

On Semistructured Data and Linguistics

Frank Morawietz

(joint work with Uwe Mönnich and Stephan Kepser)

<http://tcl.sfs.uni-tuebingen.de/~frank/>

Seminar für Sprachwissenschaft
Theoretische Computerlinguistik
Universität Tübingen

15. April 2002



Overview

- Introduction and Motivation
- Semistructured Data
- Query Languages
- Intermediate Conclusion
- Secondary Edges
- Regular Queries of Context-Sensitive Relations
- MSO Queries
- Lifting
- Reconstruction
- Conclusion and Outlook





Introduction and Motivation

The Two Cultures

- exchange of information
- The Web
 - decomposition into units: files
 - structure (internal and between documents)
- Databases
 - relational database schemata or entity relationship diagrams
 - logical or abstract vs. physical implementation
 - concerns
 - query languages
 - mechanisms for concurrency control and for recovery
 - efficient implementation



Why do we have to bridge the gap?

- TV station A publishes data about political parties
 - stored: relational database
 - output is generated on demand in HTML via the formatting of the results of an SQL query
- Party B wants to analyze A's data
 - access only to A's HTML pages
 - write software for the analysis (parsing of HTML)
- Problems:
 - brittle
 - need to download many pages to get at the desired information



What we have ...

● The Web

- global infrastructure and standards for document exchange
- presentation of hypertexts
- user interfaces for document retrieval (IR techniques)
- XML (data with structure)

● Databases

- storage techniques and query languages with efficient methods for large amounts of (rigidly) structured data
- data models and methods for structuring data
- mechanisms for maintaining integrity and consistency
- new model: semistructured data (relaxation of the demands on the structuring)



Basic Architecture

- Databases: client/server architecture
 - client issues a query
 - processing, optimization and execution
 - answer is returned by the server
- Web: “multitier” architecture
 - lowest level: data sources (database or file servers, legacy systems, “anything that produces data”) in a common logical data model and a common format (XML)
 - highest level: user interfaces or analysis software
 - middleware: transforms, integrates, adds data (data warehousing, mediator systems)
 - queries: translation, processing, translation



Data Abstraction

- databases: three levels
 - physical
 - logical
 - external
 - web
 - no distinction between physical and logical layer
 - is there a reason why we should make this distinction?
 - Example: database can be a single XML file, or a compressed version, or a relational database, or ...
- ⇒ XML as both physical and logical representation



Data Diversity

- Heterogeneous data
 - How is the data conceptually represented?
 - Example: addresses – sorted by name? sorted by town?
 - Reconciliation at the logical level is possible
- ⇒ XML for the logical level



Research Context

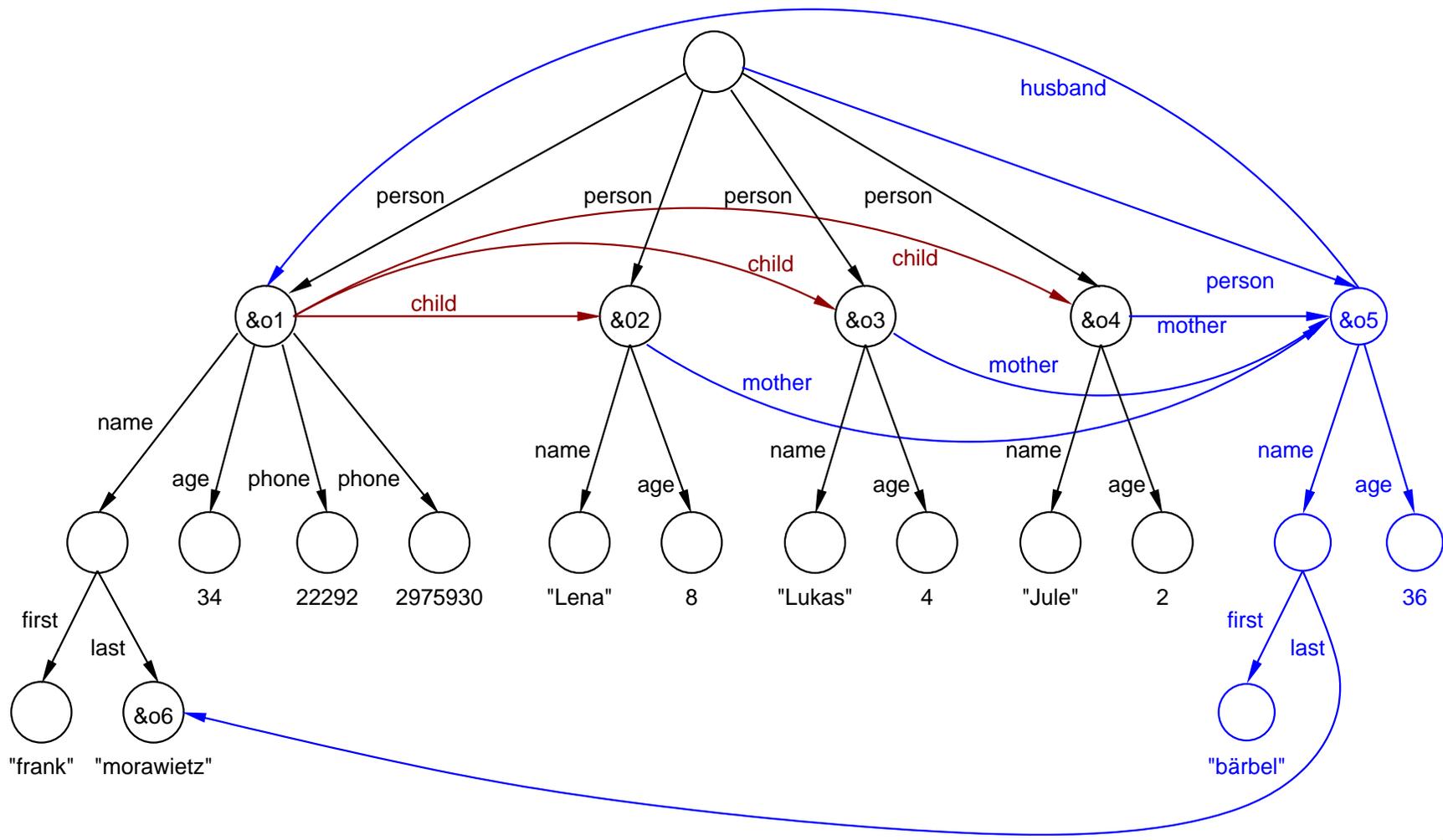
- **Semistructured Data**
 - incompleteness and irregularity of data
 - self-describing
 - integration of several kinds of data
 - no relational schemata, but structured
- **XML**
 - universal format for data exchange
 - data model: labeled graphs/trees
 - query languages





