

# *Grammatical Framework*

## *A Type-Theoretical Grammar Formalism*

AARNE RANTA

*Department of Computing Science,  
Chalmers University of Technology and the University of Gothenburg,  
41296 Gothenburg, Sweden  
(e-mail: aarne@cs.chalmers.se)*

---

### **Abstract**

Grammatical Framework (GF) is a special-purpose functional language for defining grammars. It uses a Logical Framework (LF) for a description of abstract syntax, and adds to this a notation for defining concrete syntax. GF grammars themselves are purely declarative, but can be used both for linearizing syntax trees and parsing strings. GF can describe both formal and natural languages. The key notion of this description is a grammatical object, which is not just a string, but a record that contains all information on inflection and inherent grammatical features such as number and gender in natural languages, or precedence in formal languages. Grammatical objects have a type system, which helps to eliminate run-time errors in language processing. In the same way as a LF, GF uses dependent types in abstract syntax to express semantic conditions, such as well-typedness and proof obligations. Multilingual grammars, where one abstract syntax has many parallel concrete syntaxes, can be used for reliable and meaning-preserving translation. They can also be used in authoring systems, where syntax trees are constructed in an interactive editor similar to proof editors based on LF. While being edited, the trees can simultaneously be viewed in different languages. This paper starts with a gradual introduction to GF, going through a sequence of simpler formalisms till the full power is reached. The introduction is followed by a systematic presentation of the GF formalism and outlines of the main algorithms: partial evaluation and parser generation. The paper concludes by brief discussions of the Haskell implementation of GF, existing applications, and related work.

---

### **Capsule Review**

Type Theory provides logical foundation for functional programming, at the levels of higher order abstract syntax and semantics, both from the operational aspect with computational reduction rules, and from the denotational aspect with constructive realization structures. Concrete syntax however has usually a second-class status, through ad-hoc low-level parsers and printers. With his Grammatical Framework, Arne Ranta bridges this gap by providing a general computational linguistics formalism with dependent types where multilingual support is consistent with semantical design. This system has a powerful Haskell implementation with a pleasant structure editor. This reference paper is a comprehensive presentation of an original synthesis of recent research in proof theory, programming methodology and computational linguistics, which provides an elegant platform for literate programming.

---

## 1 Introduction: the goals of GF

The Grammatical Framework (GF) is a *grammar formalism*, i.e. a language for defining grammars. The development of GF started as a notation for *type-theoretical grammars* (Ranta, 1994), which use Martin-Löf's (1984) type theory to express the semantics of natural language. The first implementation was released in 1998 at Xerox Research Centre Europe in Grenoble, with focus on multilingual authoring via a type-theoretical pivot language. After the first publication (Mäenpää & Ranta, 1999), GF has developed into a functional programming language, whereby its notation has been completely revised, but it has preserved downward compatibility.

The focus of GF has thus shifted from an initial theoretical idea to practical applications. In this paper, we try to make explicit the theory that has proved useful for practical applications. The goal of GF is to serve both linguists, who want a high-level and reliable grammar formalism, and programmers, who want an elegant and efficient tool for building natural-language applications.

### 1.1 A logical framework with concrete syntax

When describing or implementing a language, it is customary to distinguish between its *abstract syntax*, i.e. the hierarchical structure of the language, and its *concrete syntax*, i.e. what the language looks like as it is read and written. The idea is that notions such as type checking and semantics are better defined on the level of abstract syntax, without the clutter of concrete syntax details.

In programming language design, concrete syntax is usually kept as simple as possible, with some concessions allowed to tradition, e.g. to include standard mathematical notations. In linguistics, the situation is different: the linguist has to take a natural language as it is, and describe it the best she can. It is not common to reach a level at which one can make a neat distinction between abstract and concrete syntax, or to discuss questions like type checking and semantics with the same precision as in programming languages.

GF was born from a synthesis between the computer science and linguistics ways of thinking: what about if we take an abstract syntax, with all type checking and semantics, and try to define a concrete syntax that looks exactly as we want, including natural languages? We took a powerful formalism for abstract syntax, a *Logical Framework* (LF),<sup>1</sup> and extended it with a notation for concrete syntax. With this formalism, it became possible to define all aspects of a language at once, which is an advantage for language implementation. At the same time, it gives a new perspective on natural languages, since it makes precise semantical notions applicable to those fragments of natural language that are recognized by GF grammars.

The practical issue of adding concrete syntax to LF was already addressed in the Mathematical Vernacular project of de Bruijn (1994). The goal was to make LF proof systems more accessible to users. An early implementation of this idea is a

<sup>1</sup> The logical framework of GF is a version constructive type theory, as are LF (Harper *et al.*, 1993), ALF (Magnusson & Nordström, 1994) and Coq (The Coq Development Team, 1999).

program that translates proofs in the logical framework Coq into an English-like notation (Coscoy *et al.*, 1995). However, the translation is part of the implementation of Coq itself: thus, even though the abstract syntax of new mathematical concepts is user-definable in Coq, their concrete syntax is not.

A step towards user-defined concrete syntax is taken in Isabelle (Paulson, 2002), which has a *mixfix* notation to define the concrete syntax of functions. Mixfix is a generalization of infix declarations – which of course are rudimentary concrete-syntax definitions – into full context-free rules. To some extent, mixfix notation is enough even for natural language. For instance, if we want to define an English notation for the length function, whose abstract syntax is given by the function declaration

$$\text{length} : (A : \text{Set}) \rightarrow \text{List } A \rightarrow \text{Int}$$

it is enough to write:<sup>2</sup>

$$\text{length } x = \text{``the''} ++ \text{``length''} ++ \text{``of''} ++ x$$

However, if we want to express *length* in correct German, we need to inflect it in different cases, put its argument into the dative case prefixed by the preposition *von*, and tell what gender it has. All this is done by a GF definition

$$\text{length } _ = fl \text{ ``Länge'' } Fem$$

where *fl* is a concrete syntax function taking care of the details of inflection and argument case. Of course, *fl* itself is user-defined,

$$fl : Str \rightarrow Gen \rightarrow \{s : Cas \Rightarrow Str ; g : Gen\} \rightarrow \{s : Cas \Rightarrow Str ; g : Gen\} = \\ \lambda F, G, x \rightarrow \{s = \mathbf{table} \{c \Rightarrow der Sg G c ++ F ++ \text{``von''} ++ x.s ! Dat\} ; g = G\}$$

This rule uses another user-defined function, *der*, which gives the inflectional forms of the German definite article.

Rules like the ones for *length* are typically written by persons who work in LF and are experts in the mathematical theories that they are working with. Functions like *fl* and *der* require expertise in linguistics and German grammar. A good division of labour is that functions of the latter kind are provided in *resource grammars* written by linguists and can be taken for granted by those who write applications. This requires a powerful notation permitting high levels of abstraction.

### 1.2 A grammar formalism with the linearization perspective

Both computer scientists and linguists have developed grammar formalisms – declarative descriptions of language from which language-processing algorithms can be automatically generated. The best-known algorithm is *parsing*, which takes strings into syntax trees (in the case of GF: to functional terms). Many grammar formalisms are designed to permit easy generation of parsers. The rules of GF, however, have their most direct readings in the direction of *linearization*, which takes functional

<sup>2</sup> We use GF notation; Isabelle mixfix does not support argument suppression.

terms into strings. To show that also a parser can be derived from every GF grammar requires a complicated argument.

In computer science, the best-known grammar formalisms are context-free grammars (=BNF) and attribute grammars (Knuth, 1968). BNF is used for describing languages in reports, but language implementations use extensions of the declarative format with *semantic actions* written in a general-purpose programming language. YACC (Johnson, 1975) is the classical model for such formalisms. Because of semantic actions, YACC grammars cannot generally be used for linearization. BNF does not cover the type-checking aspect of languages, whereas attribute grammars are able to do some of it. In YACC, some type checking can be performed in semantic actions, but this easily becomes a mess, and separate phases are recommended instead. In logical frameworks, type checking can be neatly included in abstract syntax definitions, but there are not yet any tools for doing so in practical language implementations.

In linguistics, the tradition closest to GF is Montague grammar (Montague, 1974), which uses simple type theory to express abstract syntax. In a way, GF is just a generic framework for implementing Montague-style grammars extended with dependent types.<sup>3</sup> Another point of reference are the grammar formalisms based on unification. These formalisms include DCG (Pereira & Warren, 1980) and PATR (Shieber, 1986), which are, like GF, pure frameworks, and HPSG (Pollard & Sag, 1994) and LFG (Bresnan, 1982), which have elements of built-in linguistic theory. While postponing the discussion of related work to section 10, it can be useful to point out the major ways in which GF differs from unification grammar formalisms:

- GF has separate rules for abstract and concrete syntax.
- The primary perspective in GF is linearization.
- GF has a strong type system.
- GF is a functional language.
- GF can integrate semantics with abstract syntax.
- GF supports multilingual grammars.

A common feature between GF and formalisms like PATR and HPSG is the use of records to model complex grammatical objects.

### 1.3 A multilingual authoring system

GF can be used in batch mode for linearization, parsing, and translation. A new kind of application is inherited from type-theoretical proof editors: syntax editing natural language. This is the application that explains Xerox's interest in the GF project: with a multilingual grammar, the user of GF can edit a document in a language that she does not know, while at the same time seeing how it evolves in her own language. This activity is called *multilingual authoring* (Power & Scott, 1998; Dymetman *et al.*,

<sup>3</sup> The semantic aspects of this extension are studied in Ranta (1994). A strong application of dependent types is the analysis of pronominal reference.

2000). For instance, a letter being edited in Swedish may look like this:

*Kära [Recipient],  
jag har äran att meddela dig att du har blivit befordrad till [Position].*

The expressions in brackets are *placeholders*, yet to be filled in to complete the letter. The possible fillings are either choices from menus (generated from the grammar) or strings of English or Swedish text (parsable by the grammar). In parallel with Swedish, an English version may be generated from the same abstract source:

*Dear [Recipient],  
I have the honour to inform you that you have been promoted to [Position].*

The author can be sure that both letters come out grammatically correct and convey the same message. The type-theoretical representation works as a *pivot language*, which controls the consistency of the document, and guarantees that all translations have the same meaning. Interaction eliminates a notorious problem of automatic translation, which is that a source text written in natural language generally does not fully determine the semantic content (cf. Kay 1997).

A multilingual GF grammar is a simple and efficient way to implement translation, which works surprisingly well, due to the long distance from concrete-syntax details that is possible in abstract syntax. Most other systems use separate *transfer rules* to translate between languages. The obvious disadvantage of transfer rules is that  $n(n - 1)$  transfer modules are needed to translate between  $n$  languages, whereas GF only needs  $n + 1$  grammar modules. The disadvantage of GF's pivot language method is that translation is limited to be structure-preserving. Even if there is no GF notation to express transfer rules, GF does not preclude them: the API module (section 8.4) gives support for defining transfer rules in Haskell.

## 2 Context-free rewrite grammars

This section starts a series of more and more powerful fragments of GF. Most GF concepts and applications already make sense in this first fragment, but only a limited class of GF grammars can be written in it. Sections 3 and 4 extend the abstract syntax part of GF, whereas section 5 independently extends the concrete syntax part.

### 2.1 From context-free grammars to context-free rewrite grammars

Context-free rewrite grammars are a generalization of context-free grammars, arising from a distinction between two aspects: abstract syntax and concrete syntax. Consider a context-free rule

$$f. C ::= t_1 \dots t_m$$

where  $f$  is the *rule label*,  $C$  is a *category symbol*, and each  $t_i$  is either a category symbol (a *nonterminal*) or a string (a *terminal*). The abstract syntax aspect of this rule is a *function declaration*, declaring  $f$  as a function whose value type is  $C$  and argument types are the nonterminals  $C_1 \dots C_n$  among  $t_1 \dots t_m$ :

$$\mathbf{fun} f : C_1 \rightarrow \dots \rightarrow C_n \rightarrow C$$

The concrete syntax aspect is a *linearization rule*, which assigns a *linear pattern* to the function  $f$ :

$$\mathbf{pattern} \ f \ x_1 \ \dots \ x_n = t_1 ++ \dots ++ t_m$$

Each  $t_i$  is either one of the variables  $x_j$  (a nonterminal) or a string (a terminal).

As an example of splitting a context-free rule into a function declaration and a linearization rule, consider the rule

$$\mathit{Div. Prop} ::= \mathit{Exp} \ \text{''is''} \ \text{''divisible''} \ \text{''by''} \ \mathit{Exp}$$

Splitting gives the rule pair

$$\begin{aligned} \mathbf{fun} \ \mathit{Div} &: \mathit{Exp} \rightarrow \mathit{Exp} \rightarrow \mathit{Prop} \\ \mathbf{pattern} \ \mathit{Div} \ x \ y &= x ++ \text{''is''} ++ \text{''divisible''} ++ \text{''by''} ++ y \end{aligned}$$

Abstract syntax rules alone define a system of *syntax trees*, i.e. functional terms formed by using rule labels as constants. Given two more function declarations,

$$\begin{aligned} \mathbf{fun} \ \mathit{two} &: \mathit{Exp} \\ \mathit{sum} &: \mathit{Exp} \rightarrow \mathit{Exp} \rightarrow \mathit{Exp} \end{aligned}$$

we can form the syntax tree

$$\mathit{Div} (\mathit{sum} \ \mathit{two} \ \mathit{two}) \ \mathit{two}$$

whose type is *Prop*. Given the linearization rules

$$\begin{aligned} \mathbf{pattern} \ \mathit{sum} \ x \ y &= \text{''the''} ++ \text{''sum''} ++ \text{''of''} ++ x ++ \text{''and''} ++ y \\ \mathit{two} &= \text{''two''} \end{aligned}$$

we have a correspondence between this tree and the string

$$\mathit{the \ sum \ of \ two \ and \ two \ is \ divisible \ by \ two}$$

## 2.2 *Permutation, suppression and reduplication*

To represent context-free rules in context-free rewrite grammar, the full expressive power of linear patterns is not needed, but only the special case in which the sequence of nonterminals in the linear pattern corresponds one-to-one to the arguments of the function. The full format extends this special case in three ways:

- **Permutation:** constituent order may be changed.
- **Suppression:** constituents may be hidden.
- **Reduplication:** constituents may be repeated.

