

1

Modular Grammar Design with Typed Parametric Principles

1.1 Introduction

Extensible Dependency Grammar (XDG) (Debusmann et al., 2004) is a general framework for dependency grammar, with multiple levels of linguistic representations called *dimensions*. Its approach, motivated by the dependency grammar paradigm (Tesnière, 1959, Mel'čuk, 1988), is articulated around a description language for multi-dimensional attributed labeled graphs. XDG is a generalization of *Topological Dependency Grammar* (TDG) (Duchier and Debusmann, 2001).

For XDG, a grammar is a constraint that describes the valid linguistic signs as n -dimensional attributed labeled graphs, i.e. n -tuples of graphs sharing the same set of attributed nodes, but having different sets of labeled edges. It is central to XDG that all aspects of these signs are stipulated explicitly by *principles*: the class of models for each dimension, additional properties that they must satisfy, how one dimension must relate to another, and even lexicalization.

Yet, no formal account set in the XDG framework has so far explained what exactly these principles are, nor how they can be brought to bear on specific dimensions. In this paper, we show how an XDG grammar can be formally assembled from modular components called *parametric principles*. This yields a modular and compositional approach to grammar design. Compositional coherence is ensured by a *type system* whose primary novelty is to accommodate the notion of multi-dimensional

graphs. Instantiation of parametric principles not only imposes grammatical constraints, but, through the type system, also determines the necessary structure of grammatical signs. In this perspective, a grammar framework is simply a library of parametric principles such as the one offered by the XDG *Development Kit* (XDK) (Debusmann and Duchier, 2004).

1.2 Extensible Dependency Grammar

We briefly illustrate the XDG approach with an example of the German subordinate sentence “(dass) einen Mann Maria zu lieben versucht”.¹ Figure 1 shows an analysis with two dimensions: ID models grammatical function and LP word-order using *topological fields* (Duchier and Debusmann, 2001, Gerdes and Kahane, 2001)

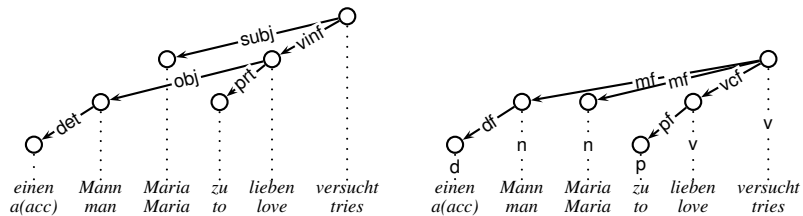


FIGURE 1 Example XDG analysis, ID left, LP right

Both dimensions share the same set of nodes (circles) but have different edges. On ID, *Maria* is subject (*subj*) of control verb *versucht*, and *Mann* object (*obj*) of *lieben*. On LP, *Mann* and *Maria* are both in the *Mittelfeld* (*mf*) of *versucht*, i.e. *Mann* has *climbed* to the finite verb.

1.3 Type System

Formalization. Let \mathcal{V} be an infinite set of node variables, \mathcal{L} a finite set of labels, and \mathcal{A} a set of attributes. An *attributed labeled graph* (V, E, A) consists of a finite set of nodes $V \subseteq \mathcal{V}$, a finite set of labeled directed edges $E \subseteq V \times V \times \mathcal{L}$ between them, and an assignment $A : V \rightarrow \mathcal{A}$ of attributes to nodes. An *n-dimensional attributed labeled graph* $((V, E_1, A_1), \dots, (V, E_n, A_n))$, or *multigraph*, is a *n*-tuple of labeled attributed graphs (V, E_i, A_i) over the same set of nodes V .