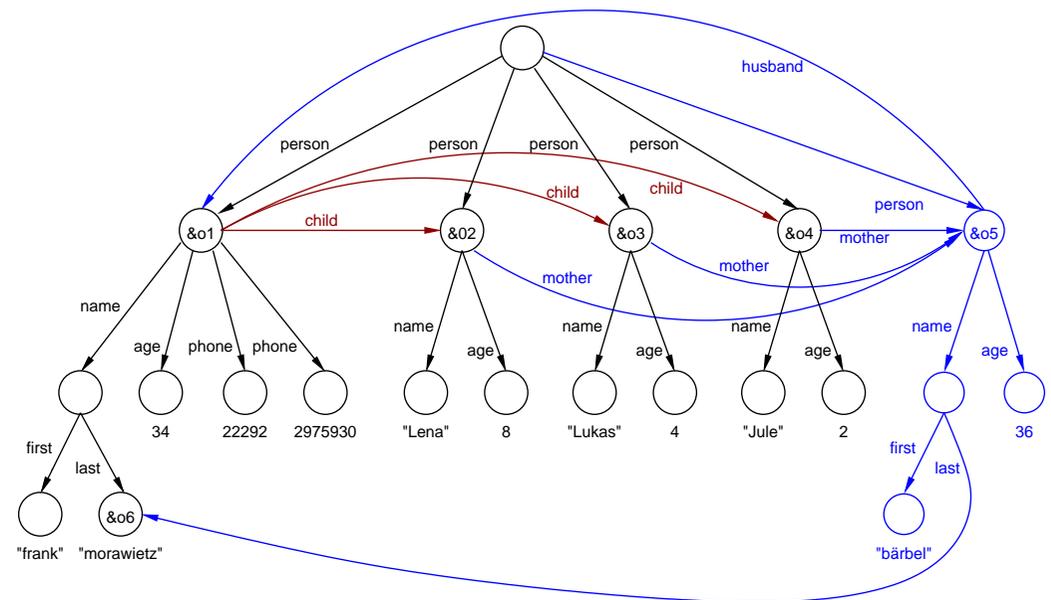
Semistructured Data

The Data Model: Syntax



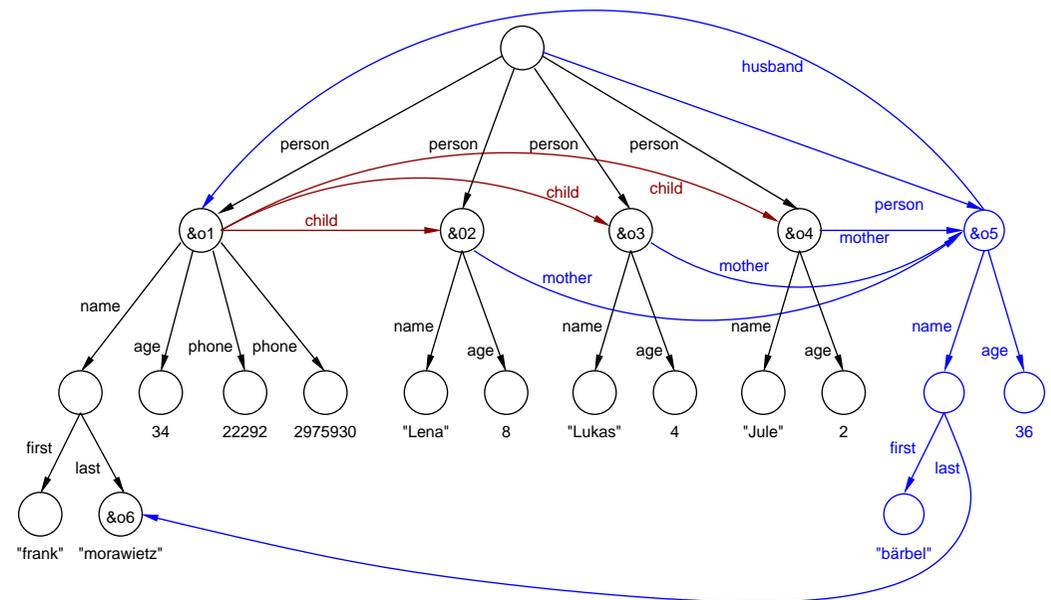
The Data Model: Syntax

{name: "Lena", age: 8}



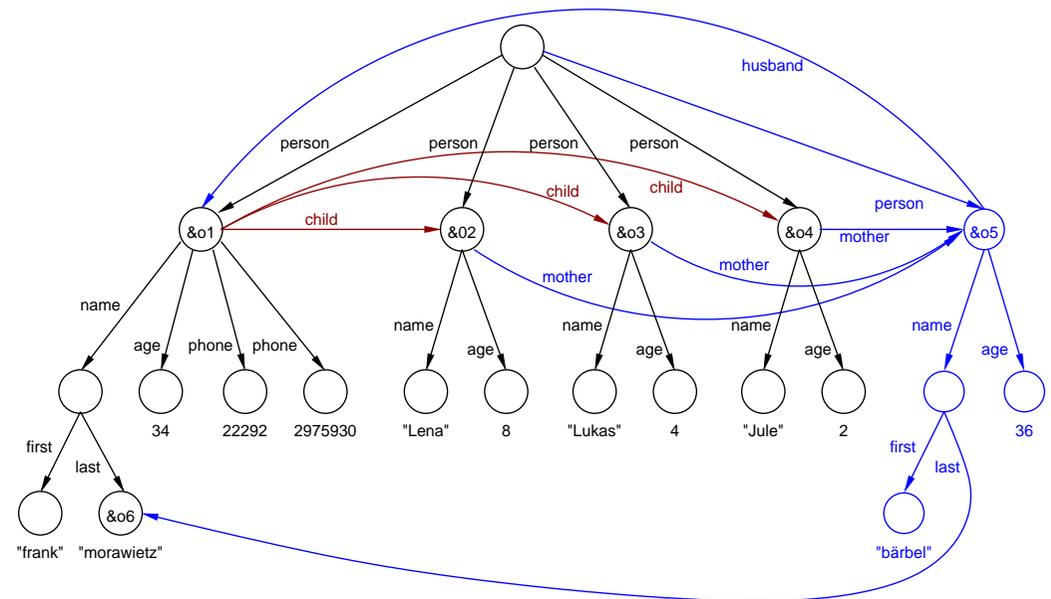
The Data Model: Syntax

