

The TIGER language.

A Description Language for Syntax Graphs.

- Formal Definition -

Esther König and Wolfgang Lezius

Institut für Maschinelle Sprachverarbeitung (IMS)
University of Stuttgart, Germany
www.ims.uni-stuttgart.de/projekte/TIGER

April 22, 2003

Contents

1	Introduction	3
1.1	Acknowledgements and References	3
2	Syntax Graph Model	4
2.1	Feature Records	4
2.2	Nodes	4
2.3	Syntax Graphs	5
2.4	Syntax Graph Relations	7
2.5	Syntax Graph Models	7
3	Syntax and Interpretation of the TIGER Language	8
3.1	Elementary Symbols	8
3.2	Feature Values	10
3.3	Feature Constraints	10
3.4	Node Descriptions	11
3.5	Graph Descriptions	11
3.6	Satisfaction of local constraints	14
3.7	Corpus Definition	14
3.7.1	Feature Values	14
3.7.2	Feature Constraints	14
3.7.3	Node Descriptions	15
3.7.4	Graph Descriptions	15

3.7.5	Interpretation of a corpus definition	15
3.8	Queries	16
3.9	Types	16
3.10	Template Definitions	16
3.11	Semantic Consequence	17
4	TIGER calculus	18
4.1	Feature Values	20
4.1.1	Disjunctive Normal Form	20
4.1.2	Disjunction reduction	20
4.1.3	Consistency check	20
4.2	Feature Constraints	22
4.2.1	Disjunctive Normal Form	22
4.2.2	Disjunction reduction	23
4.2.3	Feature normalization	23
4.2.4	Consistency check	24
4.3	Node Descriptions	26
4.3.1	Disjunctive Normal Form	26
4.3.2	Node description normalization	26
4.3.3	Node guess	27
4.4	Graph Descriptions	27
4.4.1	Disjunctive Normal Form	27
4.4.2	Disjunction reduction	27
4.4.3	Graph description normalization	27
4.4.4	Consistency check	28
4.4.5	Templates	30
4.4.6	Queries	31

Chapter 1

Introduction

This report gives a formal definition of a description language for *syntax graphs*, called the TIGER language. This formal definition complements the informal introduction to the TIGER language which is given in the TIGERSearch user's manual [12].

Syntax graphs are close relatives to (syntax) trees, to feature structures or dependency graphs. Syntax graphs are syntax trees with two additions: Edges may carry labels, and crossing edges are permitted. On the other hand, syntax graphs differ from feature structures due to the following two properties: The terminal nodes are ordered (like in syntax trees). Albeit edges may cross, two edges must not join in a common node. This means that the structure sharing mechanism of feature structures is ruled out (at least in the kernel of the TIGER language).

Chap. 2 defines the syntax graph model. The syntax of the TIGER language and its interpretation with respect to syntax graphs is given in Chap. 3. A calculus for formula reduction is introduced in Chap. 4, which serves as the basis for the implementation of the TIGER language.

1.1 Acknowledgements and References

The design and the formal definition of the TIGER language has been heavily influenced by other work on formal languages, in particular

- unification-based formalisms, cf. [5], [4], [6], [9], [8], [11], [15]
- tree description languages, cf. [1], [7], [13], [14],
- corpus query languages, cf. [2], [3]

and by some of the literature cited in the publications above.

Our work has been partially funded by the DFG (project TIGER) and by the Ministry of Science, Research and the Arts of the State of Baden-Württemberg (project DEREKO).

Chapter 2

Syntax Graph Model

One component of a *syntax graph model* is the set of *syntax graphs*. A (syntax) graph G is a relation over a set of *nodes* V . Whereas, from the point of view of the graph, a node is an atomic element, a node can still have its own life, i.e. it can possess internal structure. For example, in the case of TIGER syntax graphs, a node is a pair of a *node identifier* and a *feature record*.

2.1 Feature Records

According to common practice in computational linguistics, *features* are functions. In our case, their range is restricted to constants. For this reason, we prefer the term *feature record* instead of *feature structure* (where values of features are again arbitrary feature structures).

Definition 1 (Feature Record)

A feature record F is a set of feature maps $F = \{f_1, \dots, f_m\}$ where $f_i : \{d \in D\} \rightarrow C_i$ and C_i is a set of constants. D is a non-empty set, the domain of feature records \mathcal{F} .

2.2 Nodes

Definition 2 (Node)

A node $v \in V$ is a pair $v = \langle \nu, F \rangle$ of a node identifier ν and a feature record F .

This means that a node is unique due to its identifier.

2.3 Syntax Graphs

The following definition specifies the properties of syntax graphs $G \in \mathcal{G}$, where \mathcal{G} is the set of syntax graphs.

Definition 3 (Syntax Graph)

A syntax graph G is a sextuple $G = (V_{NT}, V_T, L_G, E_G, O_G, R_G)$ where

1. V_{NT} is a possibly empty set of nonterminal or inner nodes.
2. V_T is a non-empty set of terminal nodes or leaves.
Notation: $V_G = V_{NT} \cup V_T$ denotes the set of all nodes of a graph G .
3. L_G is a (possibly empty) set of edge labels.
4. E_G are the labeled, directed edges of G . E_G is a set of relations $l, l \in L_G$, with $\langle V_{NT}, (V_{NT} \cup V_T) \rangle$,
5. O_G is a bijective order transformation $O_G : V_T \rightarrow \{1, \dots, |V_T|\}$
I.e. for any pair of terminal nodes $v_i, v_j \in V_T, v_i \neq v_j$, it can be decided whether $O_G(v_i) < O_G(v_j)$ or $O_G(v_i) > O_G(v_j)$.
6. $R_G \in (V_{NT} \cup V_T)$ is the root node.

G is a graph with the following restrictions:

1. G is a directed acyclic graph with exactly one root node R_G . A root node is a node without incoming edges.
2. For all nodes $v \in V_G$ different from the root node R_G , there is exactly one edge leading into v .
3. Each nonterminal node $v \in V_{NT}$ must have at least one successor, i.e.
 $v \in V_{NT}$ iff $\exists l, v' \langle v, v' \rangle \in l$.

Note that a single node $v \in V_T$ is trivially a syntax graph.

Similar to nodes, there must be a means to distinguish graphs. Since the identifier of a node is unique in the whole domain of nodes, we can simply take the identifier of the root node as the *identifier of a graph*. Furthermore, it follows that each node v_i belongs to a unique graph G_i .

There are few convenient definitions wrt. a syntax graph G .

Definition 4 (Path)

A path $p = \langle v_0, \dots, v_n \rangle$ in a syntax graph G is an n -tuple of nodes with $\langle v_0, v_1 \rangle \in l$, for some l , and either $v_1 = v_n$ or $p' = \langle v_1, \dots, v_n \rangle$ is a path, $0 < n$. The number n is called the length of the path p .

Definition 5 ($>_n$, dominance with distance n , $n \geq 0$)

A node v_1 dominates a node v_2 with distance n ($v_1 >_n v_2$) if there is a path $p = \langle v_1, \dots, v_n \rangle$ in G .

Definition 6 ($>_{@l}$, left corner)

A node $v_2 \in V_T$ is the left corner with respect to (the subgraph dominated by) a node v_1 ($v_1 >_{@l} v_2$) if

- either $v_2 = v_1$
- or $v_1 >_n v_2$
and for all nodes $v_i \in V_T$ with $v_1 > * v_i$ it holds that $O_G(v_i) > O_G(v_2)$.

The right corner v_2 of a node v_1 ($v_1 >_{@r} v_2$) is defined analogously.

The (generalized) precedence of two nodes v_1 and v_2 is defined with respect to the respective leftmost terminal nodes ('left corners') dominated by v_1 and v_2 . Since the precedence of terminal nodes is a total order, the order relation between two left corners is always defined.

Definition 7 ($.n$, precedence with distance n , $n \geq 0$)

A node v_1 precedes a node v_2 with distance n ($v_1 .n v_2$) if $n \geq 0$ and for the left corner v_3 of v_1 ($v_1 >_{@l} v_3$) with respect to the left corner v_4 of v_2 ($v_2 >_{@l} v_4$) it holds that $n = O_G(v_4) - O_G(v_3)$.

Definition 8 ($\$$, siblings)

A node v_1 is a sibling of a node v_2 ($v_1 \$ v_2$) if there exists a node v_0 which directly dominates both, v_1 and v_2 , i.e. $\exists v_0 \in V_{NT} : v_0 > v_1 \wedge v_0 > v_2 \wedge v_1 \neq v_2$.