To provide a typed account of the XDG framework, we need a satisfactory type for graphs. A first idea is $(\{\mathcal{V}\}, \{\mathcal{V} \times \mathcal{V} \times \mathcal{L}\}, \mathcal{L} \rightarrow \mathcal{A})$, where we write $\{\tau\}$ for *set of* τ , but such a type is very imprecise: it fails to express that the nodes used in the edges (2nd arg) are elements

¹(that) Mary tries to love a man – see (Duchier and Debusmann, 2001)

of the graph's set of nodes (1st arg). Since additionally element graphs of a multigraph must be defined over the *same* set of nodes, some form of dependent typing appears inescapable. Given that general systems of dependent types tend to make type-checking undecidable, we sketch instead a restricted system, using kinds, sufficient for our purpose.²

1.3.1 A system of restricted dependent types

We assume given a number of disjoint sets of symbols D_i called finite domain kinds. Among them, two are distinguished: **Nodes** and **Labels**. Much of our type and kinding systems are quite standard. For reasons of space, we only detail the parts which are original to our proposal, and write $\tau :: \kappa$ and $e : \tau$ for the judgments that type τ has kind κ and expression e has type τ , omitting kinding and typing contexts.

$\kappa ::=$	\star	top	$D_i \sqsubset \star$	$G_v \sqsubset \star$	$M_v \sqsubset \star$
	D_i	domain			
	$G_{c_1 \dots c_k}$	graph			$\frac{\tau :: \kappa \quad \kappa \sqsubset \kappa'}{\tau :: \kappa'}$
	$M_{c_1 \dots c_k}$	multigraph			
	...				

FIGURE 2 Kinds and subkinds

$\tau, \tau_i ::=$	$c_1 \dots c_k$	domain	$[f_1 : \tau_1, \dots, f_n : \tau_n]$	record
	$\{\tau\}$	set	graph $\tau_1 \tau_2 \tau_3$	graph
	(τ_1, \dots, τ_n)	tuple	$\llbracket f_1 : \tau_1, \dots, f_n : \tau_n \rrbracket$	multigraph
	$\tau_1 \rightarrow \tau_2$	function	grammar τ	grammar
	$!\tau$	singleton	...	

FIGURE 3 Types

Finite domain sum types. They are built from symbols drawn from finite domain kinds: we write $c_1|\dots|c_k$ for a finite domain sum type in kind D_i and ϵ_{D_i} for its empty sum. Here are their kinding (left) and typing (right) rules:

$$\frac{c_1, \dots, c_k \text{ symbols in } D_i}{c_1|\dots|c_k :: D_i} \qquad \frac{c_1|\dots|c_k :: D_i}{c_i : c_1|\dots|c_k}$$

Kinds $G_{c_1|\dots|c_k}$ and $M_{c_1|\dots|c_k}$ indexed by finite domain sums make it possible to have types that depend on specific subsets of **Nodes** or **Labels**

²This presentation is based on ongoing work with Andreas Rossberg.

4 /

without requiring more general term-dependent types.

Singleton types. If $c_1 | \dots | c_k :: \text{Nodes}$ is the type of nodes of a graph, then $\{c_1 | \dots | c_k\}$ is the type of a set of these nodes. This is not sufficiently precise to type the set of nodes of the graph because it doesn't express that the latter is a maximal set. To achieve this aim, we adopt and adapt a limited version of *singleton types* (Stone, 2000). We write $!(c_1 | \dots | c_k)$ for the type inhabited by the single value $\{c_1, \dots, c_k\}$. Here are the relevant kinding, typing, and subtyping rules:

$$\frac{\tau :: D_i}{!\tau :: \star} \quad \frac{!(c_1 | \dots | c_k) :: \star}{\{c_1, \dots, c_k\} : !(c_1 | \dots | c_k)} \quad !\tau \sqsubseteq \{\tau\}$$

Graph types. A graph is defined from finite domain types ν and ℓ for its nodes and labels, and a type a for its attributes. We write G_ν for the kind of a graph over node type ν . Kinding and typing rules are:

$$\frac{\nu :: \text{Nodes} \quad \ell :: \text{Labels} \quad a :: \star}{\text{graph } \nu \ell a :: G_\nu} \quad \frac{V : !\nu \quad E : \{(\nu, \nu, \ell)\} \quad A : \nu \rightarrow a}{\text{Graph } V E A : \text{graph } \nu \ell a}$$