```
{name: {first: "frank", last: "morawietz"},  
age: 34,  
phone: 22292}
```



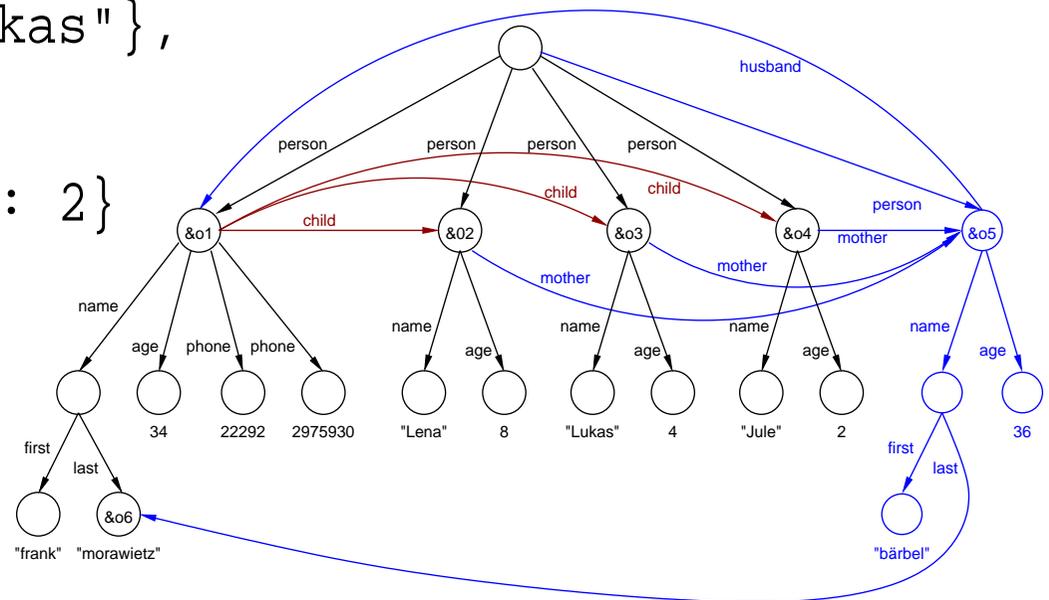
The Data Model: Syntax

```
{name: {first: "frank", last: "morawietz"},  
age: 34,  
phone: 22292,  
phone: 2975930}
```