Definition 9 (Arity n of a node with identifier ν)

The arity n of a node v with identifier ν is the number n of its daughters or outgoing edges, i.e.

$$\text{arity}(\nu) = |\{l | \langle v, v_l \rangle \in l\}|.$$

Definition 10 (Token arity n of a node with identifier ν)

The token arity n of a node v with identifier ν is the number n of terminal nodes which are dominated by v , i.e.

$$\text{tokenarity}(\nu) = |\{l \mid \langle v, \dots, v_l \rangle \text{ path in } G \text{ and } v_l \in V_T\}|.$$

Definition 11 (Continuous graph rooted in v)

The graph which is rooted in a node v with identifier ν is continuous if its terminal nodes form a continuous string, i.e. for all terminal leaves $v_i = v_1, \dots, v_n$ dominated by v it holds that

$$(\exists v_j O_G(v_j) = O_G(v_i) + 1) \text{ or } (\neg \exists v_j O_G(v_j) > O_G(v_i))$$

2.4 Syntax Graph Relations

In order to model types and templates, we need a notion of *relation* among constants, feature records, nodes, and syntax graphs. Types correspond to unary relations or predicates. Let \mathcal{U} be the universe of graphs and all their substructures, i.e. the (disjoint) set union of constants C , feature records \mathcal{F} , nodes V , and syntax graphs \mathcal{G} .

Definition 12 (Syntax Graph Relation)

An n -ary (syntax graph) relation $r \in R$ is an n -tuple ($n > 0$) $r \subseteq \mathcal{U} \times \dots \times \mathcal{U}$, where the individual subsets of \mathcal{U} must not be empty.

2.5 Syntax Graph Models

The notion of *syntax graph model* ties together all the previously introduced components.

Definition 13 (Syntax Graph Model)

A syntax graph model is a quintuple $\mathcal{M} = (C, \mathcal{F}, V, \mathcal{G}, R)$.

Chapter 3

Syntax and Interpretation of the TIGER Language

This section defines the syntax of the TIGER language in a formal manner. Along the way, the semantic interpretation of each TIGER expression will be defined with respect to a given syntax graph model \mathcal{M} (see Sec. 2). There are three levels of interpretation:

1. denotation of the *elementary symbols*
2. denotation of *local constraints*, i.e. of descriptions of feature values, feature constraints, node descriptions, graph descriptions; and the corresponding notion of constraint satisfaction
3. interpretation/satisfaction of *global definitions and constraints*, i.e. corpus definitions, template definitions, type definitions, and feature declarations.

3.1 Elementary Symbols

The inventory of elementary symbols is defined as follows:

Definition 14 (Signature)

Operators

- = (*feature value*),
- ! (*negation*), & (*conjunction*), | (*disjunction*),
- > (*direct dominance*), >* (*dominance*), >@l (*left corner*), >@r (*right corner*),
- . (*direct precedence*), .* (*precedence*),
- \$ (*siblings*), \$.* (*ordered siblings*)
- := (*type definition*)

Mark-up

;, <, >, /, domain, feature, name, range

Constants $c_1, c_2, \dots \in \mathbf{C}$

Feature Names $f_1, f_2, \dots \in \mathbf{F}$

Edge Labels $l_1, l_2, \dots \in \mathbf{L}$

Type Names $t_1, t_2, \dots, t_n, \text{Constant}, \text{FREC}, \text{Graph}, \text{Node}, \text{NT}, \text{T}, \top$ (*top*) $\in \mathbf{T}$,
and \perp (*bottom*).

Predicate Names *root, arity, tokenarity, continuous, discontinuous*

Template Names $r_1, r_2, \dots \in \mathbf{R}$

Constants have to be put into quotes "...". Type names must neither start with " nor with the #-symbol (in order to avoid confusion with variables, see below). Therefore, the sets \mathbf{C} and \mathbf{T} are disjoint.

We constrain the notion of syntax graph model (Sec. 2) by the condition that each element of the signature must have a semantic counterpart in the given syntax graph model \mathcal{M} . In order to improve readability, we omit indices, i.e. $c \in C$ means $c_{\mathcal{M}} \in C_{\mathcal{M}}$, etc.

Definition 15 (Denotation of elementary symbols)

Constants: *there is a constant $c \in C$ for each $c \in \mathbf{C}$*

Feature Names: *there is unary partial function f on D for each $f \in \mathbf{F}$*

Edge Labels: *there is an edge label $l \in L$ for each $l \in \mathbf{L}$*

Type Names: *there is a unary syntax graph relation $r_{\mathbf{t}} \in R$ for each $\mathbf{t} \in \mathbf{T}$.*

In particular, $r_{\text{Constant}} = C, r_{\text{FREC}} = \mathcal{F}, r_{\text{Node}} = \mathcal{V}, r_{\text{Graph}} = \mathcal{G}, r_{\top} = \mathcal{U}$.

In addition, $[\perp] = \emptyset$.

Template Names: *there is an n -ary syntax graph relation $r \in R$ for each n -ary template $\mathbf{r} \in \mathbf{R}$*

Furthermore, we introduce a general notion of variables. This means that the denotation of TIGER expressions has to be calculated with respect to a given variable substitution σ .

Definition 16 (Variable substitution)

Let \mathbf{X} be a set of variable names (variables) with members $\#x_1, \#x_2, \dots \in \mathbf{X}$. A variable substitution σ is a map from variable names on elements of the set \mathcal{U} . The function value $\sigma(\#x)$ is the element which is assigned to the variable $\#x$.

The sets of variable names for feature values, feature constraints, and node identifiers, respectively, must be pairwise disjoint.

Subsequently, we mean $\llbracket \gamma \rrbracket_{\mathcal{M}, \sigma}$ ('denotation of γ in model \mathcal{M} under variable substitution σ ') when writing $\llbracket \gamma \rrbracket$.

3.2 Feature Values

A feature value expression d denotes a subset of the set of constants C . A type \mathbf{t} is interpreted as a one-place relation $r_{\mathbf{t}}$, whose single argument must be a constant.

<i>Syntax</i>	<i>Explanation</i>	<i>Denotation</i> $\llbracket \cdot \rrbracket$
$d \Rightarrow c$	constant (string)	$\{c\}$
\mathbf{t}	type	$\{u \mid \langle u \rangle \in r_{\mathbf{t}}\} \cap C$
$! d$	negation	$C \setminus \llbracket d \rrbracket$
$d \ \& \ d$	conjunction	$\llbracket d_1 \rrbracket \cap \llbracket d_2 \rrbracket$
$d \ \ d$	disjunction	$\llbracket d_1 \rrbracket \cup \llbracket d_2 \rrbracket$

Variables may only be used in the place of feature values or as prefixes to feature value expressions, and of course, they must denote constants.

<i>Syntax</i>	<i>Explanation</i>	<i>Denotation</i> $\llbracket \cdot \rrbracket$
$d' \Rightarrow d$	feature value	$\llbracket d \rrbracket$
$\#x$	variable	$\{\sigma(\#x)\} \cap C$
$\#x:d$	constrained variable	$\{\sigma(\#x)\} \cap \llbracket d \rrbracket$

3.3 Feature Constraints

The syntax and denotation of a *feature constraint* k is defined below. In addition to the above mentioned feature value expressions, we also admit regular expressions over strings as feature values. However, a regular expression can neither be referred to by variable nor can it be embedded into a Boolean expression (since this could lead to very large database joins otherwise). We assume that there is a function `evalre` which evaluates a regular expression z and returns a set of constants. In the case of feature constraints, types must evaluate to feature records.

<i>Syntax</i>	<i>Explanation</i>	<i>Denotation</i> $\llbracket \cdot \rrbracket$
$\mathbf{k} \Rightarrow \mathbf{t}$	type	$\{u \mid \langle u \rangle \in r_{\mathbf{t}}\} \cap \mathcal{F}$
$\mathbf{f} = \mathbf{d}'$	feature-value pair	$\{F \mid f(F) \in \llbracket \mathbf{d}' \rrbracket\}$
$\mathbf{f} \neq \mathbf{d}$	negated feature-value	$\llbracket \mathbf{f} = \mathbf{d} \rrbracket$
$\mathbf{f} = /z/$	regular expression	$\{F \mid f(F) \in \text{evalre}(z)\}$
$\mathbf{f} \neq /z/$	neg. regular expression	$\{F \mid f(F) \notin \text{evalre}(z)\}$
$! \mathbf{k}$	negation	$F \setminus \llbracket \mathbf{k} \rrbracket$
$\mathbf{k} \& \mathbf{k}$	conjunction	$\llbracket \mathbf{k}_1 \rrbracket \cap \llbracket \mathbf{k}_2 \rrbracket$
$\mathbf{k} \mid \mathbf{k}$	disjunction	$\llbracket \mathbf{k}_1 \rrbracket \cup \llbracket \mathbf{k}_2 \rrbracket$

Similarly to feature values, feature constraints can be prefixed by variables. And finally, a (prefixed) feature constraint comes with square brackets at its outermost level.