Permutation is important if we want to have concrete syntaxes sharing an abstract syntax. For instance, it permits giving the same abstract syntax to prefix and infix notation, or to adjectival modification in English (prefix: *even number*) and French (postfix: *nombre pair*). Suppression is needed if we want a syntax tree to carry more information than the corresponding string, as in proof-carrying documents

(section 4.3). Reduplication shows that the formalism is more powerful than context-free grammars: for instance, the copy language of the universal language  $U$  (over some alphabet)

$$\{xx|x \in U\}$$

is context-sensitive, but it is encoded by the context-free rewrite grammar

$$\mathbf{fun} \ f : U \rightarrow S; \mathbf{pattern} \ f \ x = x \ ++ \ x$$

### 2.3 Linearization

To linearize a syntax tree

$$f \ a_1 \ \dots \ a_n$$

the linearization algorithm reads the linear pattern given in the rule

$$\mathbf{pattern} \ f \ x_1 \ \dots \ x_n = t_1 \ ++ \ \dots \ ++ t_m$$

and scans the sequence  $t_1 \ ++ \ \dots \ ++ t_m$  from left to right:

$$(f \ a_1 \ \dots \ a_n)^o = s_1 \ ++ \ \dots \ ++ s_m$$

where

$$s_i = \begin{cases} a_j^o & \text{if } t_i = x_j \text{ (the } j\text{'th argument)} \\ s & \text{if } t_i = s \text{ (string)} \end{cases}$$

The algorithm assumes that expressions are in the *full application form*, i.e. functions are endowed by all their arguments: otherwise it is in general not possible to fill the linear pattern in a meaningful way. To express functions, lambda abstraction must be used (section 3).<sup>4</sup>

### 2.4 Type checking and syntax editing

The type checker uses abstract syntax to look up the types of functions and verifies that they are used in accordance with their types. Rather than formulating the algorithm explicitly, we give the typing rule of syntax trees:

$$\frac{f : C_1 \rightarrow \dots \rightarrow C_n \rightarrow C \quad a_1 : C_1 \dots a_n : C_n}{f \ a_1 \ \dots \ a_n : C}$$

This is of course the same as the full application rule of typed lambda calculus.

Syntax editing takes place in a *state*, which consists of a tree being edited and the subtree that is the current *focus*. The tree may be *incomplete*, that is, contain *placeholders* (also called *metavariables*), which are yet to be filled by subtrees. An example is the following arithmetical proposition, where the focus is marked by an asterisk (\*) and placeholders by question marks (?):

$$\mathit{Div}(\mathit{sum} \ \mathit{two} \ *?) \ ?$$

<sup>4</sup> Full application is usually required in all functional languages whenever other kinds of functions are used than prefixes: one cannot write `(if b then 5 else)` to replace `(\x -> if b then 5 else x)`, and omitting two arguments would lead to the completely bizarre `(if b then else)`. One of the rare exceptions are Haskell's infix sections, such as `(4+)` (equal to `(\x -> 4 + x)`).

The linearization of this tree is

*the sum of two and (\* ?) is divisible by ?*

An important aspect of editing in GF is that it is possible to switch between the tree representation and its linearizations, even when the tree is incomplete.

For an efficient implementation of editing commands, we represent trees in a form in which the types of all subtrees are shown. The above tree is then represented as

<i>Div</i>	:	<i>Prop</i>
<i>sum</i>	:	<i>Exp</i>
	<i>two</i>	<i>Exp</i>
*	?	<i>Exp</i>
	?	<i>Exp</i>

The editor uses the typing rule for full applications to annotate each subtree with the value type of its function head.

The most important editing command is *refinement*, which replaces the metavariable in focus with a function. If the function takes arguments, like *sum*, refinement introduces new metavariables:

*the sum of two and (\* the sum of ? and ?) is divisible by ?*

The value type of the function must of course match the type of the focus metavariable. The editor guarantees this by maintaining a *menu* of type-correct refinements extracted from the grammar.

The editor focus is analogous to the cursor in a string editor: it marks the place to which editing commands apply. The analogue of cursor movements are *navigation commands*, which shift the focus without changing the tree. GF uses a *zipper* (Huet, 1997) to represent the editor state. As shown by Huet, navigation commands can be implemented efficiently for the zipper. Editing commands are efficient, as well, if the nodes locally contain all information that is needed when executing them – in particular, the types of all subtrees.

Refinement is a purely *top-down* editing command, inherited to GF from the proof editor ALF (Magnusson & Nordström, 1994). The zipper implementation has made it easy to add a generalized *bottom-up* command, the *local wrap*. The local wrap embeds the focus subtree, which need not be a metavariable, in a function application. If the focus is a tree

$t : A$

then any function

$f : \dots \rightarrow A \rightarrow \dots \rightarrow A$

can be used to replace  $t$  with  $(f ? \dots t \dots ?)$ ; i.e. one argument place of  $f$  is filled by  $t$ , and the other places by metavariables. For instance, wrapping the focus term of

*Div two (\*one)*

in the first argument place of  $sum : Exp \rightarrow Exp \rightarrow Exp$  results in the tree

*Div two (\*sum one ?)*

In the corresponding English sentence, the word *one* gets embedded in the phrase *the sum of one and ?*. Such an insertion is of course trivial in a text editor working on strings, but not available in most tree-based structure editors. Yet it is a functionality that document authors expect from an editor: they want to make local modifications in a document without having to destroy and rebuild parts of it.

The wrap operation is similar to *tree adjoining* in the grammar formalism TAG (Tree Adjoining Grammars) (Joshi, 1985). A special case is wrapping the top node, which is pure bottom-up editing: there we can relax the requirement that the value type of the function be the same as the type of the focus.

## 2.5 Parsing

While linearization in context-free rewrite grammars is straightforward, parsing is a search problem. We reduce it to parsing in context-free grammars, which has a complete solution by, for example, the Earley algorithm (Earley, 1970). Context-free parsing is completed by postprocessing that involves a rearrangement of subtrees.

**Translation into context-free rules.** To each pair of a typing judgement and a linearization rule

$$\begin{cases} \text{fun } f : C_1 \rightarrow \dots \rightarrow C_n \rightarrow C \\ \text{pattern } f x_1 \dots x_n = t_1 ++ \dots ++ t_m \end{cases}$$

we assign the context-free rule

$$f_p. C ::= c_1 \dots c_m$$

where

$$c_i = \begin{cases} C_j & \text{if } t_i = x_j \quad (\text{nonterminal}) \\ s & \text{if } t_i = s \quad (\text{terminal}) \end{cases}$$

The function  $f$  is indexed by a *profile*  $p$ , which is a list of lists of integers

$$[p_1, \dots, p_n]$$

where

$$p_i = [j \mid j \in \{1, \dots, k\}, t'_j = x_i]$$

where  $t'_1 \dots t'_k$  is the sequence of nonterminals in the **pattern** clause. In other words, each item  $p_i$  in the profile tells what places the  $i$ th argument occupies in the linear pattern. For instance, the pair of rules

$$\text{fun } f : A \rightarrow B \rightarrow C \rightarrow D; \text{pattern } f x y z = y ++ \text{"kuin"} ++ y ++ \text{"on"} ++ z$$

generates the rule

$$f_{[[], [1,2], [3]]}. D ::= B \text{"kuin"} B \text{"on"} C$$

**Postprocessing context-free parse trees.** A parse tree produced by the context-free parser may have a wrong number of arguments in wrong places and even in inconsistent ways (because of reduplication). The transformation  $R$  of parse trees

into proper functional terms is performed by reference to the profile of the function head:<sup>5</sup>

$$(f_p c_1 \dots c_m)^R = f a_1 \dots a_n$$

where

$$a_i = \begin{cases} c_k^R & \text{if } k \in p_i \text{ and } x_i \text{ is consistently represented} \\ ? & \text{if } p_i \text{ is empty} \end{cases}$$

If an argument is suppressed in linearization, this operation thus introduces a metavariable to represent it. If different occurrences of an argument are not represented consistently,<sup>6</sup> the operation of rearrangement fails.

**Context-free parsing.** The choice of context-free parsing algorithm is often the main efficiency issue when processing a language. Since GF grammars are implemented as first-class data objects in Haskell, well-known analyses and transformations (Hopcroft & Ullman, 1979) can be applied to them. For instance, even though the Earley algorithm is applicable to all grammars, some grammars may turn out to permit deterministic LR(1) parsing (Knuth, 1965), which can then be chosen for efficiency.

Some pathological rules cause problems for all context-free parsing algorithms. A cyclic rule

$$f. C ::= C$$

generates the infinite sequence of parse trees

$$t, f t, f(f t), f(f(f t)), \dots$$

for any tree  $t$  of type  $C$ . Cyclic rules are sometimes generated from innocent-looking GF rules, such as

$$\mathbf{fun} f : A \rightarrow C \rightarrow C; \mathbf{pattern} f x y = y$$

where the argument  $x$  is suppressed in linearization.

## 2.6 Semantic definitions and computation

The abstract syntax of a GF grammar can be thought of as a semantic model of the language, especially in those cases where GF is used as a syntactic annotation language for a logical framework. In such a model, we often want not only to declare functions but also to define them. To this end, GF has the form of judgement

$$\mathbf{def} f x_1 \dots x_n = t$$

where  $f : C_1 \rightarrow \dots \rightarrow C_n \rightarrow C$  and  $t : C$  in the context  $x_1 : C_1, \dots, x_n : C_n$ . An example is

$$\mathbf{def} \mathit{double} x = \mathit{sum} x x$$

<sup>5</sup> In context-free rewrite grammar, profiles are an optimization that make it possible to avoid lookups in the grammar at the postprocessing stage. In full GF, profiles are indispensable (cf. section 7.2).

<sup>6</sup> Consistent means here that the trees are the same. In a more general setting, it means that they are unifiable (cf. section 7.2).

Definitions are used for *computing* trees, and they define thereby a notion of *equality* between trees. This equality is called *semantic equality*, since it does not affect linearization: even though the terms *two* and *double one* are equal by definition, they are linearized as two different strings. We say that strings resulting from semantically equal trees are *paraphrases* of each other.

### 3 Variable bindings

Many interesting languages have variable-binding operations. For instance, predicate calculus has the universal quantifier  $\forall$  forming propositions such as

$$\forall x.x + 0 = x$$

where the variable  $x$  is bound in the subformula  $x + 0 = x$ . The context-free syntax of universally quantified propositions is

$$Prop ::= \text{''}\forall\text{'' } Var \text{''} \cdot Prop$$

This rule, however, does not capture the fact that the variable is bound in the subformula. The two parts of the structure, *Var* and *Prop*, are not constituents in the same sense: the *Prop* is an argument, whereas the *Var* is a binding. The distinction between arguments and bindings is fundamental for abstract syntax operations, such as type checking and computation – and in particular for syntax editing – and has to be stated somewhere, either by separate rules, or by using *higher-order abstract syntax* instead of context-free syntax. We will now explain how GF implements higher-order abstract syntax and corresponding concrete syntax.

#### 3.1 The abstract syntax of bindings

In higher-order abstract syntax, variable-binding operators are treated as functions that take functions as arguments. For instance, the universal quantifier has just one argument: a function from expressions to propositions:

$$\text{fun } Univ : (Exp \rightarrow Prop) \rightarrow Prop$$

This function has a second-order type. The general form of a type is now

$$A_1 \rightarrow \cdots \rightarrow A_n \rightarrow C$$

where each  $A_i$  is a type and  $C$  is a category.<sup>7</sup> Objects of function types can be formed by  $\lambda$ -abstraction, of the form

$$\lambda x_1, \dots, x_n \rightarrow b$$

We require that the number of  $\lambda$ -bound variables be the same as the number of argument types; this is known as the  $\eta$  *long normal form* of  $\lambda$ -terms.

For example, the formula  $\forall x.x + 0 = x$  has the syntax tree

$$Univ(\lambda x \rightarrow Eq(\text{sum } x \text{ Zero}) x)$$

<sup>7</sup> We could say:  $C$  is a type, but the resulting notion of type would be equivalent. As we formulate it now, we emphasize that each type has a basic type as its value type.

### 3.2 The concrete syntax of bindings

In the  $\eta$  long normal form, every subtree whose type is a function type has a bound variable for each of its argument types. We collect the variable symbols  $x_1, \dots, x_n$  and the linearization  $b$  of the body into a *record*

$$\{v_1 = x_1 ; \dots ; v_n = x_n ; s = b\}$$

which is the linearization of the whole tree. We use the record label  $s$  for the body, and the labels  $v_1, v_2, v, \dots$  for the variables; if there is just one variable, we use  $v$ .

For example, the linearization of  $\lambda x \rightarrow Eq(\text{sum } x \text{ Zero}) x$  is

$$\{v = "x" ; s = "x" ++ "++" ++ "++" ++ "0" ++ "=" ++ "x"\}$$

The linearization rule of the function *Univ* is

$$\mathbf{lin} \text{ Univ } P = \{s = "\forall" ++ P.v ++ "." ++ P.s\}$$

The general form of a linearization rule is now

$$\mathbf{lin} f x_1 \dots x_n = \{s = t_1 ++ \dots ++ t_m\}$$

where each  $t_i$  has one of the forms "*foo*" (terminal string),  $x_j.s$  (nonterminal body),  $x_j.v_k$  (nonterminal variable). The linearization algorithm distinguishes three cases:

$$\begin{aligned} \text{application: } (f a_1 \dots a_n)^o &= t(x_1 := a_1^o, \dots, x_n := a_n^o) \\ &\quad \text{if } \mathbf{lin} f x_1 \dots x_n = t \\ \text{abstraction: } (\lambda z_1 \rightarrow \dots \rightarrow \lambda z_n \rightarrow b)^o &= \{v_1 = 'z_1'; \dots ; v_n = 'z_n'\} ** b^o \\ \text{variable: } x^o &= \{s = 'x'\} \end{aligned}$$

This definition uses an operation  $**$  for conjoining records, and a substitution operation ( $x := a$ ). It also presupposes a symbol-printing operation producing a string ' $x$ ' from a variable symbol  $x$ .<sup>8</sup> The full normal form of linearization is obtained by these rules, substitutions, and the record projection rule

$$\{\dots ; r = t ; \dots\}.r = t$$

Using records instead of strings as values of linearization is crucial for maintaining the *compositionality* of linearization: for each function  $f$ , the linearization rule assigns a concrete-syntax function  $f'$  such that

$$(f a_1 \dots a_n)^o = f' a_1^o \dots a_n^o$$

where we denote the linearization of a tree  $t$  by  $t^o$ . Thus the linearization of a tree depends only upon the *linearizations* of its subtrees, not on the subtrees themselves. Compositionality guarantees a natural correspondence between abstract and concrete syntax. It also helps to make the implementation of GF efficient (section 5.4).