Multigraph types. A multigraph is defined from graphs over the same finite domain ν of nodes. We write M_ν for its kind. Here are the kinding and typing rules:

$$\frac{g_1 :: G_\nu \quad \dots \quad g_n :: G_\nu}{\llbracket f_1 : g_1, \dots, f_n : g_n \rrbracket :: M_\nu} \quad \frac{G_1 : g_1 :: G_\nu \quad \dots \quad G_n : g_n :: G_\nu}{\llbracket f_1 = G_1, \dots, f_n = G_n \rrbracket : \llbracket f_1 : g_1, \dots, f_n : g_n \rrbracket :: M_\nu}$$

Grammars. An XDG *grammar* is a set of predicates over the same multigraph type:

$$\frac{\tau :: M_\nu}{\text{grammar } \tau :: \star} \quad \frac{\tau :: M_\nu \quad S : \{\tau \rightarrow \text{Bool}\}}{\text{Grammar } S : \text{grammar } \tau}$$

Note that, in order to match our intuitions about grammars, the grammar type should be polymorphic in the finite domain type for nodes, otherwise the number of nodes is fixed. This can be achieved by extending the kinding system to admit *kind schemes* G_δ and M_δ where δ is a domain variable. An XDG *framework* is a set, also called a library, of principle templates.

1.4 Typed Templates

Attributes are usually given in the form of attribute/value matrices and principles parametrized by values which can be found at specific feature paths. For example, on the syntactic dimension $(V, E_{\text{ID}}, A_{\text{ID}})$ the agreement tuple assigned to each word must be one of those licensed by its lexical entry:

$$\text{Agr } \llbracket \text{id} = \text{Graph } V E_{\text{ID}} A_{\text{ID}} \rrbracket = \quad \forall v \in V : A_{\text{ID}}(v).\text{agr} \in A_{\text{ID}}(v).\text{lex.agsr}$$

We can generalize this into a reusable principle by abstracting over feature paths using access functions:

$$\text{Elem } D \ F_1 \ F_2 \ M \quad = \quad \mathbf{let \ val \ Graph } \ V \ E \ A = D(M) \ \mathbf{in} \\ \quad \quad \quad \forall v \in V : F_1(A(v)) \in F_2(A(v)) \\ \quad \quad \quad \mathbf{end}$$

but this is not very legible and, for notational convenience, we explore here an alternative that we call *templates*:

$$\text{Elem} \langle d, p_1, p_2 \rangle \llbracket d = \mathbf{Graph } \ V \ E \ A \rrbracket \quad = \quad \forall v \in V : A(v).p_1 \in A(v).p_2$$

where d, p_1, p_2 are *feature path variables*. A feature path is a (possibly empty) sequence of features. We write π for a path, ϵ for the empty path, and $\pi_1\pi_2$ for the concatenation of π_1 and π_2 . It is straightforward to generalize the language by allowing feature paths or feature path variables in types and patterns (and by extension in record ‘dot’ access) where previously only features were allowed. The intuition of such an extension lies in the congruence $[\epsilon : \tau] \equiv \tau$, $[\pi_1\pi_2 : \tau] \equiv [\pi_1 : [\pi_2 : \tau]]$, and in the interpretation of a dot access $.\pi$ as a postfix function of type $[\pi : \tau] \rightarrow \tau$.

Note that, if we write $\mathbf{graph } \ \nu \ \ell \ a$ for the type of Elem’s argument graph, it is our intention that type inference should require a to match the pattern $[p_1 : \tau, p_2 : \{\tau\}, \dots]$. This can be achieved either with a type system supporting record polymorphism or by adopting an open-world semantics³ for records and multigraphs, à la ψ -terms of LIFE. For simplicity, in this article we choose the latter.