<i>Syntax</i>	<i>Explanation</i>	<i>Denotation</i> $\llbracket \cdot \rrbracket$
$\mathbf{k}' \Rightarrow []$	empty constraint	\mathcal{F}
$[\mathbf{k}]$	feature constraint	$\llbracket \mathbf{k} \rrbracket$
$[\# \mathbf{x}]$	variable	$\{ \sigma(\# \mathbf{x}) \} \cap \mathcal{F}$
$[\# \mathbf{x} : \mathbf{k}]$	constrained variable	$\{ \sigma(\# \mathbf{x}) \} \cap \llbracket \mathbf{k} \rrbracket$

3.4 Node Descriptions

A *node description* is a pair of a node identifier \mathbf{n} and a (prefixed) feature constraint \mathbf{k}' where either the identifier or the feature constraint can be omitted. The denotation of a node description is actually a single node since the nodes in the syntax graph model come with unique identifiers. However, for reasons of conformity, we let the denotation of a node description be a singleton *set*.

<i>Syntax</i>	<i>Explanation</i>	<i>Denotation</i> $\llbracket \cdot \rrbracket$
$\mathbf{n} \Rightarrow \mathbf{c}$	constant node identifier	$\{c\}$
$\# \mathbf{x}$	variable node identifier	$\{ \sigma(\# \mathbf{x}) \} \cap C$
$\mathbf{v} \Rightarrow \mathbf{n}$	unconstrained node	$\{ \langle c, F \rangle \mid c \in \llbracket \mathbf{n} \rrbracket; F \in \mathcal{F} \}$
\mathbf{k}'	unidentified node	$\{ \langle c, F \mid c \in C; F \in \llbracket \mathbf{k}' \rrbracket \rangle \}$
$\mathbf{n} : \mathbf{k}'$	node description	$\{ \langle c, F \rangle \mid c \in \llbracket \mathbf{n} \rrbracket; F \in \llbracket \mathbf{k}' \rrbracket \}$

3.5 Graph Descriptions

Obviously, a *graph description* should denote a set of graphs G . The details of the definition are given in Fig. 3.5. First, we have to define the syntax of template calls (or template heads).

<i>Syntax</i>	<i>Explanation</i>	<i>Denotation</i> $\llbracket \cdot \rrbracket$
$\mathbf{u} \Rightarrow \#x$ $\#x, \mathbf{u}$	argument parameter arguments	$\sigma(\#x)$
$\mathbf{r}' \Rightarrow \mathbf{r}(\mathbf{u})$	template call/head	$\begin{cases} \mathcal{G} & \text{if } \langle \sigma(\#x_1), \dots, \sigma(\#x_n) \rangle \in r \\ \emptyset & \text{otherwise} \end{cases}$

<i>Syntax</i>	<i>Explanation</i>	<i>Denotation</i> $[\cdot]$
$g \Rightarrow \text{root}(n)$	root identification	$\{G \mid \llbracket n \rrbracket = R_G\}$
$\text{arity}(n, i_1)$	arity constraint	$\{G \mid \text{arity}(\llbracket n \rrbracket) = i_1\}$
$\text{arity}(n, i_1, i_2)$	arity interval	$\{G \mid i_1 \leq \text{arity}(\llbracket n \rrbracket) \leq i_2\}$
$\text{tokenarity}(n, i_1)$	# of terminal leaves	$\{G \mid \text{tokenarity}(\llbracket n \rrbracket) = i_1\}$
$\text{tokenarity}(n, i_1, i_2)$	interval for # of terminals	$\{G \mid i_1 \leq \text{tokenarity}(\llbracket n \rrbracket) \leq i_2\}$
$\text{continuous}(n)$	continuity constraint	$\{G \mid \text{continuous}(\llbracket n \rrbracket)\}$
$\text{discontinuous}(n)$	discontinuity constraint	$\mathcal{G} \setminus \{G \mid \text{continuous}(\llbracket n \rrbracket)\}$
r'	template call	
v	node description	$\{G \mid \llbracket v \rrbracket \subseteq V_G\}$
$v >_1 v$	labeled direct dominance	$\{G \mid l(v_1) = v_2\}$
$v > v$	direct dominance	$\{G \mid v_1 >_1 v_2\}$
$v >^* v$	dominance	$\{G \mid \exists i (v_1 >_i v_2)\}$
$v >_{i_3} v$	dominance, distance i_3	$\{G \mid v_1 >_{i_3} v_2\}$
$v >_{i_3, i_4} v$	dominance, distance $i_3 \dots i_4$	$\{G \mid \exists i (i_3 \leq i \leq i_4 \wedge v_1 >_i v_2)\}$
$v . v$	direct precedence	$\{G \mid v_1 .1 v_2\}$
$v .^* v$	precedence	$\{G \mid \exists i (v_1 .i v_2)\}$
$v >@l v$	left corner	$\{G \mid v_1 > @l v_2\}$
$v >@r v$	right corner	$\{G \mid v_1 > @r v_2\}$
$v .i_3 v$	precedence, distance i_3	$\{G \mid v_1 .i_3 v_2\}$
$v .i_3, i_4 v$	precedence, distance $i_3 \dots i_4$	$\{G \mid \exists i (i_3 \leq i \leq i_4 \wedge v_1 .i v_2)\}$
$v !R v$	negated node relation	$\mathcal{G} \setminus \{G \mid v_1 R v_2\}$
$v \$ v$	siblings	$\{G \mid v_1 \$ v_2\}$
$v \$.^* v$	siblings with precedence	$\{G \mid v_1 \$ v_2 \wedge \exists i (v_1 .i v_2)\}$
$g \& g$	conjunction	$\llbracket g_1 \rrbracket \cap \llbracket g_2 \rrbracket$
$g \mid g$	disjunction	$\llbracket g_1 \rrbracket \cup \llbracket g_2 \rrbracket$

Figure 3.1: Syntax and denotation of graph descriptions. $v_1 \in \llbracket v_1 \rrbracket$; $v_2 \in \llbracket v_2 \rrbracket$; $\llbracket v_1 \rrbracket, \llbracket v_2 \rrbracket \subseteq V_G$; n node identifier; $i_1, i_2 \geq 0$; $i_3, i_4 \geq 1$; R unnegated node relation ($\$.^*$ excluded)

3.6 Satisfaction of local constraints

A local constraint (or description), like a feature constraint or a syntax graph description, denotes a set of structures, and, if ultimately specific, a singleton set, i.e. one single structure. Satisfaction of a local constraint in a syntax graph model \mathcal{M} is therefore defined with reference to the denotation of the constraint in \mathcal{M} .

Definition 17 (Satisfaction, $\mathcal{M}, \sigma \models \gamma$, $\mathcal{M} \models \gamma$)

A constraint γ ($\gamma = \mathbf{d}', \mathbf{k}, \mathbf{v}, \mathbf{g}$) is satisfied in $\mathcal{M} = (C, \mathcal{F}, V, \mathcal{G}, R)$ under variable substitution σ ($\mathcal{M}, \sigma \models \gamma$) if $\llbracket \gamma \rrbracket_{\mathcal{M}, \sigma} \neq \emptyset$.

A constraint γ is satisfied in \mathcal{M} ($\mathcal{M} \models \gamma$) if there is some variable substitution σ such that $\mathcal{M}, \sigma \models \gamma$.

Now, we can turn to the global constraints and definitions.

3.7 Corpus Definition

A corpus definition \mathbf{k} is a non-empty list of restricted graph descriptions or *corpus graphs* \mathbf{g}^r . The restricted syntax is defined subsequently.