The **pattern** format used for context-free rewrite grammars is a special case of the **lin** format, as defined in section 6.4.2.

<sup>8</sup> One case is missing from this definition: variable applied to arguments. This case is needed if the abstract syntax uses third- or higher-order functions. GF then produces an *ad hoc* linearization where the symbol  $x$  is prefixed to the linearizations of the arguments in parentheses. In practice, functions of higher order than the second are rare in abstract syntax; strictly speaking, however, one should conclude that GF only supports user-defined concrete syntax for second-order abstract syntax.

### 3.3 Parsing bindings

From the context-free point of view, variable bindings are constituents: the parser reading the input looks for items of a certain shape (e.g. “ $x$ ”) that match a particular nonterminal  $Var$ . The context-free rule generated from  $Univ$  looks as follows:

$$Univ_{([\ ] , [2])} Prop ::= \text{“}\forall\text{” } Var \text{”} Prop$$

The right-hand side is as expected. The profile now contains information about where in the parse tree the bound variables are found: we extend the profiles of section 2.5 to lists of pairs  $(b, c)$  where  $c$  is a profile item in the old sense, telling what places the constituent occupies, and  $b$  is a list of items telling what places each of the bound variables occupy.

When producing context-free grammars, we introduce a category  $Var$  distinct from all categories in the context-free grammar, with some rules for recognizing variables. We add a rule

$$var_C. C ::= Var$$

for each category  $C$ . Finally, when postprocessing the parse tree, a suppressed binding (i.e. one with the profile item  $[\ ]$ ) is not replaced by a metavariable (?), but by a fresh ordinary variable. Notice that the unification phase of postprocessing is first-order, since variables in bindings are treated as ordinary arguments.

### 3.4 Type checking and syntax editing bindings

Type checking syntax trees with bindings is the same thing as monomorphic type checking in simply typed lambda calculus. In addition to the full application rule of section 2.4, we have the full abstraction rule

$$\frac{(x_1 : A_1, \dots, x_n : A_n) \quad c : C}{\lambda x_1, \dots, x_n \rightarrow c : A_1 \rightarrow \dots \rightarrow A_n \rightarrow C}$$

For syntax editing, we continue to use the zipper. Nodes are extended to contain the actual bindings. Each binding shows the type of the variable, which makes it easy to look up the type. For instance, a tree for the logical formula

$$\forall x. x = x$$

looks as follows:

$$\begin{array}{ll} Univ & : Prop \\ (x : Exp) Equal & : Prop \\ x & : Exp \\ x & : Exp \end{array}$$

## 4 Dependent types

*Dependent types* are types that depend on objects. We will give examples of two uses of them: first a grammar that defines the type checker of a small programming

language, and secondly a grammar of proofs, which leads to the notion of *proof-carrying documents*.

#### 4.1 Typed expressions

The following judgements define a category *Typ* of datatypes, and the category *Exp* of expressions, which depends on *Typ*.

$$\begin{array}{l} \mathbf{cat} \text{ } Typ \\ \text{ } Exp \text{ } Typ \end{array}$$

Examples of datatypes are integers, booleans, and lists:

$$\begin{array}{l} \mathbf{fun} \text{ } Int, Bool : Typ \\ \text{ } List : Typ \rightarrow Typ \end{array}$$

To define expressions, we use generalized function types where the value type depends on the argument:<sup>9</sup>

$$\begin{array}{l} \mathbf{fun} \text{ } Zero : Exp \text{ } Int \\ \text{ } True : Exp \text{ } Bool \\ \text{ } Nil : (A : Typ) \rightarrow Exp \text{ } A \\ \text{ } Cons : (A : Typ) \rightarrow Exp \text{ } A \rightarrow Exp (List \text{ } A) \rightarrow Exp (List \text{ } A) \\ \text{ } append : (A : Typ) \rightarrow (\_ \_ : Exp (List \text{ } A)) \rightarrow Exp (List \text{ } A) \end{array}$$

For example,

$$Cons \text{ } Int \text{ } Zero (Nil \text{ } Int)$$

is a valid syntax tree of type  $Exp (List \text{ } Int)$ , whereas

$$Cons \text{ } Int \text{ } Zero (Cons \text{ } Bool \text{ } True (Nil \text{ } Int))$$

is not a valid tree of any type. The grammar thus expresses not only the syntactic well-formedness of the language but also its well-typedness.

Notice how dependent types define the functions *Nil*, *Cons*, and *append* as *polymorphic*: their *Exp* arguments can be expressions of any types, in virtue of the *Typ* argument. The same technique is used in *monomorphic type theory* (Nordström *et al.*, 1990) and in the dependently typed programming language Cayenne (Augustsson, 1998). Polymorphism in concrete syntax results from argument suppression: for instance, the Haskell notation for lists is defined by the rules

$$\begin{array}{l} \mathbf{lin} \text{ } Nil\_ = \{s = "[]"\} \\ \text{ } Cons\_x \text{ } y = \{s = x.s ++ ":" ++ y.s\} \\ \text{ } append\_x \text{ } y = \{s = "(" ++ x.s ++ " ++ " ++ y.s ++ ")\} \end{array}$$

<sup>9</sup> In all variable-binding constructions of GF, the wildcard *\_* can serve as a bound variable, if the variable is not used.

When completed with semantic definitions,

$$\begin{aligned} \mathbf{def} \text{ append\_}(Nil\_)y &= y \\ \text{append } A(Cons\ a\ x)y &= Cons\ A\ a(\text{append } A\ x\ y) \end{aligned}$$

the grammar defines a complete parser, pretty-printer, type checker, interpreter, and syntax editor for this little language.<sup>10</sup>

## 4.2 Curry–Howard isomorphism

The definition of type-correct expressions in abstract syntax is a fairly simple application of dependent types. A more involved one, and originally the main motivation of logical frameworks, is to define logical calculi by formulating inference rules as declarations of proof functions. The structure is the same as with types and expressions: we have a basic type *Prop* of propositions and the dependent type *Proof**A* of proofs of a proposition *A*. The idea to treat propositions as types of proofs is known as the *Curry–Howard isomorphism*. An example is implication à la Martin-Löf (1984): the formation, introduction and elimination rules come out as follows:

$$\begin{aligned} \mathbf{fun} \text{ Impl} &: Prop \rightarrow Prop \rightarrow Prop \\ \text{ImplI} &: (A, B : Prop) \rightarrow (Proof\ A \rightarrow Proof\ B) \rightarrow Proof(\text{Impl}\ A\ B) \\ \text{ImplE} &: (A, B : Prop) \rightarrow Proof(\text{Impl}\ A\ B) \rightarrow Proof\ A \rightarrow Proof\ B \end{aligned}$$

If we now want to express formal proofs in natural language, we simply give linearization rules that produce texts, e.g.

$$\begin{aligned} \mathbf{lin} \text{ ImplI } A\ B\ b &= \{s = \text{"assume"} ++ A.s ++ \text{"."} ++ b.s ++ \text{"."} ++ \\ &\quad \text{"Hence"} ++ \text{"if"} ++ A.s ++ \text{"then"} ++ B.s\} \end{aligned}$$

On the top level of mathematical texts, we use a category *Text* for textual units such as theorems with or without proofs:

$$\begin{aligned} \mathbf{fun} \text{ ThmProof, ThmOmit} &: (A : Prop) \rightarrow Proof\ A \rightarrow Text \\ \mathbf{lin} \text{ ThmProof } A\ a &= \{s = \text{"Theorem."} ++ A.s ++ \text{"Proof."} ++ a.s ++ \text{"QED"}\} \\ \text{ThmOmit } A\ _ &= \{s = \text{"Theorem."} ++ A.s ++ \text{"Proof."} ++ \text{"Omitted."}\} \end{aligned}$$

The typing of these functions forces the proof in *ThmProof* really to prove the theorem; in *ThmOmit*, a proof must exist even though it is not shown. Thus anyone who uses the GF syntax editor to build proof texts is forced to making them correct.

## 4.3 Proof-carrying documents

Besides mathematical texts, dependent types and the Curry-Howard isomorphism are useful for other kinds of texts, to guarantee semantic properties. Consider, for

<sup>10</sup> In GF's parameter system (section 5), we could moreover define a precedence parameter to regulate the use of parentheses.

instance, texts describing train connections:

*To get from Gothenburg to Hamburg, first take train 487 to Copenhagen and then change to train 36.*

The semantic well-formedness conditions for this text are that train 487 runs from Gothenburg to Copenhagen, that train 36 runs from Copenhagen to Hamburg, and that train 487 arrives at Copenhagen before train 36 leaves. All of these conditions can be concisely expressed by a grammar with a dependent type *Train A B* of train connections from the city *A* to the city *B*, and a type of proofs of the fact that one train arrives before another train leaves. New connections are generated by the rule

$$\text{fun } \textit{Connect} : (A, B, C : \textit{City}) \rightarrow (a : \textit{Train A B}) \rightarrow (b : \textit{Train B C}) \rightarrow \\ \textit{Before A B C a b} \rightarrow \textit{Train A C}$$

It is easy to write a linearization rule for *Connect* generating texts like the example above. Linearization hides the proof of the *Before* condition, but anyone who uses the GF syntax editor to build the text is obliged to give a proof in order for the text to be complete. We call this idea *proof-carrying documents*, with a reference to proof-carrying code (Necula, 1997).

#### 4.4 Concrete syntax and dependent types

Little need be said about the concrete syntax of dependent types, since linearization rules look precisely the same as without them. In parsing rules, the arguments of dependent categories are just ignored. The context-free parsing phase thus ignores type dependencies. It accepts ill-formed expressions such as “0 : *True* : []”, which the subsequent type checking phase rejects.

It is possible to improve the performance of the GF parser by integrating parsing and type checking: errors are then detected at an earlier stage. Some amount of integration is necessary if the grammar has syntactically dummy *coercion rules* like

$$\text{fun } \textit{coerce} : \textit{Exp Int} \rightarrow \textit{Exp Float} ; \text{lin } \textit{coerce } x = x.$$

The corresponding context-free rule is cyclic,

$$\textit{coerce. Exp} ::= \textit{Exp},$$

and produces an infinity of parse trees, at most one of which is type-correct.<sup>11</sup>

#### 4.5 Type checking and syntax editing dependent types

Type checking with dependent types is harder than without them, since it involves computation of expressions. For instance, a proof that 2 is even is also a proof that  $1 + 1$  is even. In the presence of variables and metavariables, moreover, it cannot always be decided if an expression has a given type: whether a proof of *Even 2* is

<sup>11</sup> Luo and Callaghan (1999) investigate coercion as a central phenomenon of informal mathematical language and suggest an algorithm for resolving it.

also a proof of  $Even(1+?)$  depends on the value of  $?$ . Therefore, what the type checker returns is not a boolean value but a set of *constraints*, which are equalities between terms. The value *True* corresponds to the empty set of constraints. The value *False* corresponds to the situation where some of the constraints is impossible, e.g.  $1 = 0$  (cf. Magnusson & Nordström, 1994).

In general, constraints contain metavariables that appear in different positions in the tree. Because of this, metavariables have to carry unique identifiers; we use subindexed question marks for this. For instance, parsing the expression

$$0 : 1 : []$$

in the grammar of section 4.1 creates an incomplete term of an incomplete type:

$$Cons ?_0 Zero(Cons ?_1 One(Nil ?_2)) : Exp ?_3$$

The type checker can easily find out the following constraints:

$$?_0 = ?_1 = ?_2 = Int, ?_3 = List ?_0$$

In this example, a simple constraint-solving mechanism is enough to automatically instantiate all the metavariables. Such is usually the case for hidden type arguments corresponding to polymorphism: the user of an editor does not need to fill in these arguments.

Even if constraints remain unsolved, they can be helpful in syntax editing. since they narrow down available choices. For instance, in a menu of refinements for  $? : Exp Bool$ , functions whose value type is  $Exp Int$  are not shown.

Formal rules for dependent types are given in section 6.2. As type checking algorithm, we have used the one in (Coquand, 1996), which we have modified so that it type-annotates terms into trees used by the zipper editor. In addition to bindings, function body, and value type, as in section 3.4, the information stored in a node includes the constraints created when type checking that node.

## 5 Extending concrete syntax

The values returned by linearization have so far been strings and records of strings. In this section, we generalize this to records that may also contain string-valued finite functions, *tables*, as well as *parameters*. Different concrete syntaxes may use different record types for one and the same type in the abstract syntax. This extension is essential to keep abstract syntax independent of language-dependent features such as inflection.

### 5.1 Parameters, tables, and records

A *parameter type* is a finite set of *parameter values*, on which the linearization of an expression may depend. For example, the type of grammatical numbers in English has two values: the singular and the plural. Expressions linearized as English common nouns have two forms: e.g. *Int* has the singular form *integer* and the plural form *integers*.

We write

$$\mathbf{param} \text{ Num} = \text{Sg} \mid \text{Pl}$$

to define the parameter type *Num*. The type

$$\text{Num} \Rightarrow \text{Str}$$

is the type of string-valued tables on *Num*, and the expression

$$\mathbf{table} \{ \text{Sg} \Rightarrow \text{"integer"} ; \text{Pl} \Rightarrow \text{"integers"} \}$$

gives such a table in explicit form. The *selection* operation (!) is used for applying tables to arguments:

$$\mathbf{table} \{ \text{Sg} \Rightarrow \text{"integer"} ; \text{Pl} \Rightarrow \text{"integers"} \} ! \text{Pl} = \text{"integers"}$$

Concrete syntax assigns to every category in abstract syntax a *linearization type*: for instance, the linearization type of *CN* is given by the judgement

$$\mathbf{lincat} \text{ CN} = \{ s : \text{Num} \Rightarrow \text{Str} \}$$

Linearization rules for expressions of the category *CN* must have this value type. An example of such a rule is

$$\mathbf{lin} \text{ Int} = \{ s = \mathbf{table} \{ \text{Sg} \Rightarrow \text{"integer"} ; \text{Pl} \Rightarrow \text{"integers"} \} \}$$

In German, common nouns do not only depend on number, but also on case (Nominative, Accusative, Genitive, Dative). The linearization type thus has a two-argument table, which we “curry” into a table of tables:

$$\mathbf{lincat} \text{ CN} = \{ s : \text{Num} \Rightarrow \text{Case} \Rightarrow \text{Str} \}$$

In yet other languages, there may be three numbers (Arabic has the dual) or fifteen cases (Finnish has – well...). In context-free (rewrite) grammars, all this variation would have to be expressed by unrelated rules, which would make it impossible to use a common abstract syntax.