We write $\langle p_1, p_2 \rangle \rightsquigarrow \tau$ for the type of a template abstracting over feature path variables p_1 and p_2 ; t may contain occurrences of p_1 and p_2 . We write $\tau_1 :: \kappa_1, \dots, \tau_n :: \kappa_n \Rightarrow \tau$ to express kinding constraints on the free type variables of τ . The type of Elem is then:

$$\text{Elem} : \langle d, p_1, p_2 \rangle \rightsquigarrow \nu :: \text{Nodes}, \ell :: \text{Labels}, \tau :: \star \Rightarrow \\ \llbracket d : \mathbf{graph } \ \nu \ \ell \ [p_1 : \tau, p_2 : \{\tau\}] \rrbracket \rightarrow \text{Bool}$$

1.5 Parametric Principles

We now illustrate how typed templates are intended to be used to define *parametric principles*.⁴ We write $v \xrightarrow{l}_E v'$ for an edge $(v, v', l) \in E$, $v \rightarrow_E v'$ for one with any label, $v \xrightarrow{+}_E v'$ for the transitive closure, and $v \rightarrow_E D$ for the type-raised relation where $D = \{v' \mid \forall v \rightarrow_E v'\}$, etc. . .

Tree Principle. It stipulates the set L of edge labels and requires that 1) each node has at most one incoming edge, 2) there is precisely

³no closed arities

⁴For lack of space, we present only a few principles. For further information, the reader is referred to e.g. (Debusmann et al., 2004, Debusmann and Duchier, 2004).

6 /

one node with no incoming edge (one root), and 3) there are no cycles:

$$\begin{aligned} \text{Tree}\langle d \rangle L \llbracket d = \text{Graph } V \ E \ A \rrbracket = \\ \forall v \in V \quad : \quad \forall M \subseteq V : M \xrightarrow{E} v \Rightarrow |M| \leq 1 \quad \wedge \\ \exists ! v \in V \quad : \quad \forall M \subseteq V : M \xrightarrow{E} v \Rightarrow |M| = 0 \quad \wedge \\ \forall v \in V \quad : \quad \forall D \subseteq V : v \xrightarrow{+E} D \Rightarrow v \notin D \end{aligned}$$

Here is the type constraint assigning singleton type $!l$ to L :

$$\text{Tree} : \langle d \rangle \rightsquigarrow \nu :: \text{Nodes}, \ell :: \text{Labels} \Rightarrow !\ell \rightarrow \llbracket d : \text{graph } \nu \ \ell \ _ \rrbracket \rightarrow \text{Bool}$$

Valency Principle. Incoming and outgoing edges must comply with the *in* and *out* valencies which stipulate, for each label $\ell \in L$, how many edges labeled with ℓ are licensed:

$$\begin{aligned} \text{Valency}\langle d, p_{\text{in}}, p_{\text{out}} \rangle L \llbracket d = \text{Graph } V \ E \ A \rrbracket = \\ \forall v \in V \forall \ell \in L \quad : \quad \forall M \subseteq V : M \xrightarrow{\ell} v \Rightarrow |M| \in A(v).p_{\text{in}}.\ell \quad \wedge \\ \forall v \in V \forall \ell \in L \quad : \quad \forall D \subseteq V : v \xrightarrow{\ell} D \Rightarrow |D| \in A(v).p_{\text{out}}.\ell \end{aligned}$$

Writing \mathbb{N} for the type of natural numbers, here is the type constraint:

$$\begin{aligned} \text{Valency} : \langle d, p_{\text{in}}, p_{\text{out}} \rangle \rightsquigarrow \nu :: \text{Nodes}, \ell :: \text{Labels} \Rightarrow \\ !\ell \rightarrow \llbracket d : \text{graph } \nu \ \ell \ [p_{\text{in}} : \ell \rightarrow \{\mathbb{N}\}, p_{\text{out}} : \ell \rightarrow \{\mathbb{N}\}] \rrbracket \rightarrow \text{Bool} \end{aligned}$$