3.7.1 Feature Values

In corpus definitions, feature values \mathbf{d}^r must not include types or template calls. Boolean operators are excluded (negation and disjunction) or do not make sense (conjunction of constants).

$$\begin{aligned} \mathbf{d}^r &\Rightarrow \mathbf{c} \\ \mathbf{d}'^r &\Rightarrow \mathbf{d}^r \end{aligned}$$

3.7.2 Feature Constraints

In corpus definitions, feature values \mathbf{d}^r cannot be regular expressions. On the level of feature constraints \mathbf{k}^r , the empty feature constraint, and the Boolean operators $!$ and $|$ are excluded.

$$\begin{aligned} \mathbf{k}^r &\Rightarrow \mathbf{f} = \mathbf{d}'^r \\ &\quad | \quad \mathbf{k}^r \ \& \ \mathbf{k}^r \\ \mathbf{k}'^r &\Rightarrow [\mathbf{k}^r] \end{aligned}$$

3.7.3 Node Descriptions

A node description v^r in a corpus definition must be a pair of a constant node identifier and non-empty feature constraint.

$$\begin{aligned} n^r &\Rightarrow c \\ v^r &\Rightarrow n^r : k'^r \end{aligned}$$

3.7.4 Graph Descriptions

In a corpus definition, a graph description g^r must be conjunctions of basic node relations and fixed arity constraints.

$$\begin{aligned} g^r &\Rightarrow \text{root}(n) \\ &| \text{arity}(n, a) \\ &| v^r > 1 \ v^r \\ &| v^r \cdot v^r \\ &| g^r \ \& \ g^r \end{aligned}$$

Additional restrictions on corpus graphs are:

1. The precedence relation on terminal nodes must be a total order.
2. For each node, its arity has to be fixed.
3. There is a single node, the root node, which does not occur on the righthand side of any labeled direct dominance expression.
4. Each node (except for the root node) occurs exactly once as the righthand side of a labeled direct dominance expression.

3.7.5 Interpretation of a corpus definition

Whereas a conjunction of local constraints means the intersection of the respective denotations, this will not work e.g. for corpus definitions, since we would need a set of sets of syntax graphs as a component of the syntax graph model. In order to avoid that kind of complication, for global constraints, we immediately define the notion of satisfaction in a model \mathcal{M} .

<i>Syntax</i>	<i>Explanation</i>	<i>Interpretation</i>
$k \Rightarrow g^r$	single graph corpus	$\mathcal{M} \models g^r$
$g^r \ k$	corpus	$\mathcal{M} \models g^r$ and $\mathcal{M} \models k$

Note (for future extensions of corpus graphs): there may be distinct variable substitutions σ, σ' for testing the satisfaction of the individual corpus graphs $\mathbf{g}^r, \mathbf{g}^{r'}$, i.e. the scope of a variable is a single corpus graph.

3.8 Queries

A *query* is simply a graph description.

3.9 Types

A type hierarchy is defined by a set of *type definitions* or type axioms \mathbf{a} . A 'subtype' \mathbf{t}' may be either a real type or a constant. A type \mathbf{t} is the disjoint set union of its subtypes \mathbf{t}_i' . A *feature declaration* postulates the existence of a total function for its domain type. A type symbol \mathbf{t} must be defined at most once, i.e. it may not occur more than once as the lefthand side of a type definition. A type symbol must occur exactly once in a righthand side of a type definition. For a feature symbol \mathbf{f} , there must be exactly one feature declaration.

<i>Syntax</i>	<i>Explanation</i>	<i>Interpretation</i>
$\mathbf{t}' \Rightarrow \mathbf{c}$ \mathbf{t}		
$\mathbf{s} \Rightarrow \mathbf{t}'$ \mathbf{t}', \mathbf{s}	subtype subtypes	
$\mathbf{a} \Rightarrow \mathbf{t} := \mathbf{s} ;$	type definition	$[[\mathbf{t}]] = \bigcup_{1 \leq i < j \leq n} [[\mathbf{t}_i']]$ and $[[\mathbf{t}_i']] \cap [[\mathbf{t}_j']] = \emptyset,$ $1 \leq i < j \leq n$
<code><feature name="f" domain="t" range="t"></code> <code></feature></code>	feature declaration	f is a total function $[[\mathbf{t}_1]] \times [[\mathbf{t}_2]]$

3.10 Template Definitions

The syntax and semantics of a template defining clause \mathbf{s} is defined as follows.

<i>Syntax</i>	<i>Explanation</i>	<i>Interpretation</i>
$\mathbf{s} \Rightarrow \mathbf{r}' \leftarrow \mathbf{g} ;$	clause	$\exists \sigma: \text{if } \mathcal{M}, \sigma \models \mathbf{g} \text{ then } \mathcal{M}, \sigma \models \mathbf{r}'$

The scope of a variable covers a defining clause.

Note that the body of a clause is an embedded query. This means that graph descriptions always occur with the same 'polarity', i.e. only in queries, but never in template heads. Therefore, graph descriptions are interpreted in a uniform way (always in the 'if'-clause of the interpretation of a template definition).

3.11 Semantic Consequence

Definition 18 (TIGER database)

Let \mathbf{k} be a corpus definition, \mathbf{S} a set of template definitions, and \mathbf{A} a set of type definitions and feature declarations. A TIGER database DB is a triple $\langle \mathbf{k}, \mathbf{S}, \mathbf{A} \rangle$.

Definition 19 (Consequence)

A query \mathbf{g} is a consequence of a TIGER database $\text{DB} = \langle \mathbf{k}, \mathbf{S}, \mathbf{A} \rangle$ if for all models \mathcal{M} it holds that

$$\text{if } \mathcal{M} \models_{\mathbf{A}, \mathbf{S}} \mathbf{k} \text{ then } \mathcal{M} \models \mathbf{g}.$$

As before, the scope of a variable does not exceed a graph description, i.e. a corpus graph or a query \mathbf{g} .

Chapter 4

TIGER calculus

As a step towards interpreter for the TIGER language, we present a set of inference rules, the TIGER calculus. The calculus is a straightforward adaptation of the resolution calculus for Horn clauses in logic programming languages. Each graph in a corpus definition corresponds to a fact in the *database*, say with predicate name `graph` whose argument is the graph under consideration. Similarly, a query can be thought of as being enclosed by a `graph`-predicate name (template calls remain separate). The database includes also the template defining clauses. Queries (and template subgoals) are handled by an adaptation of the resolution rule. Instead of first order terms, we plug in graph descriptions (with embedded feature constraints). Accordingly, term unification is replaced by constraint solving. The TIGER calculus constitutes a naive interpreter for the TIGER language with the goal of keeping the correctness proofs with respect to the formal semantics as simple as possible. For an efficient implementation, certain portions of the calculus have to be redesigned. But then, it is sufficient to give correctness proofs with respect to the original TIGER calculus, not with respect to the formal semantics of the TIGER language.

Subsequently, we assume that there is a given TIGER database, i.e. a corpus definition \mathbf{k} , a set of template definitions \mathbf{S} , and a set of type definitions and feature declarations \mathbf{A} which contains at least the type definition

$$\top := \text{Constant}, \text{FRec}, \text{Node}, \text{Graph}; \quad (4.1)$$

Instead of a variable substitution σ , now a *constraint store* Σ and *store* Θ of *variable renamings* are used to keep track of variable substitutions. A constraint store Σ is a set of pairs $\langle \#x, \phi \rangle$, where the 'constraint' ϕ is either a feature value formula \mathbf{d} , a feature constraint \mathbf{k} , or a pair $\langle \mathbf{c}, \#x \rangle$ of a constant node identifier \mathbf{c} (or the \top -symbol) and a 'feature constraint variable' $\#x$ (i.e. $\{ \sigma(\#x) \} \cap \mathcal{F} \neq \emptyset$). We assume that, initially, the constraint store Σ contains a pair $\langle \#x, \top \rangle$ resp. $\langle \#x, \langle \top, \#x' \rangle \rangle$ ($\#x'$ fresh feature constraint variable) for every possible variable $\#x$. Or, more practically, in the case that a membership test for a variable $\#x$ in Σ would fail, one assumes the existence of the pair $\langle \#x, \top \rangle$ resp. $\langle \#x, \langle \top, \#x' \rangle \rangle$ in Σ ($\#x'$ fresh feature constraint variable).