In addition to parameters that produce different forms, expressions may have parameters as *inherent features*. For instance, German common nouns have a gender (Masculine, Feminine, Neuter) associated to them, but not as inflection forms: any noun inherently has just one gender. Inherent features are expressed by record fields in linearization types. Here is an amended rule for German common nouns:

$$\mathbf{lincat} \text{ CN} = \{ s : \text{Num} \Rightarrow \text{Case} \Rightarrow \text{Str} ; g : \text{Gen} \}$$

The record linearizing a tree contains all linguistic information concerning the expression: its inflection table and its inherent features. Such information is what we normally find in dictionaries. In grammar, this information is not only needed for individual words, but for arbitrarily complex phrases. For instance, in English, when a common noun is modified by an adjective, the resulting complex common noun can still be inflected in number:

$$\mathbf{fun} \text{ Mod} : \text{Adj} \rightarrow \text{CN} \rightarrow \text{CN}$$

$$\mathbf{lin} \text{ Mod } F \text{ A} = \{ s = \mathbf{table} \{ \text{Sg} \Rightarrow F.s ++ A.s ! \text{Sg} ; \text{Pl} \Rightarrow F.s ++ A.s ! \text{Pl} \} \}$$

All **fun** rules must have **lin** rules of matching linearization types. This can be checked at compile time, before the grammar is used. It is also easy to check that the rules are complete – essentially, that all tables have values for all elements of their argument types. Using a type system to prevent run-time errors is one of the key ideas that GF has inherited from functional programming languages.

### 5.2 Hierarchical parameters

Parameter types are like data types in Haskell and other functional languages, with the restriction that they must be finite. *Hierarchical* parameter types are permitted, and they are, in fact, often very appropriate. To give an example, French verbs, as presented in the authoritative *Bescherelle* (1997), have three persons, two numbers, two genders, four (non-composite) tenses, and six modes. But the inflection tables display only 51 (non-composite) verb forms, not 288, which would be the case if the forms were simply cross-products of all parameters. The reason is that many combinations do not exist. A natural way of describing this parameter system is by using parameter types whose constructors have arguments from other parameter types. The following system is a straightforward GF formalization of the *Bescherelle*:<sup>12</sup>

```

param Nombre = Sg | Pl
      Personne = P1 | P2 | P3
      Genre = Masc | Fem
      Temps = Pres | Imparf | Passe | Futur
      TSubj = SPres | SImparf
      TPart = PPres | PPasse Genre Nombre
      NImper = SgP2 | PIP1 | PIP2
      VForm = Inf | Indic Temps Nombre Personne | Cond Nombre Personne
              | Subj TSubj Nombre Personne | Imper NImper | Part TPart

```

### 5.3 Discontinuous constituents

In all examples so far, the linearization of a tree has been a record with one principal string or string-valued table, stored in a field labelled *s*. Now we consider cases where the linearization consists of separate parts, which can change order and get other expressions inserted between them. Such expressions are called *discontinuous constituents*.

A famous example of discontinuous constituents is German verb phrases. A verb phrase is a complex expression consisting of a verb and its complements, in English e.g. *loves Mary* in *John loves Mary*. The analysis of a sentence into a noun phrase (the

<sup>12</sup> We have used this type system in a complete GF implementation of the *Bescherelle* conjugations; see GF Homepage (Ranta, 2002). Huet (2000) uses CAML datatypes in the same way in his morphology of Sanskrit.

subject) and a verb phrase (the predicate) is motivated both by logic (Aristotelian or modern) and by linguistic facts such as the conjunction *John loves Mary and hates Bill*. In German, however, the verb phrase (*liebt Maria*) cannot be found in all uses of the sentence (*Johann liebt Maria*). For instance, in the conditional

*wenn Johann Maria liebt, liebt Johann Maria*

the verb phrase is used with a reverse word order in the antecedent, and dissolved into two parts in the succedent. Those linguists who still believe that *liebt Maria* is a constituent of the sentence, have to treat it as a discontinuous constituent.

In GF, discontinuous constituents are records with more than one string-valued fields. The German linearization type of verb phrases can be defined as

$$\mathbf{lincat} \ VP = \{s_1 : Agr \Rightarrow Str ; s_2 : Str\}$$

consisting of the verb part  $s_1$  and the complement part  $s_2$ . The verb part depends on agreement features, such as number and person. It receives them from the subject of the sentence, which has them as inherent features. The sentence-forming predication rule

$$\mathbf{fun} \ Pred : NP \rightarrow VP \rightarrow S$$

is linearized under a three-valued parameter that produces different strings for direct, inverse, and subordinate sentences:

$$\mathbf{lin} \ Pred \ N \ V = \{s = \mathbf{table} \ \{Dir \Rightarrow N.s ++ V.s_1 ! N.a ++ V.s_2 ; \\ Inv \Rightarrow V.s_1 ! N.a ++ N.s ++ V.s_2 ; \\ Sub \Rightarrow N.s ++ V.s_2 ++ V.s_1 ! N.a\}\}$$

The complementation rule forms a verb phrase from a transitive verb (*TV*) and a noun phrase:

$$\mathbf{fun} \ Compl : TV \rightarrow NP \rightarrow VP$$

$$\mathbf{lin} \ Compl \ V \ N = \{s_1 = V.s ; s_2 = N.s ! Acc\}$$

Given the noun phrases *Johann* and *Maria* and the transitive verb *Lieben*, we can form the syntax tree

*Pred Johann(Compl Lieben Maria)*

which has a linearization producing three forms: *Johann liebt Maria*, *liebt Johann Maria*, and *Johann Maria liebt*.

Parsing discontinuous constituents will be explained as a part of the full GF parsing algorithm (section 7.2). We just notice that discontinuous constituents make it possible to define intricate non-context-free languages, such as

$$\{a^n b^n c^n \mid n = 1, 2, \dots\}$$

This language is defined by the category  $S$  of the following GF grammar:

**cat**  $S ; Aux$

**fun**  $exp : Aux \rightarrow S ; first : Aux ; next : Aux \rightarrow Aux$

**lineat**  $Aux = \{s_1 : Str ; s_2 : Str ; s_3 : Str\}$

**lin**  $exp x = \{s = x.s_1 ++x.s_2 ++x.s_3\}$

$first = \{s_1 = "a" ; s_2 = "b" ; s_3 = "c"\}$

$next x = \{s_1 = "a" ++x.s_1 ; s_2 = "b" ++x.s_2 ; s_3 = "c" ++x.s_3\}$

#### 5.4 Canonical GF

We have extended context-free grammars into grammars where linearizations of syntax trees are records of tables of grammatical objects. The operational semantics of these grammars will be explained in terms of computation rules for table selections and record projections in section 6.3. From this perspective, linearization is similar to evaluation in a functional programming language.

However, GF has a computational model that is simpler than evaluation in functional language, since it does not involve substitutions for variables. The only variables that are present in the right-hand-side  $t$  of a linearization rule

$$\mathbf{lin} \ f \ x_1 \ \dots \ x_n = t$$

are  $x_1 \dots x_n$ , which stand for the linearizations of the arguments of  $f$ . The substitution of values for these variables can be performed in the same way as selections and projections: as look-up followed by simple replacement. Linearization as a whole is a single inorder traversal of the syntax tree.

We refer to the GF concrete-syntax notation so far introduced as *canonical GF*. In the next section, we will go far beyond canonical GF by adding functions and pattern matching. This extension is important for the usability of GF. For the implementation, however, the important thing is that the rich notation can be compiled back into canonical GF. Even though linearization could be performed directly as evaluation on the rich notation, it is much more efficient to perform *partial evaluation* on the grammar and use canonical GF at runtime. Moreover, it is from canonical GF that parsers are derived. Partial evaluation and parsing will be explained in section 7.

#### 5.5 Abstraction mechanisms

Linguists, just like functional programmers, like to work with strong generalizations and on a high level of abstraction. GF makes accessible to linguists two abstraction mechanisms of functional programming: *function definitions* and *pattern matching*.

Function definitions in GF are called *operation definitions* to distinguish them from the **fun** judgements of abstract syntax. An example is the operation that produces

regular common nouns in English:

**oper**  $regCN : Tok \rightarrow \{s : Num \Rightarrow Str\} = \lambda c \rightarrow \{s = Sg \Rightarrow c ; Pl \Rightarrow c + "s"\}$ <sup>13</sup>

The linearization rule of the datatype expression *Int* can now be concisely written

**lin**  $Int = regCN "integer"$

Pattern matching is used in tables: branches can be defined not only for constructor expressions, but also for *patterns*, which may contain variables and wildcards (-). For instance, the following table defines the English adjectival modification rule by using a pattern variable *n* for number:

**lin**  $Mod F A = \{s = \mathbf{table} \{n \Rightarrow F.s ++ A.s ! n\}\}$

It is possible to expand this table into the fully explicit form shown in section 5.1; using patterns, however, captures the generalization that it is the noun part that receives the number of the whole phrase.

Function types ( $A \rightarrow B$ ) and table types ( $A \Rightarrow B$ ) have many common properties: both allow currying, full and partial application, and formation by abstraction. There are important differences, however:

- Tables, but not functions, are restricted to finite argument types.
- Tables, but not functions, can be formed by case analysis.<sup>14</sup>
- Tables, but not functions, are values in canonical GF.

The partial evaluation algorithm (section 7.1) shows in fact that

- Functions can always be eliminated from linearization rules.

### 5.6 Resource grammars

The intended use of GF is to build natural-language fragments on top of semantic models, such as mathematical theories. This makes it possible to minimize the size of grammars and avoid many irrelevant linguistic problems. For instance, a French grammar for mathematics does not need to define all the 51 French verb forms, but two is often enough.

However, the *ad hoc* way of defining grammars may lead to duplication of work: if different parts of verb conjugation are needed in different applications, one cannot use the conjugation defined for one grammar as a resource for another grammar. And, of course, this style of grammar-writing favours linguistically unmotivated solutions.

The idea of *resource grammars* is to define common and unproblematic parts of concrete syntax – such as inflection tables – independently of abstract syntax. Using a resource grammar needs some care, however: a grammar encoding 51 forms of one thousand verbs is a heavy tool for actually dealing with two forms of ten verbs.

<sup>13</sup> We use + instead of ++ between strings to say that they belong to the same token (*Tok*).

<sup>14</sup> If the argument type of a function *f* is a parameter type, case analysis is of course possible in the form  $f = \lambda x \rightarrow \mathbf{table} \{...\} ! x$ .

When used in a naïve way, it produces enormous runtime systems. What makes resource grammars practical is type-driven partial evaluation (section 7.1). Suppose we only need two forms of French verbs – say, the indicative and subjunctive of third person singular present tense. The linearization type of verbs is then

$$\mathbf{lincat} \text{ Verb} = \{s : \text{Mode} \Rightarrow \text{Str}\}$$

Assume that we have a resource grammar with complete forms of conjugation à la *Bescherelle*, in a form like

$$\mathbf{oper} \text{ tenir} : \text{VForm} \Rightarrow \text{Str} = t$$

where *VForm* is the parameter type with 51 values. When evaluated, *t* yields the full table for the verb *tenir* (cf. section 5.2). Now, to use objects of this type as linearizations of verbs in the category *Verb*, all we need is an interface operation

$$\mathbf{oper} \text{ useVerb} : (\text{VForm} \Rightarrow \text{Str}) \rightarrow \{s : \text{Mode} \Rightarrow \text{Str}\} = \lambda t \rightarrow \\ \{s = \mathbf{table} \{ \text{Ind} \Rightarrow t ! (\text{Indic Pres Sg P3}) ; \text{Sub} \Rightarrow t ! (\text{Subj Pres Sg P3}) \} \}$$

Linearizations can then be defined compactly, for instance,

$$\mathbf{lin} \text{ Tenir} = \text{useVerb} \text{ tenir}$$

and the result is a two-element table with the forms *tient*, *tienne*, since the term is  $\eta$ -expanded with respect to the expected linearization type and then evaluated. If some other set of forms is needed, all that has to be changed is the definition of the interface operation *useVerb*.

Resource grammars are an obvious way to define morphology and lexicon, and they can often be compiled from existing resources or created by using general-purpose programming languages. But the idea also makes sense for syntax. For instance, the German linearization types for sentences, noun phrases, verb phrases, and transitive verbs, and the predication and complement rules (section 5.3) can be written as operations:<sup>15</sup>

$$\begin{aligned} \mathbf{oper} \text{ S} & : \text{Type} = \{s : \text{Ord} \Rightarrow \text{Str}\} \\ \text{NP} & : \text{Type} = \{s : \text{Case} \Rightarrow \text{Str} ; n : \text{Agr}\} \\ \text{VP} & : \text{Type} = \{s_1 : \text{Agr} \Rightarrow \text{Str} ; s_2 : \text{Str}\} \\ \text{TV} & : \text{Type} = \{s : \text{Agr} \Rightarrow \text{Str}\} \\ \text{Pred} & : \text{NP} \rightarrow \text{VP} \rightarrow \text{S} = \dots \\ \text{Compl} & : \text{TV} \rightarrow \text{NP} \rightarrow \text{VP} = \dots \end{aligned}$$

More functions can be defined in terms of these basic operations:

$$\begin{aligned} \text{Pred1} & : \text{VP} \rightarrow \text{NP} \rightarrow \text{S} = \lambda F, x \rightarrow \text{Pred } x \text{ F} \\ \text{Pred2} & : \text{VP} \rightarrow \text{NP} \rightarrow \text{NP} \rightarrow \text{S} = \lambda F, x, y \rightarrow \text{Pred } x (\text{Compl } F \text{ } y) \end{aligned}$$

<sup>15</sup> If the linguist prefers to write her grammar using **fun**, **cat**, **lin**, and **lincat**, as in Section 5.3, the **oper** definitions can be extracted automatically from them.