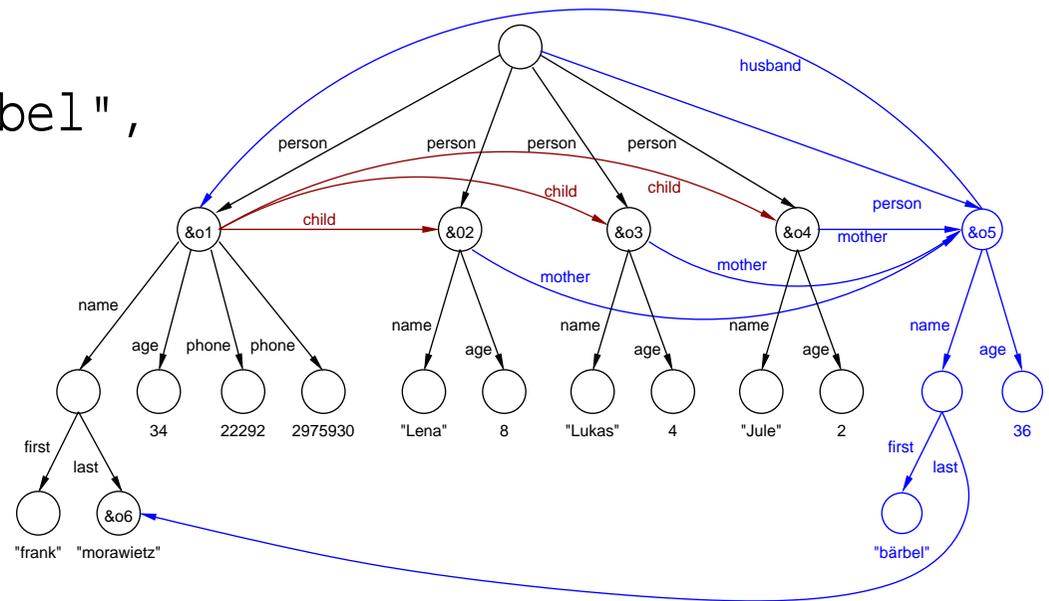
The Data Model: Syntax

```
{  
  person: {name: "Lena", age: 8},  
  
  person: {age: 4, name: "Lukas"},  
  
  person: {name: "Jule", age: 2}  
}
```



The Data Model: Syntax

```
{  
person:  
  &o2{name: "Lena", age: 8, mother: &o5 },  
person:  
  &o5{name: {first: "Bärbel",  
            last: &o6},  
       age: 36,  
       husband: &o1  
},
```



...

}



The Data Model: Syntax

- Convention: implicit OIDs
- $\{a: \{b: 3\}, a: \{b: 3\}\}$
- coding of relational databases (e.g., via rows)
- edge labeled graph

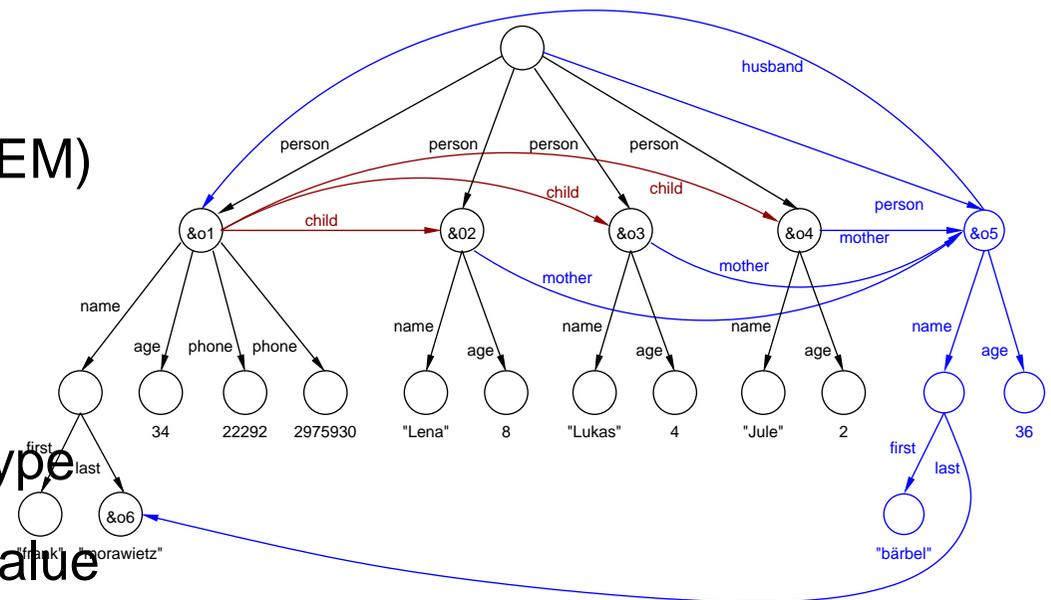
⇒ Object Exchange Model (OEM)

label: string

oid: object identifier

type: complex or a specific type

value: list of OIDs or atomic value



Semistructured Data and XML

```
{ person: {name: {first: "frank", last: "morawietz"},  
           age: 34,  
           phone: 22292}}
```

```
<person>  
  <name>  
    <first> frank </first>  
    <last> morawietz </last>  
  </name>  
  <age> 34 </age>  
  <phone> 22292 </phone>  
</person>
```



Semistructured Data in Linguistics

Linguistic Databases: collections of linguistic data

- relational databases (e.g., typological databases)
- lexicons
- corpora (e.g., books, newspaper articles, transcriptions, ...)