The store Θ of variable renamings $\langle \#x_1, \#x_2 \rangle$ initially is an empty set. We assume that

there exists a function $\text{deref}(\#x) \rightarrow \#x$ which is defined as follows:

$$\text{deref}(\#x_1) := \begin{cases} \#x_1 & \text{if } \langle \#x_1, \#x_2 \rangle \notin \Theta \\ \text{deref}(\#x_2) & \text{if } \langle \#x_1, \#x_2 \rangle \in \Theta \end{cases}$$

In order to illustrate how complex the treatment of variables in the scope of negation is, we add reduction rules for these cases, as well (but we refrain from the implementation of this part of the calculus.) For variables in the scope of negation, one more data structure is required: a store Δ of 'dif' constraints¹ $\langle \#x_1, \#x_2 \rangle$. Initially, the dif store is empty. A dif constraint $\langle \#x_1, \#x_2 \rangle$ (on the level of feature values) is interpreted as

$$\llbracket \#x_1 \rrbracket \cap (C \setminus \llbracket \#x_2 \rrbracket)$$

We chose the above interpretation from a variety of ways how two variables can be considered different:

- different variable names

Two variables denote two (physically) non-identical structures.

- inequality

The denotations of two variables are not equivalent, i.e.

$$!((\#x_1 \mid \#x_2) \& (!\#x_2 \mid \#x_1)).$$

- disjointness of denotations

$$!(\#x_1 \& \#x_2)$$

- the denotation of variable $\#x_2$ is a negative constraint on variable $\#x_1$:

$$\#x_1 \& !\#x_2$$

The last alternative fits the formal semantics of variables which occur in the scope of negation, e.g. in the formula $!(f=\#x_2)$, where $\#x_1$ stands for $!(\#x_2)$. Note that the substitution $\sigma(\#x_2)$ ignores whether $\#x_2$ occurs in a positive or in a negative context.

Subsequently, the tests for type inclusion and type membership refer to the semantic notions directly. The definition of a rule system for type inclusion etc. seems somewhat superfluous vis-a-vis the standard computer implementation of these tests. When implementing strict type hierarchies, a type is usually represented by a bit vector which contains a positive bit for each constant member of the type, i.e. such a bit vector is a direct implementation of the denotation of a given type.

¹This name has been adopted from [10] although the interpretation may be different from the one in [10].

4.1 Feature Values

4.1.1 Disjunctive Normal Form

Double negation elimination

$$!(\neg d) \implies d$$

de Morgan rules

$$\neg(d_1 \ \& \ d_2) \implies (\neg d_1) \mid (\neg d_2)$$

$$\neg(d_1 \mid d_2) \implies (\neg d_1) \ \& \ (\neg d_2)$$

Distribution rule

$$(d_1 \mid d_2) \ \& \ d_3 \implies (d_1 \ \& \ d_3) \mid (d_2 \ \& \ d_3)$$

4.1.2 Disjunction reduction

$$c \mid c \implies c$$

$$c \mid (\neg c) \implies \top$$

$$(\neg c) \mid (\neg c) \implies (\neg c)$$

$$c \mid t \implies t \quad \text{if } c \in \llbracket t \rrbracket$$

$$c \mid (\neg t) \implies \begin{cases} \top & \text{if } \llbracket t \rrbracket = \{c\} \\ (\neg t) & \text{if } c \notin \llbracket t \rrbracket \end{cases}$$

$$(\neg c) \mid t \implies \begin{cases} \top & \text{if } \llbracket t \rrbracket = \{c\} \\ (\neg c) & \text{if } c \notin \llbracket t \rrbracket \end{cases}$$

$$(\neg c) \mid (\neg t) \implies \begin{cases} (\neg c) & \text{if } c \in \llbracket t \rrbracket \\ \top & \text{if } c \notin \llbracket t \rrbracket \end{cases}$$

$$t_1 \mid t_2 \implies \begin{cases} t_2 & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \\ t_1 & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \end{cases}$$

$$t_1 \mid (\neg t_2) \implies \begin{cases} \top & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \\ (\neg t_2) & \text{if } \llbracket t_1 \rrbracket \cap \llbracket t_2 \rrbracket = \emptyset \end{cases}$$

$$(\neg t_1) \mid (\neg t_2) \implies \begin{cases} (\neg t_1) & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \\ (\neg t_2) & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \\ \top & \text{if } \llbracket t_1 \rrbracket \cap \llbracket t_2 \rrbracket = \emptyset \end{cases}$$

4.1.3 Consistency check

Constant-constant merge

$$\begin{aligned}
c_1 \ \& \ c_2 &\Longrightarrow \begin{cases} c_1 & \text{if } c_1 = c_2 \\ \perp & \text{otherwise} \end{cases} \\
c_1 \ \& \ (!c_2) &\Longrightarrow \begin{cases} c_1 & \text{if } c_1 \neq c_2 \\ \perp & \text{otherwise} \end{cases} \\
(!c) \ \& \ (!c) &\Longrightarrow (!c)
\end{aligned}$$

Type check

$$\begin{aligned}
c \ \& \ t &\Longrightarrow \begin{cases} c & \text{if } c \in \llbracket t \rrbracket \\ \perp & \text{otherwise} \end{cases} \\
c \ \& \ (!t) &\Longrightarrow \begin{cases} c & \text{if } c \notin \llbracket t \rrbracket \\ \perp & \text{otherwise} \end{cases} \\
(!c) \ \& \ t &\Longrightarrow t \quad \text{if } c \notin \llbracket t \rrbracket \\
(!c) \ \& \ (!t) &\Longrightarrow (!t) \quad \text{if } c \in \llbracket t \rrbracket
\end{aligned}$$

Type-type merge

$$\begin{aligned}
t_1 \ \& \ t_2 &\Longrightarrow \begin{cases} t_1 & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \\ t_2 & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \\ \perp & \text{otherwise} \end{cases} \\
t_1 \ \& \ (!t_2) &\Longrightarrow \begin{cases} t_1 & \text{if } \llbracket t_2 \rrbracket \cap \llbracket t_1 \rrbracket = \emptyset \\ \perp & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \end{cases} \\
(!t)_1 \ \& \ (!t_2) &\Longrightarrow \begin{cases} (!t_2) & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \\ (!t_1) & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \end{cases}
\end{aligned}$$

Lemma 1 (Completeness)

If $\llbracket d \rrbracket = \emptyset$ then $d \xRightarrow{*} \perp$

Proof: (sketch)

Assume that there exists a case where $\llbracket d \rrbracket = \emptyset$ but not $d \xRightarrow{*} \perp$.

However, all syntactic combinations (of Boolean operators and types of subformulas) have been covered by the reduction rules. Furthermore, whenever a lefthand side of a reduction rule is inconsistent, its lefthand side is an inconsistent formula, as well. In particular, a clash \perp will be propagated in the cases $(c \ \& \ \perp)$, $(c \ \& \ (!\top))$, $((!c) \ \& \ \perp)$, $((!c) \ \& \ (!\top))$, $(\perp \ \& \ t)$, $(t \ \& \ (!\top))$, $((!\top) \ \& \ (!t))$, $(\perp \ | \ \perp)$, $(\perp \ | \ (!\top))$, $((!\top) \ | \ (!\top))$ by the corresponding constant-constant merge rule, type check rule, or disjunction reduction rule. Hence, we get a contradiction to the assumption.

Lemma 2 (Soundness)

If $d \xRightarrow{*} \perp$ then $\llbracket d \rrbracket = \emptyset$

Proof: (sketch)

Assume that $d \xRightarrow{*} \perp$ but not $\llbracket d \rrbracket = \emptyset$.

However, in every reduction rule, the denotation of the lefthand side equals the denotation of the righthand side, therefore this assumption cannot hold. For example, the treatment of $(t_1 \ \& \ (!t_2))$ is sound, since types are organized in a hierarchy. There are three cases:

1. $\llbracket t_2 \rrbracket \cap \llbracket t_1 \rrbracket = \emptyset$:
 $(C \setminus \llbracket t_2 \rrbracket) \cap \llbracket t_1 \rrbracket = C \cap \llbracket t_1 \rrbracket = \llbracket t_1 \rrbracket$
2. $\llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket$:
 $(C \setminus \llbracket t_2 \rrbracket) \cap \llbracket t_1 \rrbracket = (C \setminus \llbracket t_2 \rrbracket) \cap \llbracket t_2 \rrbracket = \emptyset$
3. $\llbracket t_2 \rrbracket \subset \llbracket t_1 \rrbracket$ cannot be reduced to a more concise form

Lemma 3

$\llbracket d \rrbracket = \emptyset \iff d \xRightarrow{*} \perp$

4.2 Feature Constraints

The function $\text{evalre}(z)$ takes a regular expression z and returns a disjunction of constants $(c_1 \mid \dots \mid c_n)$. As an 'internal' representation, we admit the formula $!\#x$.