The writer of an application grammar, e.g. for mathematics can use these operations without knowing anything about German word order and agreement. If she has decided that propositions are linearized as *S*, one-place predicates as *VP*, and two-place predicates as *TV*, all she has to know is which verbs (from the resource grammar) are used for each predicate.<sup>16</sup> For instance, to linearize the one-place convergence predicate and the two-place intersection predicate, she writes

$$\begin{aligned} \mathbf{lin} \textit{ Converge} &= \textit{Pred1 konvergieren} \\ \textit{Intersect} &= \textit{Pred2 schneiden} \end{aligned}$$

In this way, a division of labour is achieved between authors of resource grammars, who are experts in linguistic rules, and authors of application grammars, who are experts in the domain of application.

## 6 The GF language

This section gives a concise definition of the GF formalism. The framework-level notions of type checking and evaluation are specified by inference rules. The notation we use is exactly the same as the notation recognized by the GF parser, with the exception of a handful of non-ASCII symbols: in ASCII-written GF source code, we replace  $\lambda$  by  $\backslash$ ,  $\rightarrow$  by  $->$ , and  $\Rightarrow$  by  $=>$ .<sup>17</sup>

### 6.1 Grammars and judgements

A GF grammar is a sequence of *judgements*. Judgements are divided into two sorts: those of abstract syntax and those of concrete syntax. Figure 1 shows the forms of judgement used in GF grammars, together with their verbal readings.

Every form of judgement has a *keyword* (such as **cat**, **param**). Every judgement ends with a semicolon (;), which we usually omit in typeset text, where we have access to layout. Using the semicolon (or layout) makes it possible to omit keywords: once a keyword appears in the code, it is read as the first word of every semicolon-separated judgement, until a keyword is encountered again.

The forms of judgement shown in Figure 1 are the ones that may appear in GF grammars. On the metalevel, we also use judgements of the forms

$$\begin{aligned} A : \textit{Type} & \quad A \textit{ is a type} \\ a : A & \quad a \textit{ is an object of type } A \\ a = b & \quad a \textit{ is definitionally equal to } b \end{aligned}$$

with their usual Logical Framework meanings (in e.g. (Nordström *et al.*, 1990)).

<sup>16</sup> Including adjectives as possible linearizations of predicates would not be a problem: the adjective-verb distinction could be hidden in slightly more general linearization types.

<sup>17</sup> Some GF structures that are supported by the actual implementation are left out, since we consider them experimental; we refer to documentation in Ranta (2002) for such features.

Abstract syntax.

<b>cat</b> $C \Gamma$	$C$ is a category depending on the context $\Gamma$
<b>fun</b> $f : A$	$f$ is a function of type $A$
<b>def</b> $a = b$	$a$ is defined as $b$
<b>data</b> $C = f_1 \mid \dots \mid f_n$	$C$ has the constructors $f_1, \dots, f_n$

Concrete syntax.

<b>param</b> $P = C_1 \Gamma_1 \mid \dots \mid C_n \Gamma_n$	$P$ is a parameter type with the constructors $C_1$ with context $\Gamma_1, \dots, C_n$ with context $\Gamma_n$
<b>linat</b> $C = L$	$C$ has the linearization type $L$
<b>lindf</b> $C = t$	$C$ has the default linearization $t$
<b>lin</b> $f = t$	$f$ has the linearization function $t$
<b>oper</b> $h : T = t$	$h$ is an operation of type $T$ , defined as $t$

Syntactic sugar: omitting keywords.

$$\mathbf{key} J ; \dots ; K \equiv \mathbf{key} J ; \dots ; \mathbf{key} K$$

Fig. 1. Forms of judgement in GF.

## 6.2 Abstract syntax

### 6.2.1 Categories, types, and functions

Judgements of **cat** and **fun** forms are used for building basic types and objects. The **cat** judgement

$$\mathbf{cat} C \Gamma$$

presupposes that  $\Gamma$  is a *context*, i.e. a sequence of variable declarations

$$(x_1 : A_1) \cdots (x_n : A_n)$$

where  $A_i : \text{Type}(x_1 : A_1) \cdots (x_{i-1} : A_{i-1})$  for every  $i = 1, \dots, n$ . The rule of *basic type formation* (figure 2) tells how types are formed from a category by instantiating the context. If  $n = 0$ , the context is empty, and  $C$  is itself a type.

The **fun** judgement

$$\mathbf{fun} f : A$$

presupposes that  $A$  is a type. It generates an object  $f$  of type  $A$ , to which the rules of application and abstraction apply in accordance with the type  $A$ , as well as the  $\beta$  and  $\eta$  conversion rules. These rules are shown in figure 2. They are more or less the standard rules of logical frameworks with dependent types, such as (Nordström *et al.*, 1990).

Syntactic sugar in figure 3 is of two opposite kinds: variable elimination and type factorization. As usual in dependently typed languages, variables can be eliminated from contexts and function types whenever there are no dependencies on them. The resulting notation is similar to simply typed languages, such as Haskell. On the other hand, the use of variables allows type factorizations that are not possible in Haskell. For instance, the following abbreviation is useful if  $A$  is complex:

$$(\_, \_ : A) \rightarrow B \equiv A \rightarrow A \rightarrow A \rightarrow B.$$

Basic type formation.

$$\frac{\mathbf{cat} \ C \ (x_1 : A_1) \cdots (x_n : A_n) \ a_1 : A_1 \ \dots \ a_n : A_n \ (x_1 := a_1, \dots, x_{n-1} := a_{n-1})}{C \ a_1 \dots a_n : \mathit{Type}}$$

Basic object formation.

$$\frac{\mathbf{fun} \ f : A}{f : A}$$

Function type formation, application, and abstraction.

$$\frac{A : \mathit{Type} \quad \frac{(x : A) \quad B : \mathit{Type}}{(x : A) \rightarrow B : \mathit{Type}}}{(x : A) \rightarrow B : \mathit{Type}} \quad \frac{f : (x : A) \rightarrow B \quad a : A}{f \ a : B(x := a)} \quad \frac{(x : A) \quad b : B}{\lambda x \rightarrow b : (x : A) \rightarrow B}$$

$\beta$  and  $\eta$  conversion.

$$(\lambda x \rightarrow b) \ a = b(x := a) \quad \frac{c : (x : A) \rightarrow B}{c = \lambda x \rightarrow (c \ x)}$$

Definition expansion.

$$\frac{f \ a_1 \dots a_n = t \ \gamma_1 \dots \gamma_n}{\text{for the first } \mathbf{def} \ f \ p_1 \dots p_n = t \ \text{such that } p_1 < \gamma_1 > a_1, \dots, p_n < \gamma_n > a_n}$$

Fig. 2. Rules for types and objects in abstract syntax.

Variable elimination.

$$\begin{aligned} (\_ : A) \Gamma &\equiv (x : A) \Gamma && \text{if } \Gamma \text{ does not depend on } x \\ A \ \Gamma &\equiv (\_ : A) \Gamma \\ (\_ : A) \rightarrow B &\equiv (x : A) \rightarrow B && \text{if } B \text{ does not depend on } x \\ A \rightarrow B &\equiv (\_ : A) \rightarrow B \\ \lambda \_ \rightarrow b &\equiv \lambda x \rightarrow b && \text{if } b \text{ does not depend on } x \end{aligned}$$

Factorization.

$$\begin{aligned} \mathbf{fun} \ f, \dots, g : A &\equiv \mathbf{fun} \ f : A ; \dots ; g : A \\ (x, \dots, y : A) &\equiv (x : A) \cdots (y : A) \\ (x, \dots, y : A) \rightarrow B &\equiv (x : A) \rightarrow \cdots \rightarrow (y : A) \rightarrow B \\ \lambda x, \dots, y \rightarrow b &\equiv \lambda x \rightarrow \cdots \rightarrow \lambda y \rightarrow b \end{aligned}$$

Fig. 3. Syntactic sugar for abstract syntax.

### 6.2.2 Normal forms of abstract syntax types and objects

The type of any **fun** function  $f$  has the form

$$(x_1 : A_1) \rightarrow \cdots \rightarrow (x_n : A_n) \rightarrow A$$

where  $A$  is a basic type  $C \ t_1 \dots t_m$  where  $C$  is a category. With reference to this form, we say that  $A_1, \dots, A_n$  are the *argument types* of  $f$ , that  $A$  is its *value type*, and that

$C$  is its *value category*. The *full application* of  $f$  has the form

$$f a_1 \dots a_n$$

whose type is  $A(x_1 := a_1 \dots x_n := a_n)$ . A term of a function type is in  $\beta\eta$  normal form, if it is an abstraction

$$\lambda z_1 \rightarrow \dots \rightarrow \lambda z_n \rightarrow b$$

and  $b$  is an application of a constant or a variable or a metavariable, with all arguments in  $\beta\eta$  normal form. We can use the  $\beta$  and  $\eta$  conversion rules to bring any well-typed term into this form.<sup>18</sup>

### 6.2.3 Metavariables

There is an infinite supply of metavariables

$$?_0, ?_1, ?_2, \dots$$

which can be terms of any type. Metavariables are generated in parsing and in interactive editing, and they do not appear in GF grammars. Since they are generated directly in  $\eta$ -expanded form  $\lambda x_1, \dots, x_m \rightarrow ?_k$ , the metavariable  $?_k$  itself has always a basic type  $C a_1 \dots a_n$ .

### 6.2.4 Definitions

A judgement of the form

$$\mathbf{data} C = f_1 \mid \dots \mid f_n$$

presupposes that  $C$  is a category and that  $f_1, \dots, f_n$  are **fun** functions, such that the value category of each  $f_i$  is  $C$ . What the judgement says is that  $f_1, \dots, f_n$  are *constructors* of the category  $C$ . Like in ALF (Magnusson & Nordström, 1994), constructors can be added incrementally, by new **data** judgements.

A judgement of the form

$$\mathbf{def} f p_1 \dots p_m = d$$

presupposes that  $f$  is a **fun** function but not a constructor, and that  $d$  is an object of type determined by the types of  $f$  and  $p_1, \dots, p_m$ .<sup>19</sup> The arguments  $p_1, \dots, p_m$  are *patterns*, i.e. terms formed from variables, the wildcard  $\_$ , and constructors. Those *def* judgements that have one and the same  $f$  form the *implicit definition* of that  $f$ . They determine how applications of  $f$  are *computed* by using *pattern matching*. Matching is performed in the order in which the equations appear in the grammar, and the patterns may overlap; the pattern-matching rules are the same as the ones for concrete syntax in figure 6.

<sup>18</sup> Thus a function  $f$  alone is a term in normal form only if its type is a basic type: in the general case, the normal form of the term  $f$  is  $\lambda z_1 \rightarrow \dots \rightarrow \lambda z_n \rightarrow f z_1 \dots z_n$ .

<sup>19</sup> Cf. the definition of pattern contexts in figure 6, modified for dependent types.

Functions  $f$  that are neither constructors nor defined implicitly are *primitive notions*. The lexical rules of GF make no distinction between constructors, defined functions, and primitive notions.

### 6.3 Concrete syntax types and expressions

#### 6.3.1 Tokens and strings

The type *Str*, informally called “strings”, is actually a type of lists of *tokens*, which are objects of type *Tok*. Tokens in normal form are quoted strings (“foo”). The *agglutination*  $t + u$  of two tokens is also a token. Token lists are built from the empty list [] and from tokens by means of *concatenation* ++ (Table 4).

We treat *Tok* and *Str* as abstract types, which can be instantiated by any types that support the aforementioned methods. A simple model is one in which tokens are strings and token lists are lists of strings. In this model, + is string concatenation and ++ is list concatenation. Expressions for tokens can also be used as expressions for singleton token lists, and are thus overloaded.<sup>20</sup>

#### 6.3.2 Parameters and parameter types

A *parameter type* (*PType*)  $P$  is defined by a *parameter declaration*

$$\mathbf{param} P = C_1 \Gamma_1 \mid \dots \mid C_n \Gamma_n$$

where each  $\Gamma_i$  is a *parameter context*, i.e. a sequence  $P_1 \dots P_m$  of parameter types. The parameter declaration introduces the *parameter constructors*  $C_1, \dots, C_n$ , which can be used as functions from their parameter contexts to  $P$ .

The parameter declarations of a grammar may not be recursive, nor mutually recursive. As a consequence, every parameter type  $P$  is finite, and we can form the list of all *parameter values* of type  $P$ ,

$$V_P = [1_P, 2_P, \dots, n_P]$$

generated by a left-to-right enumeration of the parameter values of type  $P$ .<sup>21</sup>

#### 6.3.3 Record types and records

The rules for labelled records in figure 4 are completely standard. Record labels have local scopes, and their name space is distinct from identifiers. For a record  $r$  to be of the type  $R$ , it is enough that every label of  $R$  is given a value of appropriate type in  $r$ . The order of fields does not matter, nor do superfluous fields in  $r$ . In notation, fields in record types and records can be factorized.

<sup>20</sup> Another model of tokens are word descriptions obtained from a morphological analyser, e.g. “point+Noun+Pl” instead of “points”. This demands separate morphology passes before parsing and after linearization, but is an efficient way to implement parsing when the lexicon is large.

<sup>21</sup> This is similar to the derivation of Enum class instances in Haskell, but more powerful, since it applies not only to enumerated types but also to disjunctive and conjunctive types.

The types of tokens and strings and parameter types.

$$\text{Tok, Str, PType} : \text{Type}$$

Tokens, the empty string, tokens as strings, agglutination, and concatenation.

$$\text{"foo"} : \text{Tok} \quad [] : \text{Str} \quad \frac{t : \text{Tok}}{t : \text{Str}} \quad \frac{t, u : \text{Tok}}{t + u : \text{Tok}} \quad \frac{s, t : \text{Str}}{s ++ t : \text{Str}}$$

Parameter types and constructors.

$$\frac{\text{param } P = \dots}{P : \text{PType}} \quad \frac{P : \text{PType}}{P : \text{Type}} \quad \frac{\text{param } P = \dots \mid C P_1 \dots P_n \mid \dots}{C : P_1 \rightarrow \dots \rightarrow P_n \rightarrow P}$$

Record type formation.

$$\frac{T_1, \dots, T_n : \text{Type}}{\{r_1 : T_1 ; \dots ; r_n : T_n\} : \text{Type}}$$

Record formation and projection.

$$\frac{t_1 : T_1 \quad \dots \quad t_n : T_n}{\{r_1 = t_1 ; \dots ; r_n = t_n\} : \{r_1 : T_1 ; \dots ; r_n : T_n\}} \quad \frac{c : \{\dots ; r : T ; \dots\}}{c.r : T}$$

Projection computation.

$$\{\dots ; r = t ; \dots\}.r = t$$

Table type formation.

$$\frac{P : \text{PType} \quad T : \text{Type}}{P \Rightarrow T : \text{Type}}$$

Table formation and selection.

$$\frac{t_1 : T \quad \Gamma_P p_1 \quad \dots \quad t_n : T \quad \Gamma_P p_n \quad [p_1, \dots, p_n \text{ exhaustive for } P]}{\mathbf{table} \{p_1 \Rightarrow t_1 ; \dots ; p_n \Rightarrow t_n\} : P \Rightarrow T} \quad \frac{c : P \Rightarrow T \quad p : P}{c ! p : T}$$

Selection computation.

$$\mathbf{table} \{\dots ; p \Rightarrow t ; \dots\} ! v = t\gamma \quad \text{for the first } p \text{ such that } p < \gamma > v$$

Local definition.

$$\frac{(x : T) \quad t : T \quad e : E}{\mathbf{let} \{x : T = t\} \mathbf{in} \quad e : E} \quad \mathbf{let} \{x : t = T\} \mathbf{in} \quad e = e(x := t)$$

Global definition.

$$\frac{\mathbf{oper} \quad h : T = t}{h : T} \quad \frac{\mathbf{oper} \quad h : T = t}{h = t}$$

Fig. 4. Types and objects of concrete syntax.