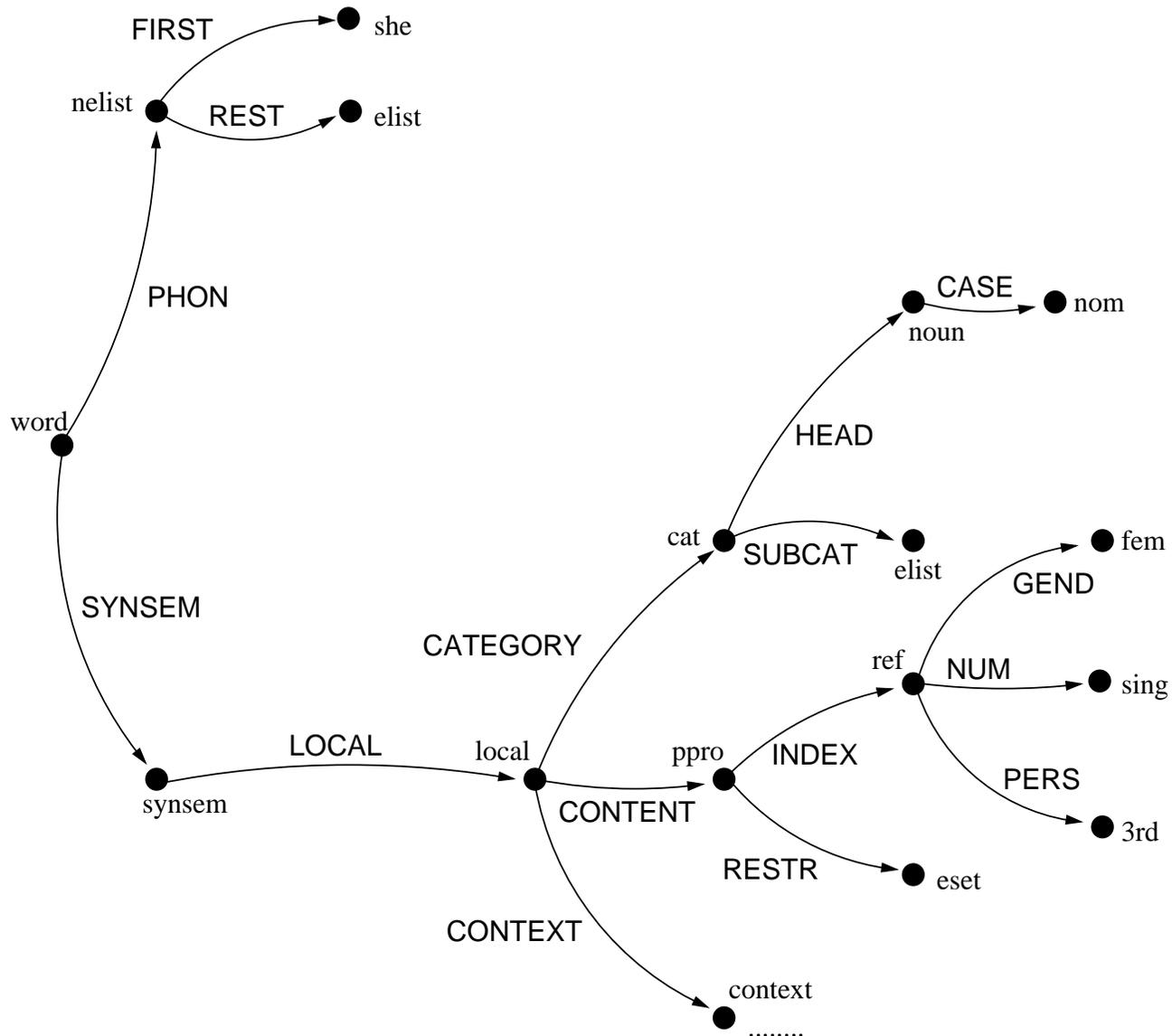


Semistructured Data in Linguistics

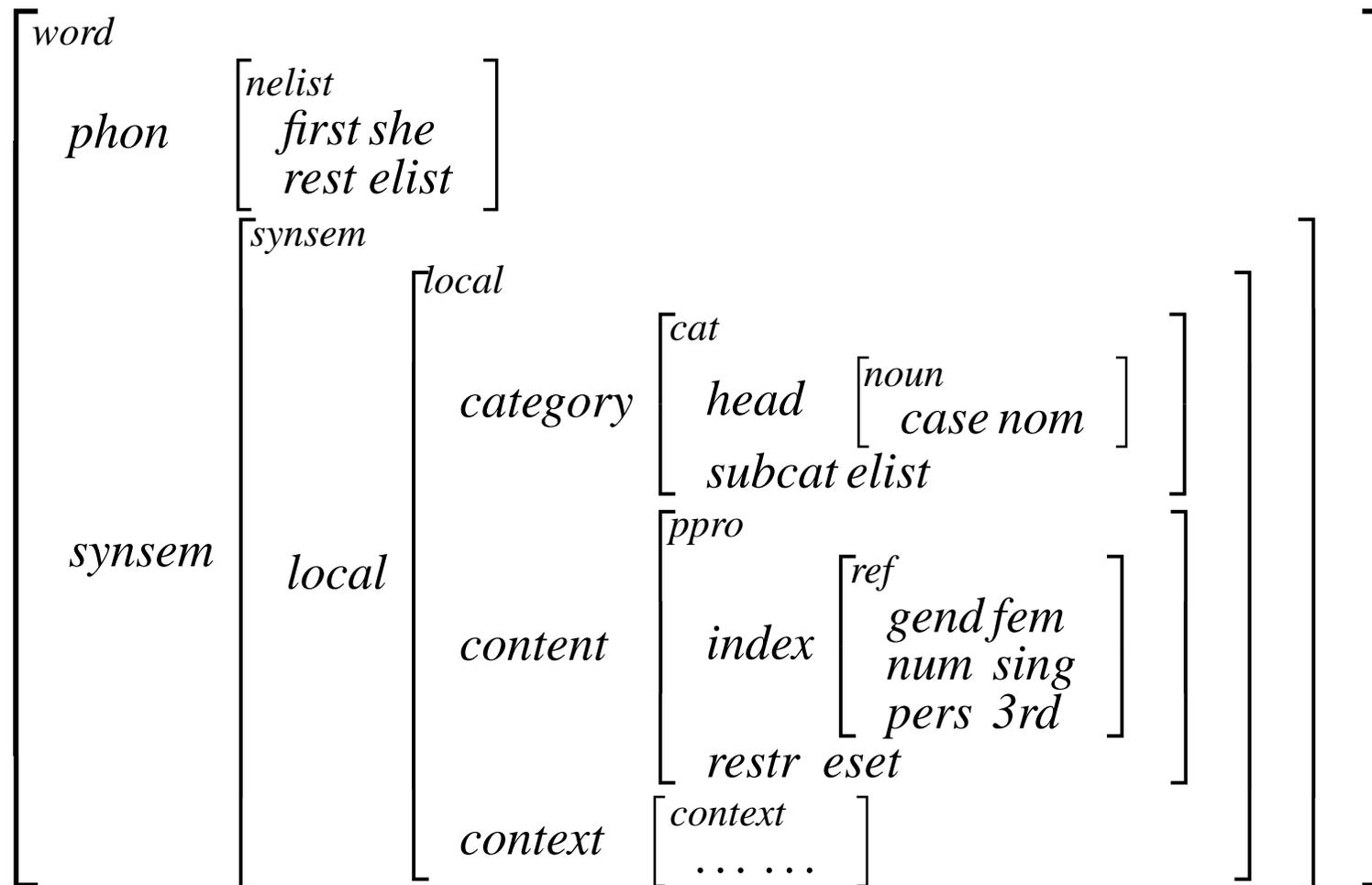
Language	WO	Prep	Case	Adj	...
German	SOV	←	Acc	←	...
English	SVO	←	Acc	←	...
Warao	SVO?	→	Acc	→	...