4.2.1 Disjunctive Normal Form

Feature Value Disjunction

$$\begin{aligned} f = (d_1 \mid d_2) &\implies (f = d_1) \mid (f = d_2) \\ f = \#x:(d_1 \mid d_2) &\implies (f = \#x:d_1) \mid (f = \#x:d_2) \end{aligned}$$

Negated feature reduction (t_1 domain of f)

$$\begin{aligned} !(f = d) &\implies (!t_1) \mid (f = !d) \\ !(f = \#x_2) &\implies (!t_1) \mid (f = !\#x_2) \\ !(f = \#x:d) &\implies !((f = \#x) \ \& \ (f = d)) \\ !(f \neq d) &\implies !(f = !d) \\ !(f = /z/) &\implies !(f = \text{evalre}(z)) \\ !(f \neq /z/) &\implies !(f = !\text{evalre}(z)) \end{aligned}$$

Double negation elimination

$$!(!k) \implies k$$

de Morgan rules

$$!(k_1 \ \& \ k_2) \implies (!k_1) \ | \ (!k_2)$$

$$!(k_1 \ | \ k_2) \implies (!k_1) \ \& \ (!k_2)$$

Distribution rule

$$(k_1 \ | \ k_2) \ \& \ k_3 \implies (k_1 \ \& \ k_3) \ | \ (k_2 \ \& \ k_3)$$

4.2.2 Disjunction reduction

$$\perp \ | \ k \implies k$$

$$t_1 \ | \ t_2 \implies \begin{cases} t_2 & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \\ t_1 & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \end{cases}$$

$$t_1 \ | \ (!t_2) \implies \begin{cases} \top & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \\ (!t_2) & \text{if } \llbracket t_1 \rrbracket \cap \llbracket t_2 \rrbracket = \emptyset \end{cases}$$

$$(!t_1) \ | \ (!t_2) \implies \begin{cases} (!t_1) & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \\ (!t_2) & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \\ \top & \text{if } \llbracket t_1 \rrbracket \cap \llbracket t_2 \rrbracket = \emptyset \end{cases}$$

The subsequent rules apply to the individual disjuncts of the Disjunctive Normal Form. For the evaluation of each disjunct, there are fresh stores Σ , Θ , Δ .

4.2.3 Feature normalization

$$f = !\#x_2 \implies f = \#x_1$$

$$\Delta := \Delta \cup \{\langle \#x_1, \#x_2 \rangle\}, \#x_1 \text{ fresh variable}$$

$$f = d \implies f = \#x \quad \Sigma := \Sigma \cup \{\langle \#x, d \rangle\}, \#x \text{ fresh variable}$$

$$f = \#x \implies t_1 \ \& \ f = \#x$$

$$\Sigma := (\Sigma \setminus \{\langle \text{deref}(\#x), d \rangle\}) \cup \{\langle \text{deref}(\#x), (\text{Constant} \ \& \ t_2 \ \& \ d) \rangle\}$$

$$t_1 \text{ domain, } t_2 \text{ range of } f$$

$$f = \#x:d_1 \implies f = \#x$$

$$\Sigma := (\Sigma \setminus \{\langle \text{deref}(\#x), d_2 \rangle\}) \cup \{\langle \text{deref}(\#x), (d_1 \ \& \ d_2) \rangle\}$$

$$f != d \implies f = !d$$

$$f = /z/ \implies f = \text{evalre}(z)$$

$$f != /z/ \implies f = !\text{evalre}(z)$$

4.2.4 Consistency check

Constraint store

$$k \implies \perp \quad \text{if } \langle \#x, \perp \rangle \in \Sigma$$

Dif store

$$k \implies \perp \\ \langle \#x_1, \#x_2 \rangle \in \Delta; \langle \text{deref}(\#x_1), d_1 \rangle, \langle \text{deref}(\#x_2), d_2 \rangle \in \Sigma \\ \text{if } (d_1 \ \& \ !d_2) \implies \perp$$

Failure propagation

$$\perp \ \& \ k \implies \perp$$

Type-type merge

$$t_1 \ \& \ t_2 \implies \begin{cases} t_1 & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \\ t_2 & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \\ \perp & \text{otherwise} \end{cases}$$

$$t_1 \ \& \ (!t_2) \implies \begin{cases} t_1 & \text{if } \llbracket t_2 \rrbracket \cap \llbracket t_1 \rrbracket = \emptyset \\ \perp & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \end{cases}$$

$$(!t)_1 \ \& \ (!t)_2 \implies \begin{cases} (!t)_2 & \text{if } \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \\ (!t)_1 & \text{if } \llbracket t_2 \rrbracket \subseteq \llbracket t_1 \rrbracket \end{cases}$$

Feature-feature merge

$$(f = \#x_1) \ \& \ (f = \#x_2) \implies f = \#x_1 \\ \Sigma := (\Sigma \setminus \{ \langle \text{deref}(\#x_1), d_1 \rangle, \langle \text{deref}(\#x_2), d_2 \rangle \}) \\ \cup \{ \langle \text{deref}(\#x_1), (d_1 \ \& \ d_2) \rangle \} \\ \Theta := \Theta \cup \{ \langle \text{deref}(\#x_2), \text{deref}(\#x_1) \rangle \}$$

For the subsequent proofs, let k' be a single disjunct of the disjunctive normal form, i.e. a disjunction-free formula.

Lemma 4 (*Completeness, single disjunct*)

$$\text{If } \llbracket k' \rrbracket = \emptyset \quad \text{then} \quad \langle k', \Sigma, \Theta, \Delta \rangle \xRightarrow{*} \langle \perp, \Sigma', \Theta', \Delta' \rangle$$

Proof: (sketch)

Similar argument as in the proof sketch of Lemma 1. Note that the reduction rules also reduce the formulas in the constraint store Σ .

In addition, one has to check that the constraint store and the store of variable renamings are handled in such a manner that all inconsistencies will be discovered:

A chain of variable renamings denotes an equivalence class of variables. The final element (the `deref`-value) of such a chain is the representative of the equivalence class.

There is at most one constraint per variable in the constraint store:

For example, the rule for feature-feature merge erases the constraints for the 'dereferenced variables' $\#x_1$ and $\#x_2$ from the constraint store and adds the accumulated constraints for $\#x_1$. The fact that $\#x_1$ and $\#x_2$ belong to the same equivalence class of variables is realized by the appropriate extension of the store of variable renamings.

Furthermore, the termination property must be checked:

It is sufficient to apply the type check for features only once to each feature in the feature constraint.

It is sufficient to carry out an inconsistency check of the dif store each time the Σ - or Θ -information on one of the variables contained in the dif store has been extended. The inconsistency check of the dif store will not create new dif constraints, since the updated Σ -constraint ($d_1 \& !d_2$) does not contain any variables.

Lemma 5 (*Soundness, single disjunct*)

If $\langle k', \Sigma, \Theta, \Delta \rangle \xRightarrow{*} \langle \perp, \Sigma', \Theta', \Delta' \rangle$ then $\llbracket k' \rrbracket = \emptyset$

Proof: (sketch)

Similar argument as in the proof sketch of Lemma 2. Check for example the equivalence of the lefthand side and righthand side of the first rule for negated feature reduction, cf. [16].

$$\begin{aligned}
& \llbracket !(f = d) \rrbracket \\
&= \mathcal{F} \setminus \llbracket f = d \rrbracket \\
&= \mathcal{F} \setminus (\llbracket t_1 \rrbracket \cap \llbracket f = d \rrbracket) \\
&= (\mathcal{F} \setminus \llbracket t_1 \rrbracket) \cup (\mathcal{F} \setminus \llbracket f = d \rrbracket) \\
&= (\mathcal{F} \setminus \llbracket t_1 \rrbracket) \cup ((\mathcal{F} \setminus \llbracket t_1 \rrbracket) \setminus \llbracket f = d \rrbracket) \\
&= \llbracket !t_1 \rrbracket \cup \llbracket f = !d \rrbracket
\end{aligned}$$

The rule for $!(f = \#x_2)$ works together with the check of the dif store, which guarantees that the denotation of the value of feature f will be a subset of $\mathcal{F} \setminus \llbracket \#x_2 \rrbracket$.