Concatenation of tokens.

$$["foo \dots bar"] \equiv "foo" ++ \dots ++ "bar"$$

Factorization.

$$\begin{aligned} \{\dots; r, \dots, s : T; \dots\} &\equiv \{\dots; r : T; \dots; s : T; \dots\} \\ \{\dots; r, \dots, s = t; \dots\} &\equiv \{\dots; r = t; \dots; s = t; \dots\} \\ \mathbf{let} \{x_1 : T_1 = t_1; \dots; x_n : T_n = t_n\} \mathbf{in} e &\equiv \mathbf{let} \{x_1 : T_1 = t_1\} \mathbf{in} \\ &\dots \mathbf{let} \{x_n : T_n = t_n\} \mathbf{in} e \end{aligned}$$

Case expression.

$$\mathbf{case} e \mathbf{of} \{\dots\} \equiv \mathbf{table} \{\dots\} ! e$$

Fig. 5. Syntactic sugar for concrete syntax.

### 6.3.4 Table types, tables and pattern matching

Tables of type  $P \Rightarrow T$  are finite functions from  $P$  to  $T$ . The argument type  $P$  must be a parameter type. The normal form of a table is a complete enumeration of argument-value pairs:

$$\mathbf{table} \{1_P \Rightarrow t_1 ; \dots ; n_P \Rightarrow t_n\}$$

For convenience, and to capture generalizations, the use of *patterns* is permitted as well.<sup>22</sup> There are three kinds of patterns, as shown in figure 6. The matching relation

$$p < \gamma > v$$

“the pattern  $p$  matches the value  $v$  with the substitution  $\gamma$ ”, is used in figure 6 to define the computation of *selections* from tables with patterns. To test whether a list of patterns is *exhaustive* for a given type  $P$ , we just test whether all values of type  $P$  are matched by them.

Patterns are matched from left to right, and they are allowed to overlap, like in Haskell. Nonlinear patterns are forbidden (i.e. patterns where a variable  $x$  occurs more than once). It is in virtue of this that pattern contexts and substitutions can be simply concatenated.

If the argument type  $P$  of a table  $t$  is known, the table can be expanded to eliminate patterns, by going through the list of all parameter values of  $P$ :

$$t = \mathbf{table} \{1_P \Rightarrow t ! 1_P ; \dots ; n_P \Rightarrow t ! n_P\}$$

It is often handy to use **case** expressions as syntactic sugar for selections, as shown in figure 5.

<sup>22</sup> Tables with patterns are syntactically similar to fn expressions in ML.

Patterns: wildcard, variable, and constructor.

$$- \quad x \quad C p_1 \dots p_n$$

Pattern contexts.

$$\Gamma_{P-} = () \quad \Gamma_{P x} = (x : P)$$

$$\Gamma_P(C p_1 \dots p_n) = \Gamma_{A_1 p_1} \dots \Gamma_{A_n p_n} \text{ if } C : A_1 \rightarrow \dots \rightarrow A_n \rightarrow A$$

Pattern matching rules.

$$- \triangleleft v \quad x <(x := v)> v \quad \frac{p_1 <\gamma_1> v_1 \quad \dots \quad p_n <\gamma_n> v_n}{C p_1 \dots p_n <\gamma_1 \dots \gamma_n> C v_1 \dots v_n}$$

Fig. 6. Patterns and pattern matching.

### 6.3.5 Functions, operation definitions, and local definitions

Concrete syntax uses the same rules for functions and function types as abstract syntax (figure 2). Since functions are not part of canonical GF (section 5.4), the  $\beta\eta$  normal form plays no role in concrete syntax.

Functions are mostly introduced in *operation definitions*, judgements of the form

$$\mathbf{oper} \ h : T = t$$

where  $T$  is any type (in the sense of concrete syntax) and  $t : T$ . *Local definitions* (**let** expressions) have a similar syntax, but they are local to expressions.

The **oper** definitions of a grammar may not be (mutually) recursive. A consequence of this is that a defined operation  $h$  can always be eliminated from a grammar by replacing it with its definition  $t$ ; this procedure is known as *inlining*.

### 6.3.6 Type variables

GF does not have polymorphism. Explicit type variables and dependent types are used instead. The type of these variables is *Type*; there is for the time being no stratification. An example is the *flip* function (familiar from Haskell):

$$\mathbf{oper} \ flip : (a, b, c : Type) \rightarrow (a \rightarrow b \rightarrow c) \rightarrow b \rightarrow a \rightarrow c = \lambda_{\rightarrow, \rightarrow, \rightarrow} f, x, y \rightarrow fyx$$

Another example is a local type definition, which is not possible in Haskell:

$$\mathbf{let} \ \{S : Type = \{s : Str\}\} \ \mathbf{in} \ S \rightarrow S \rightarrow S \rightarrow S$$

## 6.4 Concrete syntax for abstract syntax

A mathematical view of a GF grammar is that the abstract syntax defines a free algebra of syntax trees, and the concrete syntax defines a homomorphism from this algebra into a system of concrete-syntax objects. Mainly for the purpose of deriving parsers, we restrict concrete-syntax objects to certain special forms of records, captured by the notion of a linearization type.

### 6.4.1 Linearization types

A *linearization type*  $L$  is a record type usable as value type of linearization. It must have one or more fields whose types are *Str*-valued tables, the other fields having parameter types. More precisely,

- A string type is either *Str* or  $P \Rightarrow S$  where  $S$  is a string type.
- If  $S$  is a string type,  $\{s : S\}$  is a linearization type.
- If  $L$  is a linearization type, so is the type resulting from adding a field  $r : T$  where  $T$  is a parameter type or a string type.

To simplify the generation of a parser, we require that all and only the string type valued fields are labelled  $s_1, s_2, \dots$ , or  $s$ .<sup>23</sup>

When giving a linearization rule to a function whose arguments have function types, we need to know what the linearization type of such a type is. The following clauses define this notion inductively for all types:

$$\begin{aligned} (C a_1 \dots a_n)^\circ &= L, \text{ if } \mathbf{linat} C = L \\ ((x_1 : A_1) \rightarrow (x_n : A_n) \rightarrow A)^\circ &= \{v_1 : \mathit{Str}; \dots; v_n : \mathit{Str}\} ** A^\circ \end{aligned}$$

In the latter clause, we assume that the function type is in normal form, i.e. that  $A$  is a basic type. To form the linearization type of a function type, we thus add to the linearization type of the value type one field of type *Str* for each argument type. The idea is to introduce a field for each variable symbol. If  $n = 1$ , we follow the convention of adding  $v : \mathit{Str}$  without a subscript. To avoid clashes with the system-generated labels  $v, v_1, v_2, \dots$ , for bound variables, these labels are forbidden in user-defined record types.

### 6.4.2 Linearization type definitions and linearization rules

A linearization type definition for a category  $C$  is a judgement of the form

$$\mathbf{linat} C = L$$

which presupposes that  $L$  is a linearization type. A notational convention allows us to omit the **linat** judgement of a category  $C$  if the linearization type is  $\{s : \mathit{Str}\}$ .

A linearization rule for a function

$$\mathbf{fun} f : (x_1 : A_1) \rightarrow (x_n : A_n) \rightarrow A$$

is a judgement of the form

$$\mathbf{lin} f = t$$

which presupposes that

$$t : A_1^\circ \rightarrow \dots \rightarrow A_n^\circ \rightarrow A^\circ$$

A concrete syntax is *complete* w.r.t. an abstract syntax, if it contains a **linat** judgement for every **cat** judgement, and a **lin** judgement for every **fun** judgement.

<sup>23</sup> An alternative would be to type-annotate the record fields.

Pattern variables in linearization rules.

$$\mathbf{lin} f x_1 \dots x_m = t \equiv \mathbf{lin} f = \lambda x_1, \dots, x_m \rightarrow t$$

The **pattern** rule format.

$$\mathbf{pattern} f x_1 \dots x_n = t_1 ++ \dots ++ t_m \equiv \mathbf{lin} f x_1 \dots x_n = \{s = t'_1 ++ \dots ++ t'_m\}$$

where "foo'" = "foo" and  $x'_i = x_i.s$ .

Fig. 7. Syntactic sugar for linearization.

The **pattern** notation for linearization rules (Section 2) can be used if the argument and value types of  $f$  all have the linearization type  $\{s : Str\}$ . A pattern element  $t_i$  is then either a token or one of the variables  $x_j$  (see figure 7).

### 6.4.3 Default linearization

To linearize symbols not defined in the grammar (variables and metavariables), GF uses *default linearization*. It is a function that takes a string to an object of a linearization type. The default linearization of a category  $C$  is defined by the judgement

$$\mathbf{lindef} C = t$$

which presupposes

$$t : Str \rightarrow C^\circ.$$

If the grammar does not contain a judgement of this form, a default default linearization is used: for a string  $t$ , it is a record where every field of a string type has uniformly the value  $t$ , and every parameter field of type  $P$  has as its value  $1_P$ , i.e. the first value of type  $P$ .

## 7 Partial evaluation and parsing

While the linearization and typechecking algorithms follow straightforwardly from the semantics of GF, there are two other algorithms that are crucial for most practical applications of GF, and which are nontrivial: partial evaluation and parsing. Partial evaluation takes a GF grammar into a form by which linearization can be performed with the minimum of interpretational overhead. The same form is the basis of parsing, since it permits the derivation of a finite set of context-free rules.

### 7.1 Partial evaluation

Partial evaluation is evaluation at compile time, leaving evaluation at runtime less work to do (Jones *et al.*, 1993). A partial evaluator takes a program and computes it as far as it can. In general, the result is not fully evaluated, since some input of the program is unknown. Moreover, partial evaluation can be performed in different

ways, e.g. optimizing either the time needed to run the runtime program or the space needed to store it.

What we do in GF is evaluate linearization rules into the canonical form of section 5.4. This operation does not always optimize space: although it sometimes does reduce space, it may also do quite the contrary. However, the result is always a time-optimized runtime grammar.

Given a rule pair

$$\mathbf{fun} f : (x_1 : A_1) \rightarrow \dots \rightarrow (x_n : A_n) \rightarrow A ; \mathbf{lin} f = t$$

we  $\eta$ -expand the term  $t$  with respect to a sequence of *argument variables*, that is, variables standing for arguments of the function  $f$ . We denote argument variables by pairs  $\langle C, i \rangle$ , where  $C$  is the value category  $C$  of the argument type  $A_i$ , and  $i$  tells that it is the  $i$ 'th argument. Showing the category will be useful in parser generation. The result is a linearization rule

$$\mathbf{lin} f = \lambda \langle C_1, 1 \rangle, \dots, \langle C_n, n \rangle \rightarrow \{r_1 = t_1 ; \dots ; r_m = t_m\}$$

where further  $\eta$ -expansions w.r.t. the linearization type of  $A$  have produced the record form. Each  $t_j$  ( $j = 1, \dots, m$ ) is either an expression of a basic type (*Str* or a parameter type) or a fully expanded table

$$\mathbf{table} \{1_p \Rightarrow u ! 1_p ; \dots ; n_p \Rightarrow u ! n_p\}$$

where each term  $u ! k_p$  is evaluated further by inlining **oper** and **let** constants,  $\eta$ -expanding records and tables, and applying the evaluation rules for projections, selections, and function applications. Sometimes we also need transformations analogous to the elimination of maximal segments in proof theory (Prawitz, 1965). The most important such transformation is

$$(\mathbf{table} \{p \Rightarrow f ; \dots ; q \Rightarrow g\} ! e) a \triangleright \mathbf{table} \{p \Rightarrow f a ; \dots ; q \Rightarrow g a\} ! e$$

pushing an application inside a table. The transformation is needed if the selection cannot be computed: such is the case if  $e$  depends on an argument variable. There are similar rules for projection and selection. The eliminability of function applications is analogous to the subformula property of intuitionistic propositional calculus: since the type of the linearization term consists solely of records, tables, strings, and parameters, no terms of function types need appear in it.

After partial evaluation, the only remaining unknown input in linearization rules are the argument variables. Because of compositionality, they can be treated as pointers to the linearizations of subtrees.

## 7.2 Parsing

In sections 2.5 and 3.3, the parsing problem of simplified versions of GF were reduced to parsing in context-free grammars. We will now do the same to full GF.