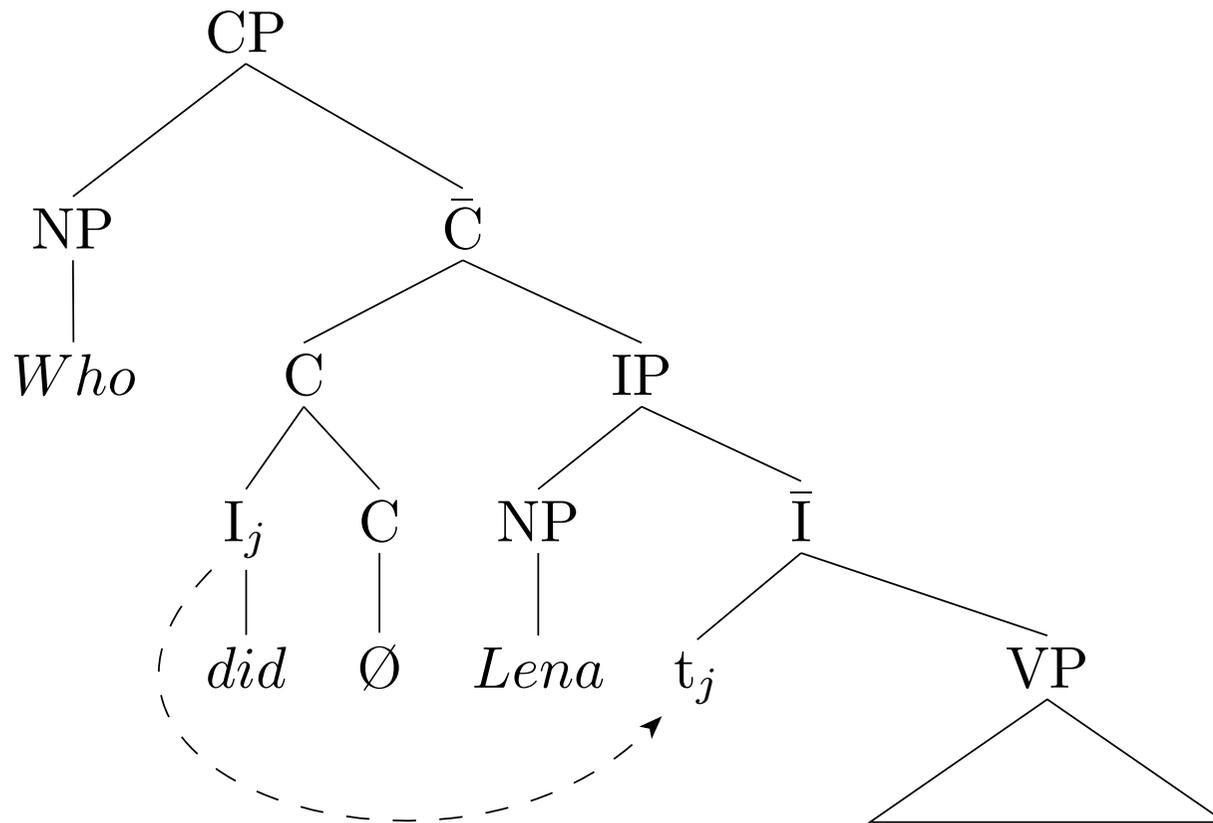
Semistructured Data in Linguistics



Semistructured Data in Linguistics



Semistructured Data in Linguistics



Who did Lena invite to her birthday party?