Lemma 6

$\llbracket d \rrbracket = \emptyset$ iff $d \xRightarrow{*} \perp$

Proof: (sketch)

Completeness:

Show that all feature constraints will be reduced to disjunctive normal form (DNF). By Lemma 4, the calculus is complete wrt. an individual disjunct of the DNF. The rules of type disjunction and disjunction reduction will reduce an inconsistent DNF to \perp .

Soundness:

Show that for each DNF rule, the lefthand side is equivalent to its righthand side. For example:

$$\begin{aligned} \llbracket f = (d_1 \mid d_2) \rrbracket &= \llbracket f = \#x : (d_1 \mid d_2) \rrbracket = \llbracket (f = \#x : d_1) \mid (f = \#x : d_2) \rrbracket \\ &= \llbracket f = \#x : d_1 \rrbracket \cup \llbracket f = \#x : d_2 \rrbracket = \llbracket f = d_1 \rrbracket \cup \llbracket f = d_2 \rrbracket = \llbracket (f = d_1) \mid (f = d_2) \rrbracket \end{aligned}$$

4.3 Node Descriptions

Subsequently, we assume that constant node identifiers occur only in corpus graphs.

4.3.1 Disjunctive Normal Form

Node description normalization

$$\begin{aligned} [k] &\implies [\#x : k] && \#x \text{ fresh variable} \\ [\#x_2 : k] &\implies \#x_1 : [\#x_2 : k] && \#x_1 \text{ fresh variable} \\ \#x_1 : [k] &\implies \#x_1 : [\#x_2 : k] && \#x_2 \text{ fresh variable} \end{aligned}$$

Feature constraint disjunction

$$\#x_1 : [\#x_2 : (k_1 \mid k_2)] \implies \#x_1 : [\#x_2 : k_1] \mid \#x_1 : [\#x_2 : k_2]$$

The remaining rules for node descriptions apply to the individual disjuncts of the Disjunctive Normal Form. We assume fresh stores Σ , Θ , Δ for each disjunct.

4.3.2 Node description normalization

$$\begin{aligned} [] &\implies \#x && \#x \text{ fresh variable} \\ c : [k] &\implies \#x_1 && \#x_1, \#x_2 \text{ fresh variables} \\ &&& \Sigma := \Sigma \cup \{\langle \#x_1, \langle c, \#x_2 \rangle \rangle, \langle \#x_2, k \rangle\} \\ \#x : [] &\implies \#x \\ \#x_1 : [\#x_2] &\implies \#x_1 && \langle \text{deref}(\#x_1), \langle c_3, \#x_3 \rangle \rangle \in \Sigma \\ &&& \Sigma := (\Sigma \setminus \{\langle \text{deref}(\#x_2), k_2 \rangle, \langle \#x_3, k_3 \rangle\}) \cup \{\langle \#x_3, (k_2 \ \& \ k_3) \rangle\} \\ &&& \Theta := \Theta \cup \{\langle \text{deref}(\#x_2), \#x_3 \rangle\} \\ \#x_1 : [\#x_2 : k_2] &\implies \#x_1 : [\#x_2] \\ &&& \Sigma := (\Sigma \setminus \{\langle \text{deref}(\#x_2), k_3 \rangle\}) \cup \{\langle \text{deref}(\#x_2), (k_2 \ \& \ k_3) \rangle\} \end{aligned}$$

4.3.3 Node guess

Here, we need the failure symbol \perp as a kind of 'node identifier' - in the case that c_1 and c_2 are not the same constants.

$$\begin{aligned}
\#x_1 &\Longrightarrow \#x_2 \\
&\text{for some } \#x_2 \text{ in the corpus graph} \\
\Sigma &:= (\Sigma \setminus \{ \langle \text{deref}(\#x_1), \langle c_1, \#x_3 \rangle \rangle, \langle \#x_2, \langle c_2, \#x_4 \rangle \rangle, \\
&\quad \langle \#x_3, k_3 \rangle, \langle \#x_4, k_4 \rangle \}) \\
&\quad \cup \{ \langle \#x_2, \langle c_1 \ \& \ c_2, \#x_4 \rangle \rangle, \langle \#x_4, k_3 \ \& \ k_4 \rangle \} \\
\Theta &:= \Theta \cup \{ \langle \text{deref}(\#x_1), \#x_2 \rangle, \langle \text{deref}(\#x_3), \#x_4 \rangle \}
\end{aligned}$$

4.4 Graph Descriptions

Subsequently, we admit the failure symbol \perp as a graph description.

4.4.1 Disjunctive Normal Form

Distribution rule

$$(\#g_1 \mid \#g_2) \ \& \ \#g_3 \quad \Longrightarrow \quad (\#g_1 \ \& \ \#g_3) \mid (\#g_2 \ \& \ \#g_3)$$

4.4.2 Disjunction reduction

$$\perp \mid \#g \quad \Longrightarrow \quad \#g$$

4.4.3 Graph description normalization

Elimination of trivial node relations

$$\#x \ \& \ \#g \quad \Longrightarrow \quad \#g$$

Node relation normalization

$$\begin{aligned}
\#x_1 \ > \ \#x_2 &\Longrightarrow \ \#x_1 \ >1 \ \#x_2 \\
\#x_1 \ . \ \#x_2 &\Longrightarrow \ \#x_1 \ .1 \ \#x_2 \\
\#x_1 \ >1 \ \#x_2 &\Longrightarrow \ (\#x_1 \ >1 \ \#x_2) \ \& \ (\#x_1 \ >1 \ \#x_2) \\
\#x_1 \ \$ \ . * \ \#x_2 &\Longrightarrow \ (\#x_1 \ \$ \ \#x_2) \ \& \ (\#x_1 \ . * \ \#x_2)
\end{aligned}$$

Corpus graph transitive closure

$$\begin{aligned}
& (\#x_1 >a \#x_2) \ \& \ (\#x_2 >b \#x_3) \quad \Longrightarrow \\
& \quad (\#x_1 >a \#x_2) \ \& \ (\#x_2 >b \#x_3) \ \& \ (\#x_1 >(a+b) \#x_3) \\
& (\#x_1 .a \#x_2) \ \& \ (\#x_2 .b \#x_3) \quad \Longrightarrow \\
& \quad (\#x_1 .a \#x_2) \ \& \ (\#x_2 .b \#x_3) \ \& \ (\#x_1 .(a+b) \#x_3) \\
\#x \quad & \Longrightarrow \quad \text{root}(\#x) \\
& \quad \text{if } \#x \text{ is the root node of a corpus graph} \\
\#x \quad & \Longrightarrow \quad \text{arity}(\#x, b) \\
& \quad \text{where } b \text{ is the number of nodes directly dominated by } \#x \\
\#x \quad & \Longrightarrow \quad \text{tokenarity}(\#x, b) \\
& \quad \text{where } b \text{ is the number of terminal leaves dominated by } \#x \\
\#x \quad & \Longrightarrow \quad \text{continuous}(\#x) \\
& \quad \text{if } \#x \text{ is the root of a continuous subgraph} \\
\#x \quad & \Longrightarrow \quad \text{discontinuous}(\#x) \\
& \quad \text{if } \#x \text{ is the root of a discontinuous subgraph}
\end{aligned}$$

The following rules apply to the individual disjuncts of the disjunctive normal form (of the query), with fresh stores Σ , Θ , Δ for each disjunct.

4.4.4 Consistency check

Constraint store

$$g \quad \Longrightarrow \quad \perp \quad \text{if } \langle \#x, \perp \rangle \text{ or } \langle \#x_1, \langle \perp, \#x_2 \rangle \rangle \in \Sigma$$

Failure propagation

$$\perp \ \& \ g \quad \Longrightarrow \quad \perp$$

Root

$$\begin{aligned}
\text{root}(\#x) \quad & \Longrightarrow \quad \#x \\
& \quad \text{if } \text{root}(\#x) \text{ in the corpus graph}
\end{aligned}$$

Arity

$$\begin{aligned}
\text{arity}(\#x, a) \quad & \Longrightarrow \quad \begin{cases} \#x & \text{if } a = b' \\ \perp & \text{if } a \neq b' \end{cases} \\
& \quad \text{if } \text{arity}(\#x, b') \text{ in the corpus graph} \\
\text{arity}(\#x, a, b) \quad & \Longrightarrow \quad \begin{cases} \#x & \text{if } a \leq b' \leq b \\ \perp & \text{if } b' < a \text{ or } b' > b \end{cases} \\
& \quad \text{if } \text{arity}(\#x, b') \text{ in the corpus graph}
\end{aligned}$$