Climbing Principle. Originating from TDG, this principle expresses that the tree-shape on dimension d_1 is a flattening of that of d_2 : 1) the dominance relation on d_1 must be a subset of that on d_2 , 2) on d_1 , each node must land on its d_2 -mother or climb higher:

$$\begin{aligned} \text{Climbing}\langle d_1, d_2 \rangle \llbracket d_1 = \text{Graph } V \ E_1 \ A_1, d_2 = \text{Graph } V \ E_2 \ A_2 \rrbracket = \\ \forall v \in V \quad : \quad \forall D_1, D_2 \subseteq V : v \xrightarrow{+E_1} D_1 \wedge v \xrightarrow{+E_2} D_2 \Rightarrow D_1 \subseteq D_2 \quad \wedge \\ \forall U_1, U_2 \subseteq V : U_1 \xrightarrow{+E_1} v \wedge U_2 \xrightarrow{*E_1} \rightarrow_{E_2} v \Rightarrow U_1 \subseteq U_2 \\ \text{Climbing} : \langle d_1, d_2 \rangle \rightsquigarrow \nu :: \text{Nodes} \Rightarrow \\ \llbracket d_1 : \text{graph } \nu \ _ \ _, d_2 : \text{graph } \nu \ _ \ _ \rrbracket \rightarrow \text{Bool} \end{aligned}$$

Lexicon Principle. XDG grammars are typically *lexicalized*: the record assignments to nodes are then partially determined by a *lexicon Lex*. The lexicon principle stipulates that each node must be assigned a *lexical entry*:

$$\text{Lexicon}\langle d \rangle \text{Lex} \llbracket d = \text{Graph } V \ E \ A \rrbracket = \forall v \in V : A(v) \in \text{Lex}$$

Here is the corresponding type constraint, stipulating that the graph on dimension d has no edges (since this dimension has only the purpose to carry the lexical entries):

$$\text{Lexicon} : \langle d \rangle \rightsquigarrow \nu :: \text{Nodes}, \tau :: \star \Rightarrow \{\tau\} \rightarrow \llbracket d : \text{graph } \nu \ \epsilon_{\text{Labels}} \ \tau \rrbracket \rightarrow \text{Bool}$$

Lookup Principle. A lexical entry is normally a record having a feature for each dimension. The *lookup* principle looks up a lexical entry's subrecord for a particular dimension and equates it with the p_{lex} feature of the node's attributes:

$$\text{Lookup}\langle d_1, d_2, p_{\text{lex}} \rangle \llbracket d_1 = \text{Graph } V \ E_1 \ A_1, d_2 = \text{Graph } V \ E_2 \ A_2 \rrbracket = \\ \forall v \in V : A_1(v).p_{\text{lex}} = A_2(v).d_1$$

$$\text{Lookup} : \langle d_1, d_2, p_{\text{lex}} \rangle \rightsquigarrow \nu :: \text{Nodes}, \tau :: \star \Rightarrow \\ \llbracket d_1 : \text{graph } \nu \ - [p_{\text{lex}} : \tau], d_2 : \text{graph } \nu \ - [d_1 : \tau] \rrbracket \rightarrow \text{Bool}$$

1.6 Example ID/LP Grammar

We now describe how the grammar of (Duchier and Debusmann, 2001) can be assembled in our typed framework of parametric principles. To better illustrate the compositionality of our approach, we adopt an incremental presentation that derives more complex grammars from simpler ones through operations of composition and restriction:

$$\begin{aligned} (\text{Grammar } S_1) ++ (\text{Grammar } S_2) &= \text{Grammar } (S_1 \cup S_2) \\ (\text{Grammar } S_1) // S_2 &= \text{Grammar } (S_1 \cup S_2) \end{aligned}$$

The grammar requires 3 dimensions: ID for syntax, LP for topology, and LEX for parametrization by a lexicon. They are respectively characterized by the following sets of labels:

$$\begin{aligned} L_{\text{ID}} &= \{\text{det, subj, obj, vbse, vpvt, vinf, prt}\} \\ L_{\text{LP}} &= \{\text{d, df, n, mf, vcf, p, pf, v, vxf}\} \\ L_{\text{LEX}} &= \emptyset \end{aligned}$$