We use the partially evaluated form of linearization functions, and consider the fields that are of string types. Each field is an  $n$ -place table ( $n \geq 0$ ) that has finitely

many possible values  $u$  of type *Str*. We call these values  $u$  the *productions* of the function  $f$ . We derive a finite set of context-free rules from every production: starting from a linearization rule

$$\mathbf{lin} f = \lambda \langle C_1, 1 \rangle, \dots, \langle C_n, n \rangle \rightarrow r$$

consider an  $s_k$ -labelled (and thus *Str*-valued) field of  $r$ . Let  $u$  be a production coming from this field. The context-free rules generated from  $u$  have the form

$$f_p \cdot C^k ::= c$$

where the value category  $C^k$  stands for the  $k$ 'th part of  $C$ , and the right-hand side  $c$  is a sequence of context-free items. The sequence  $c$  is constructed from elements of the set  $u^*$ , which is defined inductively on the structure of the production  $u$ :

Terminal:	$s^*$	=	$\{s\}$
Nonterminal:	$(\langle C, i \rangle .s_j)^*$	=	$\{C_i^j\}$
Binding:	$(\langle C, i \rangle .v_j)^*$	=	$\{Var_i^j\}$
Concatenation:	$(a ++ b)^*$	=	$\{c' \mid c \in a^*, c' \in b^*\}$
Selection:	$(a ! p)^*$	=	$a^*$
Table:	$(\mathbf{table} \{p_1 \Rightarrow b_1; \dots; p_k \Rightarrow b_k\})^*$	=	$\bigcup_{i=1}^k b_i^*$

Now, each sequence in the set  $u^*$  has the form

$$c_1 \dots c_m$$

where each  $c_l$  is either a terminal  $s$  or a category symbol  $C_i^j$  indexed by the argument position  $i$  and the discontinuous-part number  $j$ . This is not yet a sequence of context-free items, because of the presence of the position numbers: these numbers are needed for profiles. However, the formation of context-free items

$$c'_1 \dots c'_m$$

is just simplification of nonterminals: each  $Var_i^j$  becomes  $Var$ , and each  $C_i^j$  becomes  $C^j$ . The profile  $p$  is constructed by collecting, from the subsequence of nonterminals in  $c_1 \dots c_m$ , the list of positions for each argument place of  $f$ , in the same way as in section 3.3.

Profiles  $p$  in rule labels  $f_p$  were earlier just an optimization removing the need to look up  $f$  in the grammar when postprocessing parse trees. When parametric variation is introduced, profiles become indispensable: the arguments of  $f$  may be placed to different positions when the tree is linearized under different parameters. Thus it would not be enough to look up  $f$  to restore the order of constituents.

When restoring syntax trees from parse trees, we do the same as in sections 2.5 and 3.3. Discontinuous constituents bring nothing new to this, since they can be treated as reduplications. What typically happens in these reduplication instances is that different parts have metavariables in different places, so that no conflicts arise when the final result is formed by unifying the parts. For example, the German grammar of section 5.3 generates the context-free grammar (with profiles that ignore

bindings, which are empty)

$$\begin{aligned} \text{Pred}_{[[1],[2,3]]}. S &::= NP VP^1 VP^2 \\ \text{Pred}_{[[2],[1,3]]}. S &::= VP^1 NP VP^2 \\ \text{Pred}_{[[1],[2,3]]}. S &::= NP VP^2 VP^3 \\ \text{Compl}_{[[1],[0]]}. VP^1 &::= TV \\ \text{Compl}_{[[0],[1]]}. VP^2 &::= NP \end{aligned}$$

The sentence *Johann liebt Maria* has the initial parse tree

$$\text{Pred}_{[[1],[2,3]]} \text{ Johann } (\text{Compl liebten ?}) (\text{Compl ? Maria})$$

which unifies to the final tree

$$\text{Pred Johann } (\text{Compl lieben Maria})$$

Since the formation of context-free rules suppresses all parameters, the parser is over-tolerant. It could, for instance, recognize *they walks* as a valid English sentence. A strict parser is obtained by filtering away all those parse trees whose linearization does not match the input string. This arrangement of parsing has the disadvantage of being potentially inefficient: the number of rejectable parses can be exponential<sup>24</sup>. Its advantages are the simplicity of implementation and that it gives, as by-product, *grammar correction*: we can use tolerant parsing followed by linearization to correct *they walks* into *they walk*. Parsing via context-free grammars is known as *off-line parsing* in the context of unification grammars: the alternative is to perform unification at each construction step of the syntax tree.

### 7.3 The expressive power of GF

The expressive power of a grammar formalism is often characterized by its weak generative capacity – the class of sets of strings (in Chomsky hierarchy) it is capable of generating. Even though the focus in GF is on strong generative capacity (the trees it assigns to strings), its place in Chomsky hierarchy is a meaningful question.

In a trivial sense, GF is in the class 0 of unlimited languages, since we can define the universal language  $U$  of strings over any finite alphabet and encode any predicate  $P$  on  $U$  in the abstract syntax as a type of proofs. The rule pair

$$\mathbf{fun} \ f : (x : U) \rightarrow P \ x \rightarrow S \ ; \ \mathbf{lin} \ f \ x \ y = x$$

defines a string language  $S$  which is undecidable if  $P$  is. This construction is based on the suppression of an argument in linearization. However, since parsing of suppressed arguments is solved by introducing metavariables, it remains a meaningful question what happens if we consider GF without suppression. We have already seen that GF is more powerful than the class 2 of context-free languages (section 2). GF is not just *mildly context-sensitive* like e.g. TAG (Joshi, 1985), since GF can define the

<sup>24</sup> Or even infinite, if cyclic rules are present.

double copy language  $\{wewew \mid w \in \{a,b\}^*\}$ . The precise location of GF-without-suppression is an open question.

## 8 The implementation of GF

Above we have described GF as a language of its own, independently of implementation. This description is partly an abstraction from earlier implementation work, partly a specification followed in later work. In this section, we will give an outline of the implementation and some problems that we have encountered in it.

### 8.1 Overview of the code

The Haskell implementation of GF (Version 1.0) has 12k lines of source code in 95 modules. The main parts of the code are the following:

- Grammar compiler: lexer, parser, type checker, partial evaluator, parser generator.
- Command line interpreter: functions to read grammar files and use grammars in batch mode.
- Syntax editor: functions to edit GF objects interactively.

The syntax editor is based on an abstract command language built upon a zipper, and it can be used through different user interfaces: we have a line-based editor, a graphical editor written in Fudgets (Carlsson & Hallgren, 1998), an experimental speech-based editor (Ranta & Cooper, 2001), and a Java GUI client communicating with a GF server via an XML-based protocol.

### 8.2 The use of Haskell

Haskell was chosen as implementation language for two reasons: we found it to be a good general-purpose programming language (particularly good for implementing compilers for functional languages), and we wanted to connect smoothly with some other programs written in Haskell, in particular, the proof editor Alfa (Hallgren, 2000). Some of the code was translated from earlier SML programs; in general, we did not want to exploit the laziness of Haskell in any essential way.<sup>25</sup> Neither did we use Haskell's impure features such as strictness flags. Monads (IO, error, state) are used heavily, and some classes are defined to simplify function names. GF conforms to the Haskell 98 standard (Peyton Jones & Hughes, 1999), and can be compiled with all standard compilers and interpreters, on all major operating systems.<sup>26</sup>

<sup>25</sup> There is one single point where laziness would be useful: to treat infinite lists of parse trees arising in cyclic grammars (Section 2.5).

<sup>26</sup> Since the Fudgets library (Carlsson & Hallgren, 1998) requires the X window system, the Java GUI is the only graphical interface that works on Microsoft Windows.

### 8.3 Performance

There are two demanding components in GF: grammar compilation and object-language parsing. The parser of grammars was created using the Happy parser generator (Marlow, 2001), and it performs well. If the grammar is close to canonical form, type checking and partial evaluation together take less time than parsing.<sup>27</sup> However, if the grammar makes heavy use of functions and pattern matching, partial evaluation may take ten times longer than parsing. Some heuristics have helped considerably, such as topologically sorting all **oper** definitions and compiling them in dependency order, ignoring unused operations. Of course, once a grammar is ready, the compiled version can be saved in a file for rapid reuse.

The inefficiency of object-language parsing is partly due to the inherent complexity of general-purpose context-free parsing algorithms.<sup>28</sup> This can be helped in the special case of LALR(1) by using Happy parsers, which can be automatically generated from GF. More often, however, the bottleneck is postprocessing. To solve this problem, postprocessing would have to be integrated in the first parsing phase, using e.g. the semantic actions of Happy or some form of attribute or unification grammar. This is a research problem rather than an implementation issue.

### 8.4 Accessing GF

Users who do not write grammars themselves typically use GF via the graphical interactive editor. For grammar developers, and writers of batch programs, there is a command language and a shell, also permitting scripts. For instance, the following script imports an English and a French grammar, reads the file `enter.txt`, parses it as an English text, and linearizes the resulting tree in French:

```
i alarm.Eng.gf
i alarm.Fra.gf
rf enter.txt | p -lang=Eng | l -lang=Fra
```

Haskell programmers can access GF through an API (Application Programmer's Interface) module. It contains both default and customizable versions of parsing, linearization, and translation functions. This makes it possible to include GF functionalities and use GF grammars in other Haskell programs.

A library of macros is provided for creating GF grammars by Haskell programs. One way of using these macros is to define translations from other grammar formats to GF. For instance, BNF and EBNF can be used as input formats. Another use of code generation is to bypass the partial evaluator of GF: make all generalizations and abstractions in the Haskell code, and generate canonical GF directly. The next step from this idea would be to define GF as an embedded language (Hudak, 1996). However, we prefer to see GF as a language of its own, which can be used and

<sup>27</sup> Parsing a 22k-line grammar with a Swedish resource lexicon takes 4 seconds on a 1.5 GHz Pentium 4 with RedHat Linux 7.1; the rest of the compilation of this close-to-canonical grammar takes 3 seconds.

<sup>28</sup> The time they take is cubic in the length of the input string.

reasoned about independently of implementation language. Moreover, since GF has dependent types, it is not possible to rely on code generated from Haskell: at least a type checker would in any case have to be written.

## 9 Some applications of GF

GF grammars have been written for fragments of at least 20 natural languages and many formal languages. Most of these grammars serve the theoretical purpose of verifying that GF can express a particularly intricate grammatical rule, or formalize the semantics of some specific application. The following list mentions some applications that have passed the level of first experiments and become independent projects.

**Proof text editors.** These are systems in which formal proofs are interactively constructed in type theory and at the same time viewed as texts in natural language. Via a parser, natural language input is also possible. The system is extensible to user-defined constants by means of user-defined linearization rules; if a rule is not given, a default linearization is generated. Users can also extend it to new natural languages by writing GF grammars for the concrete syntax. Two implementations of proof text editors have been made in GF: one that works as a plug-in in the proof editor Alfa (Hallgren & Ranta, 2000), with support for English, French and Swedish, and another one using the generic GF interface, also supporting Finnish, Italian, and Russian.

**Software specifications.** Formal and semi-formal software specification languages, such as OCL (Warmer & Kleppe, 1999), are widely used in industry, but still wider is the use of informal specifications in natural language. A project is going on to bridge this gap by building an abstract specification language in GF, with concrete syntaxes for OCL and English (Hähnle *et al.*, 2002). The goal is to enable simultaneous production of formal and informal specifications. The editor is being integrated in an industrial CASE tool.

**Controlled language.** This is the next step from mathematical proofs via software specifications towards non-mathematical language. Controlled languages are subsets of natural languages used for technical purposes such as instruction manuals for aircraft maintenance. Today's controlled languages (e.g The Boeing Company, 2001) have neither formal grammars nor automatic checkers. But GF has been used to define prototypes where formal verification is applied to documents written in natural language. An example is a set of instructions for using an alarm system, generated in English, French, German, and Swedish, and equipped with a formal proof that the instructions preserve the system in a legal state (Johannisson & Ranta, 2001).

**Dialogue systems.** This is human-machine interaction where information is gathered by questions and answers. For instance, in a travel-agency dialogue system the machine asks where and when the customer wants to travel. The human answers all questions till enough information has been gathered to complete the booking. In order for the dialogue not to be too monotonous, the dialogue system should

be flexible and, for example, accept answers to many questions at once. Several such criteria are identified in Bohlin *et al.* (1999). Somewhat surprisingly, it turned out that the metavariable-based model of interaction in proof editors readily fulfils most of these criteria, even adding extra functionality, e.g. a better control of the continuation of a dialogue via dependent types (Ranta & Cooper, 2001).

## 10 Related work

### 10.1 Montague grammar and categorial grammars

From the linguistic point of view, GF belongs to the tradition of *Montague grammar* (Montague, 1974). For Montague, a grammar was a set of rules linearizing logically interpreted analysis trees into strings of a natural language. The focus was on semantics rather than concrete syntax. A well-known problem in Montague's syntax is the use of so-called "quantifying in" rules to linearize variable-binding operations. Unlike other parts of Montague grammars, these rules cannot be directly formalized in GF, since they are not compositional. The rules can be circumscribed, however, partly by using combinators instead of variable binding (as suggested by Steedman (1988)), partly by means of discontinuous constituents.

The distinction between abstract and concrete syntax is seldom made by linguists. It was suggested, however, by the logician Haskell B. Curry, under the headings of *tectogrammatic* and *phenogrammatic* structure (Curry, 1963). For Curry, a tectogrammatic structure is similar to a term in combinatory logic, and it can show up as different phenogrammatic structures in different languages. Neither Curry nor Montague pursued the multilingual aspect, but there is a machine translation project, Rosetta (1994), based on Montague grammar.

*Categorial grammar* shares with Montague grammar the use of a type system to explain syntactic well-formedness. However, the idea is to explain not only abstract but also concrete syntax in terms of function application. To this end, Bar-Hillel (1953) made a distinction between prefix and postfix function types,  $\beta/\alpha$  vs.  $\alpha/\beta$ . His idea was developed further by Lambek (1958), resulting in a calculus that covered an impressive fragment of English, and was eventually proved equivalent to context-free grammars (Pentus, 1993). Extensions of Lambek calculus use richer sets of connectives (Morrill, 1994), or treat it as noncommutative linear logic (Abrusci, 1990).

In functional programming, some efforts have been made to implement logical and categorial grammars. A parser for a Montague-style grammar was implemented as a part of a database query system by Frost and Launchbury (1989), in Lazy ML. The grammar used in the system can also be defined in GF. A parser for the categorial grammar of Shaumyan was implemented by Jones and Hudak (1995) in Haskell. This theory shares with Lambek calculus the use of typing rules to define concrete syntax.

### 10.2 Unification grammars

*Unification grammars* (Shieber, 1986) are a family of grammar formalisms where context-free categories are made dependent on *features*, which the parser tries to

unify. Many grammar formalisms in computational linguistics belong to this family. Definite Clause Grammar (DCG) (Pereira & Warren, 1980) is perhaps the purest and simplest of them, and it has a built-in implementation in the Prolog programming language. It is also the most widely known, because it works well in education. The biggest grammars, however, have been written in Head Driven Phrase Structure Grammar (HPSG) (Pollard & Sag, 1994).