Token arity

$$\begin{aligned} \text{tokenarity}(\#x, a) &\implies \begin{cases} \#x & \text{if } a = b' \\ \perp & \text{if } a \neq b' \end{cases} \\ &\text{if } \text{tokenarity}(\#x, b') \text{ in the corpus graph} \\ \text{tokenarity}(\#x, a, b) &\implies \begin{cases} \#x & \text{if } a \leq b' \leq b \\ \perp & \text{if } b' < a \text{ or } b' > b \end{cases} \\ &\text{if } \text{tokenarity}(\#x, b') \text{ in the corpus graph} \end{aligned}$$

Continuous

$$\begin{aligned} \text{continuous}(\#x) &\implies \#x \\ &\text{if } \text{continuous}(\#x) \text{ in the corpus graph} \end{aligned}$$

discontinuous

$$\begin{aligned} \text{discontinuous}(\#x) &\implies \#x \\ &\text{if } \text{discontinuous}(\#x) \text{ in the corpus graph} \end{aligned}$$

Node relations

$$\begin{aligned} \#x > a \#x &\implies \perp \\ \#x_1 > l \#x_2 &\implies \#x_2 \\ &\text{if } (\#x_1 > l \#x_2) \text{ in the corpus graph} \\ \#x_1 > * \#x_2 &\implies \#x_1 \\ &\text{if } (\#x_1 > a \#x_2) \text{ in the corpus graph} \\ \#x_1 > a \#x_2 &\implies \#x_1 \\ &\text{if } (\#x_1 > a \#x_2) \text{ in the corpus graph} \\ \#x_1 > a, b \#x_2 &\implies \#x_1 \\ &\text{if } (\#x_1 > a' \#x_2) \text{ in the corpus graph, } a \leq a' \leq b \\ \#x_1 > @l \#x_2 &\implies \#x_1 \\ &\text{if } \#x_2 \text{ is left corner of } \#x_1 \\ \#x_1 > @r \#x_2 &\implies \#x_1 \\ &\text{if } \#x_2 \text{ is right corner of } \#x_1 \\ \#x . \#x &\implies \perp \\ \#x_1 . * \#x_2 &\implies \#x_1 \\ &\text{if } (\#x_1 . a \#x_2) \text{ in the corpus graph} \\ \#x_1 . a \#x_2 &\implies \#x_1 \\ &\text{if } (\#x_1 . a \#x_2) \text{ in the corpus graph} \\ \#x_1 . a, b \#x_2 &\implies \#x_1 \\ &\text{if } (\#x_1 . a' \#x_2) \text{ in the corpus graph, } a \leq a' \leq b \\ \#x_1 \$ \#x_2 &\implies \#x_1 \end{aligned}$$

if $(\#x_0 > \#x_1), (\#x_0 > \#x_2)$ in the corpus graph, $\#x_1 \neq \#x_2$
 $\#x_1 !>1 \#x_2 \implies \#x_2$
 if $(\#x_1 >1 \#x_2)$ not in the corpus graph
 $\#x_1 !> \#x_2 \implies \#x_2$
 if $\forall l (\#x_1 >l \#x_2)$ not in the corpus graph
 $\#x_1 !>* \#x_2 \implies \#x_1$
 if $\forall a (\#x_1 >a \#x_2)$ not in the corpus graph
 $\#x_1 !>a \#x_2 \implies \#x_1$
 if $(\#x_1 >a \#x_2)$ not in the corpus graph
 $\#x_1 !>a,b \#x_2 \implies \#x_1$
 if $(\#x_1 >a' \#x_2)$ in the corpus graph, $a' < a$ or $b < a'$
 $\#x_1 !>@l \#x_2 \implies \#x_1$
 if $\#x_2$ is not left corner of $\#x_1$
 $\#x_1 !>@r \#x_2 \implies \#x_1$
 if $\#x_2$ is not right corner of $\#x_1$
 $\#x_1 !\$ \#x_2 \implies \#x_1$
 if $\exists \#x_0$ s.t. $(\#x_0 >1 \#x_1)$ or $(\#x_0 >1 \#x_2)$ not in the corpus graph
 $\#x_1 !. \#x_2 \implies \#x_1$ if $(\#x_1 . \#x_2)$ not in the corpus graph
 $\#x_1 !.* \#x_2 \implies \#x_1$
 if $\forall a (\#x_1 .a \#x_2)$ not in the corpus graph
 $\#x_1 !.a \#x_2 \implies \#x_1$
 if $(\#x_1 .a \#x_2)$ not in the corpus graph
 $\#x_1 !.a,b \#x_2 \implies \#x_1$
 if $(\#x_1 .a' \#x_2)$ in the corpus graph, $a' < a$ or $b < a'$

4.4.5 Templates

For the reduction of template calls, we assume the usual Prolog style deduction rule, i.e. the template call $r_1(u_1)$ is replaced by the body g of a (copy of a) template clause whose head $r_2(u_2)$ matches the template call in terms of number of arguments. Of course, the variables u_2 in the template clause have to be replaced by the corresponding variables u_1 in the template goal.

$$r_1(u_1) \implies g[u_2/u_1]$$

if there is a template clause with head $r_2(u_2)$
 and the number of arguments u_1 equals the number of arguments of u_2

4.4.6 Queries

A query q can be derived from a corpus, if the query can be derived from some graph in the corpus, i.e. if the query is compatible with some corpus graph.

$k, g \vdash q$

if $g \vdash q$

$g \vdash q$

if not $g \& q \xRightarrow{*} \perp$

Bibliography

- [1] Patrick Blackburn, Claire Gardent, and Wilfried Meyer-Viol. Talking about trees. In *Proceedings of the 6th Conference of the European Chapter of the Association for Computational Linguistics*, pages 21–29, Utrecht, 1993.
- [2] Oliver Christ. A modular and flexible architecture for an integrated corpus query system. In *Proceedings of COMPLEX'94: 3rd Conference on Computational Lexicography and Text Research (Budapest, July 7–10 1994)*, pages 23–32, Budapest, Hungary, 1994. CMP-LG archive id 9408005.
- [3] Oliver Christ, Bruno M. Schulze, Anja Hofmann, and Esther König. *Corpus Query Processor (CQP). User's Manual*. Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart, Stuttgart, Germany, 1999. www.ims.uni-stuttgart.de/CorpusWorkbench/.
- [4] Jochen Dörre. *Feature-Logik und Semiunifikation*. Dissertationen zur Künstlichen Intelligenz. infix-Verlag, Sankt Augustin, Germany, 1996.
- [5] Jochen Dörre and Michael Dorna. CUF - a formalism for linguistic knowledge representation. Deliverable R.1.2A, DYANA 2, August 1993.
- [6] Jochen Dörre, Dov M. Gabbay, and Esther König. Fibred semantics for feature-based grammar logic. *Journal of Logic, Language, and Information. Special Issue on Language and Proof Theory*, 5:387–422, 1996.
- [7] Denys Duchier and Joachim Niehren. Solving dominance constraints with finite set constraint programming. Technical report, Universität des Saarlandes, Programming Systems Lab, 1999.
- [8] Martin Emele. *Der TFS-Repräsentationsformalismus*. PhD thesis, Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart, Stuttgart, Germany, 1997. AIMS Arbeitspapiere des IMS, vol.3(6).
- [9] Martin Emele and Rémi Zajac. A fixed-point semantics for feature type systems. In *Proceedings of the 2nd International Workshop on Conditional and Typed Rewriting Systems*, Montreal, Canada, 1990.
- [10] F. Giannesini, H. Kanoui, R. Pasero, and M. van Caneghem. *Prolog*. InterÉdition, Paris, 1985.
- [11] Markus Höhfeld and Gert Smolka. Definite relations over constraint languages. LILOG-Report 53, IBM Deutschland, Stuttgart, Baden-Württemberg, October 1988.

- [12] Esther König, Wolfgang Lezius, and Holger Voormann. *TIGERSearch User's Manual*. IMS, University of Stuttgart, Stuttgart, 2003.
- [13] Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 1993.
- [14] James Rogers and K. Vijay-Shanker. Reasoning with descriptions of trees. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, Newark, Delaware, 1992.
- [15] Helmut Schmid. *YAP: Parsing and Disambiguation With Feature-Based Grammar*. PhD thesis, Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart, Stuttgart, Germany, 1999.
- [16] Gert Smolka. A feature logic with subsorts. Technical Report 33, IBM Deutschland, Institute for Knowledge-Based Systems, Stuttgart, Baden-Württemberg, 1988.