Figure 4 shows how grammars for the 3 individual dimensions (G_{ID} , G_{LP} , G_{LEX}) are stipulated by instantiation of principles. For example, for G_{ID} , instantiation of $\text{Tree}(\text{id})$ L_{ID} has three consequences: (1) signs are required to match $\llbracket \text{id} = X \rrbracket$, i.e. to have an id dimension, (2) the directed graph on that dimension must have labels in L_{ID} , (3) this graph must be a tree. $G_{\text{ID}+\text{LEX}}$ is a grammar where the ID dimension is restricted by lexical constraints from the LEX dimension: it is obtained by combining G_{ID} and G_{LEX} using the composition operator $++$, and connecting them by the $\text{Lookup}(\text{id}, \text{lex})$ principle using the restriction operator $//$. Similarly for $G_{\text{LP}+\text{LEX}}$. Finally the full grammar $G_{\text{ID}/\text{LP}}$ is obtained by combining $G_{\text{ID}+\text{LEX}}$ and $G_{\text{LP}+\text{LEX}}$ and mutually constraining them by the Climbing and Barriers principles.

8 /

$$\begin{aligned}
G_{\text{LEX}} &= \text{Grammar } \{\text{Lexicon}\langle \text{lex} \rangle \text{ Lex}\} \\
G_{\text{ID}} &= \text{Grammar } \{\text{Tree}\langle \text{id} \rangle L_{\text{ID}}, \text{Valency}\langle \text{id}, \text{in}, \text{out} \rangle L_{\text{ID}}\} \\
G_{\text{ID+LEX}} &= (G_{\text{ID}} ++ G_{\text{LEX}}) // \{\text{Lookup}\langle \text{id}, \text{lex}, \text{lex} \rangle\} \\
G_{\text{LP}} &= \text{Grammar } \{\text{Tree}\langle \text{lp} \rangle L_{\text{LP}}, \text{Valency}\langle \text{lp}, \text{in}, \text{out} \rangle L_{\text{LP}}, \\
&\quad \text{Order}\langle \text{lp}, \text{on}, \text{order}, \text{pos}, \text{self} \rangle, \\
&\quad \text{Projectivity}\langle \text{lp}, \text{pos} \rangle\} \\
G_{\text{LP+LEX}} &= (G_{\text{LP}} ++ G_{\text{LEX}}) // \{\text{Lookup}\langle \text{lp}, \text{lex}, \text{lex} \rangle\} \\
G_{\text{ID/LP}} &= (G_{\text{ID+LEX}} ++ G_{\text{LP+LEX}}) // \\
&\quad \{\text{Climbing}\langle \text{lp}, \text{id} \rangle, \text{Barriers}\langle \text{lp}, \text{id}, \text{blocks} \rangle\}
\end{aligned}$$

FIGURE 4 Defining grammars in an incremental and compositional way

1.7 Conclusion

In this paper, we have made a threefold contribution: (1) we described a novel system of restricted dependent types capable of precisely describing graphs and multi-dimensional graphs, (2) we introduced a notion of typed templates with which we could express parametric principles, (3) we showed how these could enable a modular and compositional approach to grammar design. Finally we illustrated our proposal with an example reconstruction of an earlier grammar.

References

- Debusmann, Ralph and Denys Duchier. 2004. XDG Development Kit. <http://www.ps.uni-sb.de/~rade/mogul/publish/doc/debusmann-xdk/>.
- Debusmann, Ralph, Denys Duchier, Alexander Koller, Marco Kuhlmann, Gert Smolka, and Stefan Thater. 2004. A Relational Syntax-Semantics Interface Based on Dependency Grammar. In *Proceedings of COLING 2004*. Geneva/CH.
- Duchier, Denys and Ralph Debusmann. 2001. Topological Dependency Trees: A Constraint-Based Account of Linear Precedence. In *Proceedings of ACL 2001*. Toulouse/FR.
- Gerdes, Kim and Sylvain Kahane. 2001. Word Order in German: A Formal Dependency Grammar Using a Topological Hierarchy. In *ACL 2001 Proceedings*. Toulouse/FR.
- Mel'čuk, Igor. 1988. *Dependency Syntax: Theory and Practice*. Albany/US: State Univ. Press of New York.
- Stone, Christopher Allan. 2000. *Singleton Kinds and Singleton Types*. Ph.D. thesis, Carnegie Mellon University.
- Tesnière, Lucien. 1959. *Eléments de Syntaxe Structurale*. Paris/FR: Klincksieck.