A typical example of DCG is the English predication rule,

$$S \longrightarrow NP(n) VP(n)$$

This rule expresses the condition that the subject and the verb must have the same number,  $n$ . In GF, the natural way to express the predication rule would be

$$\mathbf{fun} \text{ Pred} : NP \rightarrow VP \rightarrow S ; \mathbf{lin} \text{ Pred} N V = \{s = N.s ++ V.s ! N.n\}$$

The traditional grammar view is closer to GF than to DCG: the subject and the verb are not in symmetric relation, but the verb depends on the subject. The subject *has* a number (as inherent feature), which it *gives* to the verb (as parameter).

The advantage of treating inherent features and parameters on a par is computational: it allows a direct implementation of parsing as unification. From the descriptive point of view, DCG appears as a low-level language, which moreover does not have types. A suggestive way of parsing in GF grammars would be to compile them into a DCG, and use local unification instead of off-line parsing and postprocessing.

HPSG inherits from PATR (Shieber, 1986) the use of records to express complex grammatical objects. In HPSG, these records contain both syntactic and semantic information. For instance, the English noun form *integers* could be described by the record (in GF notation)

$$\{cat = CN ; sem = Int ; phon = "integers" ; n = Pl ; g = Neut\}$$

Records like this are called *signs* in HPSG. The information contained in a sign belongs partly to function declarations and partly to linearization rules in GF. We come close to a sign if we take a linearization record and add fields for the type and the syntax tree. However, the result is not quite the same: for *Int*, we get

$$\{cat = CN ; sem = Int ; s = \mathbf{table} \{Sg \Rightarrow "integer"; Pl \Rightarrow "integers"\} ; g = Neut\}$$

The difference reflects the characteristic fact that HPSG records are obtained by analysing strings, whereas GF records are obtained by linearizing trees. The HPSG record is, in a sense, an *instance* of the GF record: it shows one branch of a table instead of the whole table.

Another interesting feature of HPSG is that it has a type system, which helps to detect errors at compile time. As regards records, the type system has much in common with GF. But there are no function types and thus no higher-order abstractions available for grammar writers.

### 10.3 Syntax editors

As a syntax editor, GF belongs to the tradition starting from Mentor (Donzeau-Gouge *et al.*, 1975) and the Cornell Program Synthesizer (Teitelbaum & Reps, 1981). These systems were initially not frameworks but had a hard-wired object language. Later on, the Cornell system used attribute grammars (Knuth, 1968) in the same role as the GF formalism is used in the GF editor, and Mentor was extended by the formalism framework Metal (Kahn *et al.*, 1983). As for concrete syntax, these systems of course only support unambiguous programming languages, one at a time. On the abstract level, they have advanced computational features, such as stepwise forward and backward execution of code.

Proof editors are a descendant of syntax editors, and the closest to GF are those that use dependent types. GF has inherited its type theory from ALF (Magnusson & Nordström, 1994), which uses metavariables, whereas NuPRL (Constable, 1986), Coq (The Coq Development Team, 1999), and LEGO (Luo & Pollack, 1992) use tactics. All these systems support some amount of user-defined syntactic sugar, such as infix declarations, but are of course far from natural language syntax. For Coq, a natural-language interface exists (Coscoy *et al.*, 1995). It works in the direction of linearization only and cannot be extended by the user; however, it has some built-in optimizations that are not possible in compositional linearization. Even more in this direction is the proof explanation system PRex (Fiedler, 2001), which uses AI methods to adapt proof texts for individual users.

WYSIWYM (“What you see is what you mean”) is a multilingual authoring system for software manuals (Power & Scott, 1998). The user edits an abstract object which is reflected by “feedback texts” in English, French, and Italian. The grammars are hard-wired in the system and work in the direction of linearization only; the research emphasis is clearly on interaction rather than on grammars.

## 11 Conclusion

We have defined a grammar formalism GF on top of a logical framework with dependent types. The formalism is a special-purpose functional programming language, which adds the known advantages of functional languages (type checking, high abstraction level, succinctness of expression) to a simple computational model. GF grammars can be used for both parsing and generation of languages. The formalism is able to describe semantic conditions and intricate natural-language structures. It differs from earlier grammar formalisms by being based on functional programming and by having a powerful type system. The most important remaining problem is the inefficiency of the parsers generated from some GF grammars.

The main applications of GF are in domain-specific fragments of natural language, which have a semantic model that can be described in type theory. GF grammars provide natural-language interfaces to such models and make it possible to translate domain-specific language reliably via the model. GF supports interaction: it can be used as a multilingual authoring system in which texts are created in many

languages simultaneously. For future developments, an important task is to develop libraries of domain-independent resource grammars.

GF has been implemented in the functional language Haskell. The implementation follows the Haskell 98 standard and is portable to different operating systems. In addition to the separate program, GF functionalities can be accessed from other Haskell programs through an API module.

### Acknowledgements

Thanks to Markus Forsberg, Reiner Hähnle, Kristofer Johannisson, Bengt Nordström, and an anonymous referee for careful reading and valuable comments on this paper. GF itself has benefitted from contributions by many people at XRCE Grenoble, at Chalmers and Gothenburg University, at CAMS in Paris, and in the European TYPES Community. The work has been supported by the VINNOVA foundation within the project Interactive Language Technology (2001-06340).

### References

- Abrusci, M. (1990) Noncommutative Intuitionistic Linear Propositional Logic. *Zeitschrift für Mathematische Logik und Grundlagen der Mathematik*, **36**, 297–398.
- Augustsson, L. (1998) Cayenne—a language with dependent types. *Proc. ICFP'98*. ACM Press.
- Bar-Hillel, Y. (1953) A quasi-arithmetical notation for syntactic description. *Language*, **29**, 27–58.
- Bescherelle. (1997) *La conjugaison pour tous*. Hatier.
- Bohlin, P., Bos, J., Larsson, S., Lewin, I., Matheson, C. and Milward, D. (1999) *Survey of existing interactive systems*. Trindi deliverable D1.3, Gothenburg University.
- Bresnan, J. (ed). (1982) *The Mental Representation of Grammatical Relations*. MIT Press.
- Carlsson, M. and Hallgren, T. (1998) *Fudgets—Purely Functional Processes with applications to Graphical User Interfaces*. Ph.D. thesis, Department of Computing Science, Chalmers University of Technology.
- Constable, R. L. (1986) *Implementing Mathematics with the NuPRL Proof Development System*. Prentice-Hall.
- Coquand, T. (1996) An algorithm for type checking dependent types. *Science of Computer Programming*, **26**, 167–177.
- Coscoy, Y., Kahn, G. and Théry, L. (1995) Extracting text from proofs. In: Dezani-Ciancaglini, M. & Plotkin, G. (eds.), *Proc. Second Int. Conf. on Typed Lambda Calculi and Applications: LNCS 902*, pp. 109–123. Springer-Verlag.
- Curry, H. B. (1963) Some logical aspects of grammatical structure. In: Jakobson, R. (ed.), *Structure of Language and its Mathematical Aspects: Proceedings of the Twelfth Symposium in Applied Mathematics*, pp. 56–68. American Mathematical Society.
- de Bruijn, N. G. (1994) Mathematical Vernacular: a Language for Mathematics with Typed Sets. In: Nederpelt, R. (ed.), *Selected Papers on Automath*, pp. 56–68. North-Holland.
- Donzeau-Gouge, V., Huet, G., Kahn, G., Lang, B. and Lévy, J. J. (1975) A structure-oriented program editor: a first step towards computer assisted programming. *International Computing Symposium (ICS'75)*.
- Dymetman, M., Lux, V. and Ranta, A. (2000) XML and multilingual document authoring: Convergent trends. *COLING*, pp. 243–249. Saarbrücken, Germany.

- Earley, J. (1970) An efficient context-free parsing algorithm. *Comm. ACM*, **13**(2), 94–102.
- Fiedler, A. (2001) *User-Adaptive Proof Explanation*. PhD thesis, Universität des Saarlandes.
- Frost, R. and Launchbury, J. (1989) Constructing natural language interpreters in a lazy functional language. *Comput. J.* **32**(2), 108–121.
- Hallgren, T. (2000) *Home Page of the Proof Editor Alfa*. <http://www.cs.chalmers.se/~hallgren/Alfa/>
- Hallgren, T. and Ranta, A. (2000) An extensible proof text editor. In: Parigot, M. and Voronkov, A. (eds.), *LPAR-2000: LNCS/LNAI 1955*, pp. 70–84. Springer-Verlag.
- Harper, R., Honsell, F. and Plotkin, G. (1993) A Framework for Defining Logics. *J. ACM*, **40**(1), 143–184.
- Hähnle, R., Johannisson, K. and Ranta, A. (2002) An authoring tool for informal and formal requirements specifications. In: Kutsche, R.-D. and Weber, H. (eds.), *Fundamental Approaches to Software Engineering: LNCS 2306*, pp. 233–248. Springer-Verlag.
- Hopcroft, J. and Ullman, J. (1979) *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley.
- Hudak, P. (1996) Building domain-specific embedded languages. *ACM Comput. Surv.*, **28**(4).
- Huet, G. (1997) The Zipper. *J. Functional Program.* **7**(5), 549–554.
- Huet, G. (2000) *Sanskrit site*. Program and documentation, <http://pauillac.inria.fr/~huet/SKT/>.
- Johannisson, K. and Ranta, A. (2001) Formal verification of multilingual instructions. *The Joint Winter Meeting of Computing Science and Computer Engineering*. Chalmers University of Technology.
- Johnson, S. C. (1975) *Yacc — yet another compiler compiler*. Technical report CSTR-32, AT & T Bell Laboratories, Murray Hill, NJ.
- Jones, M. and Hudak, P. (1995) Using types to parse natural language. *Proceedings Glasgow Workshop on Functional Programming*.
- Jones, N. D., Gomard, C. K. and Sestoft, P. (1993) *Partial Evaluation and Automatic Program Generation*. Prentice-Hall.
- Joshi, A. (1985) Tree-adjointing grammars: How much context-sensitivity is required to provide reasonable structural descriptions. In: Dowty, D., Karttunen, L. and Zwicky, A. (eds.), *Natural Language Parsing*, pp. 206–250. Cambridge University Press.
- Kahn, G., Lang, B., Mélése, B. and Morcos, E. (1983) Metal: a formalism to specify formalisms. *Sci. Comput. Program.* **3**, 151–188.
- Kay, M. (1997) The Proper Place of Men and Machines in Language Translation. *Machine Translation*, **12**(1–2), 3–23.
- Knuth, D. (1965) On the translation of languages from left to right. *Infor. Control*, **8**, 607–639.
- Knuth, D. (1968) Semantics of context-free languages. *Math. Syst. Theory*, **2**, 127–145.
- Lambek, J. (1958) The mathematics of sentence structure. *Am. Math. Monthly*, **65**, 154–170.
- Luo, Z. and Callaghan, P. (1999) Mathematical vernacular and conceptual well-formedness in mathematical language. In: Lecomte, A., Lamarche, F. and Perrier, G. (eds.), *Logical Aspects of Computational Linguistics (LACL): LNCS/LNAI 1582*, pp. 231–250.
- Luo, Z. and Pollack, R. (1992) *LEGO Proof Development System*. Technical report, University of Edinburgh.
- Mäenpää, P. and Ranta, A. (1999) The type theory and type checker of GF. *PLI-1999: Workshop on Logical Frameworks and Meta-languages*, Paris, France.
- Magnusson, L. and Nordström, B. (1994) The ALF proof editor and its proof engine. In: *Types for Proofs and Programs: LNCS 806*, pp. 213–237. Springer-Verlag.

- Marlow, S. (2001) *Happy, The Parser Generator for Haskell*. <http://www.haskell.org/happy/>.
- Martin-Löf, P. (1984) *Intuitionistic Type Theory*. Napoli: Bibliopolis.
- Montague, R. (1974) *Formal Philosophy*. Yale University Press.
- Morrill, G. (1994) *Type Logical Grammar*. Kluwer.
- Necula, G. C. (1997) Proof-Carrying Code. *Proc. 24th ACM Symposium on Principles of Programming Languages*, pp. 106–119. Paris, France. ACM Press.
- Nordström, B., Petersson, K. and Smith, J. M. (1990) *Programming in Martin-Löf's Type Theory. An Introduction*. Oxford University Press.
- Paulson, L. (2002) *The Isabelle Reference Manual*. Available at the Isabelle homepage <http://www.cl.cam.ac.uk/Research/HVG/Isabelle/> With contributions by T. Nipkow and M. Wenzel.
- Pentus, M. (1993) Lambek grammars are context-free. *LICS*, pp. 35–43. Utrecht, The Netherlands.
- Pereira, F. and Warren, D. (1980) Definite clause grammars for language analysis—a survey of the formalism and a comparison with augmented transition networks. *Artif. Intell.*, **13**, 231–278.
- Peyton Jones, S. and Hughes, J. (1999) *Report on the Programming Language Haskell 98, a Non-strict, Purely Functional Language*. Available from <http://www.haskell.org>
- Pollard, C. and Sag, I. (1994) *Head-Driven Phrase Structure Grammar*. University of Chicago Press.
- Power, R. and Scott, D. (1998) Multilingual authoring using feedback texts. *COLING-ACL*.
- Prawitz, D. (1965) *Natural Deduction*. Stockholm: Almqvist & Wiksell.
- Ranta, A. (1994) *Type Theoretical Grammar*. Oxford University Press.
- Ranta, A. (2002) *Grammatical Framework Homepage*. [www.cs.chalmers.se/~aarne/GF/](http://www.cs.chalmers.se/~aarne/GF/).
- Ranta, A. and Cooper, R. (2001) Dialogue systems as proof editors. *IJCAR/ICoS-3*.
- Rosetta, M. T. (1994) *Compositional translation*. Kluwer.
- Shieber, S. (1986) *An Introduction to Unification-Based Approaches to Grammars*. University of Chicago Press.
- Steedman, M. (1988) Combinators and grammars. In: Oehrlé, R., Bach, E. and Wheeler, D. (eds.), *Categorial Grammars and Natural Language Structures*, pp. 417–442. Reidel.
- Teitelbaum, T. and Reps, T. (1981) The Cornell Program Synthesizer: a syntax-directed programming environment. *Commun. ACM*, **24**(9), 563–573.
- The Boeing Company (2001) *Boeing Simplified English Checker*. <http://www.boeing.com/assocproducts/sechecker/>
- The Coq Development Team. (1999) *The Coq Proof Assistant Reference Manual*. [pauillac.inria.fr/coq/](http://pauillac.inria.fr/coq/).
- Warmer, J. and Kleppe, A. (1999) *The Object Constraint Language: Precise Modelling with UML*. Addison-Wesley.