi_1, \dots, \phi_\delta)$$

where  $f$  has the following abstract typing:

$$f : (x_1 : A_1) \rightarrow \dots \rightarrow (x_\delta : A_\delta(x_1, \dots, x_{\delta-1})) \\ \rightarrow A(x_1, \dots, x_\delta)$$

From this we can create the following Horn clause (where  $t : A[\phi]$  is just syntactic sugar for a 3-tuple):

$$f(x_1, \dots, x_\delta) : A(x_1, \dots, x_\delta)[\phi] :- \\ x_1 : A_1[\phi_1], \dots, x_\delta : A_\delta(x_1, \dots, x_{\delta-1})[\phi_\delta]$$

Finally, the query  $:- x : S[\phi]$ , where  $x$  is a logic variable, will result in all possible parse trees  $x$  of category  $S$ , linearizing to the input string  $\phi$ .

### 6.4.2 Functional categories

In full GF, arguments to functions can themselves be functions. This gives rise to the question of how to linearize an “incomplete” category  $B_1 \rightarrow \dots \rightarrow B_\delta \rightarrow B$ . This is solved in GF by pairing the linearization of the result category  $B$  with linearizations of the *variable bindings* representing objects of category  $B_1, \dots, B_\delta$ .

Formally, each occurrence of a function category  $B_1 \rightarrow \dots \rightarrow B_\delta \rightarrow B$  as an argument in a typing is replaced by the new category  $\hat{B}$ , with linearization type

$$[[\hat{B}]] = [[B]] \times [[\text{Var}]]^\delta$$

where  $\text{Var}$  is a unique category for recognizing variable bindings, specified by the grammar. In GF, the default linearization type of variables is  $[[\text{Var}]] = \text{Str}$ , and it recognizes strings looking like ordinary mathematical variables (“ $x$ ”, “ $y$ ”, “ $z$ ”, ...).

For each new category  $\hat{B}$  we also need a coercion  $\hat{b}$ :

$$\hat{b} : B \rightarrow \text{Var} \rightarrow \dots \rightarrow \text{Var} \rightarrow \hat{B} \\ \hat{b}(x, y_1, \dots, y_\delta) = (x, y_1, \dots, y_\delta)$$

This conversion show that adding function arguments to abstract typings does not change the expressive power of GF, and that they are possible to handle with the parsing algorithms described in this paper.

## References

- Dymetman, Marc, Veronica Lux, and Aarne Ranta. 2000. XML and multilingual document authoring: Convergent trends. In *COLING, Saarbrücken, Germany*, pages 243–249.
- Hähnle, Reiner, Kristofer Johannisson, and Aarne Ranta. 2002. An authoring tool for informal and formal requirements specifications. In R.-D. Kutsche

- and H. Weber, eds., *Fundamental Approaches to Software Engineering*, vol. 2306 of *LNCIS*, pages 233–248. Springer.
- Hallgren, Thomas and Aarne Ranta. 2000. An extensible proof text editor. In M. Parigot and A. Voronkov, eds., *LPAR-2000*, vol. 1955 of *LNCIS/LNAI*, pages 70–84. Springer.
- Pollard, Carl. 1984. *Generalised Phrase Structure Grammars, Head Grammars and Natural Language*. Ph.D. thesis, Stanford University.
- Ranta, Aarne. 2004. Grammatical Framework, a type-theoretical grammar formalism. *Journal of Functional Programming* 14(2):145–189.
- Ranta, Aarne and Robin Cooper. 2004. Dialogue systems as proof editors. *Journal of Logic, Language and Information* To appear.
- Seki, Hiroyuki, Takashi Matsumara, Mamoru Fujii, and Tadao Kasami. 1991. On multiple context-free grammars. *Theoretical Computer Science* 88:191–229.
- Vijay-Shanker, K. David Weir, and Aravind Joshi. 1987. Characterizing structural descriptions produced by various grammatical formalisms. In *25th meeting of Association for Computational Linguistics*.