Semistructured Data in Linguistics

Penn Treebank II:

(*TOP* (*S* (*NP-SBJ* my best friend)
(*VP* gave
(*NP* me)
(*NP* chocolate)
(*NP-TMP* yesterday))
.))



Annotation

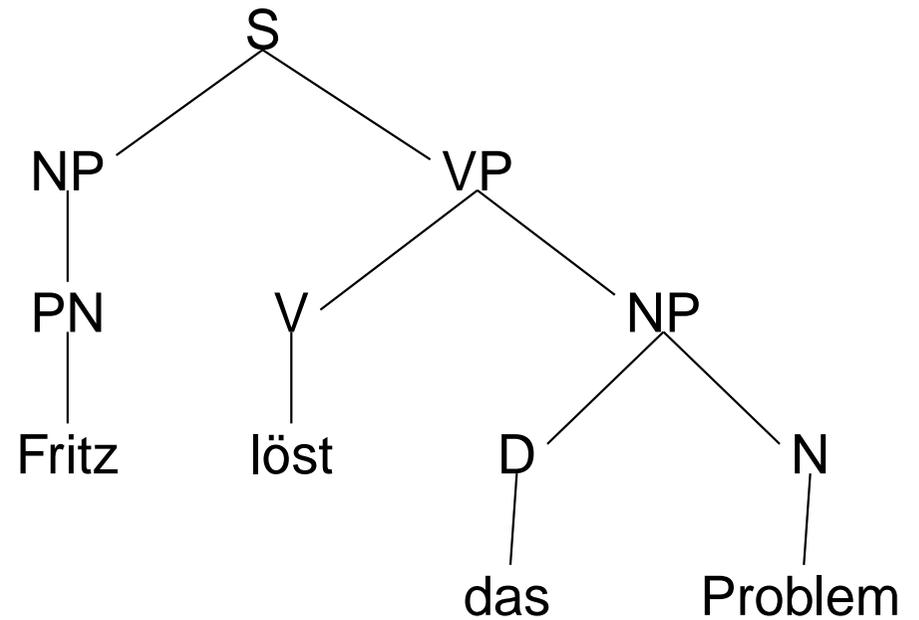
Additional information to the “raw” data

- meta information
- document structure
- part-of-speech tags
- syntactical analysis
- specific syntactic or semantic phenomena
- ...



Core XML

```
<S>  
  <NP>  
    <PN> Fritz </PN>  
  </NP>  
  <VP>  
    <V> löst </V>  
    <NP>  
      <D> das </D>  
      <N> Problem </N>  
    </NP>  
  </VP>  
</S>
```





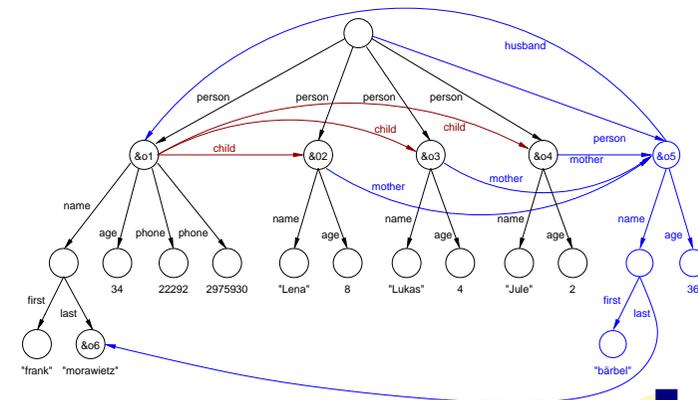
Query Languages

A simple Query Language

- path expressions ("person.name.first" yields the set of nodes with labels "frank" and "bärbel")
- path expressions result in sets of nodes
- paths via properties: regular expressions (including a wild card)
 - on the alphabet of edge labels ("person._*.name")
 - on the alphabet of the characters of edge labels (((n|N)ame)(s)?)

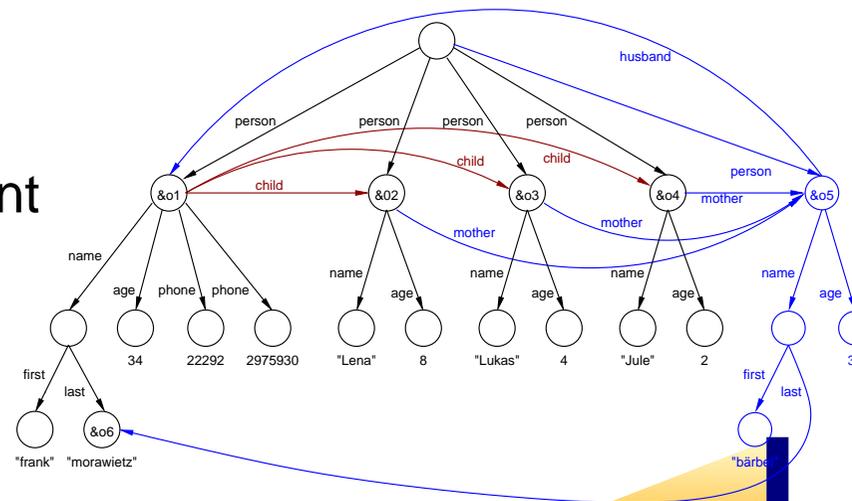
● Query syntax based on Lorel/UnQL

- select firstname: X
from person.name.(first)? X
- select firstname: X
from person.name X
where "frank" in X.first



More Complex Queries

- select person: X
from person X
X.age Y
where Y < 30
- select person: {firstname: Y, age: Z}
from person X
X.name.(first)? Y
X.age Z
- nested queries in the select statement
- natural joins



Syntactic Sugar for Queries

- default labels
- path expressions in the `select` clause
- equality predicates
- coercion
- relation names as variables
- label and path variables



Desiderata for Query Languages

- Expressive Power
 - operations (comparable to, e.g., SQL)
 - restructuring
 - completeness
- Semantics
 - query transformation
 - query optimization
- Compositionality: output can be used as input of a query
 - same data model
 - referential transparency
- Structure Consciousness
- Program Manipulation



Reintroducing Structure: Why?

- Ignoring knowledge of types loses information
- Regularities facilitates querying.
- Protection against incorrect updates.
- Simplify access and improve performance.
- Optimization of the physical representation.
- Mixing querying structured and semistructured data.
- Maybe hidden from the user?



Query Languages for XML

- combine XML syntax with query languages presented previously
- XML-QL, XQuery, XSL(T)
- Problems:
 - Order
 - infinite domain of attributes



Advanced Topics

- first-order interpretation
- reconstruction
 - Datalog
 - Graphical Languages: GraphLog, Complete Answer Aggregates
 - Structural Recursion
 - Bisimulation
- Strudel/StruQL



TGrep (Penn Treebank)

NP	anything called just NP
TOP	anything called TOP. Since every sentence is enclosed in TOP brackets, this matches everything.
VP < NP	VP immediately dominates an NP
VP <3 NP	VP has NP as its 3rd child
NP > VP	an NP that has a VP parent
VP << dog	VP contains the word dog (VP dominates dog, directly or indirectly)
S \!<< dog	S does not contain the word dog
NP \$.. NP	two NPs have a common parent (first NP precedes a sister also called NP)
/^NP/	// encloses a grep-style regular expression



TGrep (Penn Treebank)

```
S <1 / ^NP/ < (VP < (NP $.. NP))
```

Get all Ss that start with an NP (not necessarily the subject) and that dominate a VP that in turn has two NP children – in other words, sentences with what might be double-object VPs.



TGrep2

- Rather than simply having a set of required relationships and a set of prohibited relationships, nodes can have full boolean expressions of relationships to other nodes.
- Nodes can be given unique labels and may then be referred to by those labels in the pattern specification or in selecting trees for printing.
- Patterns are no longer restricted to simple tree architectures. The use of node labels and segmented patterns allows links in a pattern to form back-edges as well, permitting cycles of links.
- Customizable output formats allow a variety of information to be reported in a flexible manner.
- Multiple search patterns may be specified and one can retrieve the first subtree matching any pattern, the first subtree matching each pattern, or all subtrees matching all patterns.
- Subtrees can be reported using a code rather than by printing the whole structure. The trees themselves can later be retrieved using the codes.
- A variety of new links have been added and the immediately-precedes link now has a more conventional meaning.
- TGrep2 corpus files are substantially smaller than tgrep corpora.



FXGrep (Neumann/Seidl)

- XML input as a document forest representing the logical structure
 - A text node represents one or more adjacent pieces of character data.
 - A PI node represents a processing instruction and is annotated with the target name of the processing instruction. It has a single child: a text node giving the text of the processing instruction.
 - An element node represents an element and is annotated with the element type and the attributes specified for the element. The children of an element node are a sequence of nodes constructed from the content of the element.
- query implemented by consecutive runs of two forest automata (for simple queries a single run suffices)



FXGrep: Pattern Language

<code>/doc/section/title</code>	paths
<code>/doc//title</code>	descendant
<code>//title</code>	descendant anywhere
<code>//(section/ subsection/)title</code>	or
<code>//title/"airplane"</code>	content
<code>//title[toc='yes']</code>	attributes
<code>//section[_ subsection _]/title</code>	titles of sections with subtitles
<code>a[_ b#c _]/d</code>	context qualifier

Alternatively with grammar patterns.



FXGrep: Pattern Language

```
//section[!_ subsection _][_(//footnote/'automata')_]/title
```

The titles of all sections that do not have a child of type `subsection` and that do have a descendant of type `footnote` which contains the text `automata`.



FXGrep: Theory

- words are monadic trees
- ranked trees are terms with internal nodes being operators and leaves being constants
- unranked trees represent document structure, nodes are elements
- unranked trees can be represented by ranked trees, but not uniquely
- trees can be specified with grammars, regular expressions, formulas in fixpoint or monadic second-order logic, automata

specification	automaton
grammar	polynomial
regular expressions	polynomial
fixpoint formulas	exponential
MSO formulas	non-elementary

- Nice animated slides by Helmut Seidl (Uni Trier)



Tree Automata

- ranked case
 - work bottom-up from the frontier to the root
 - transitions according to the states of the daughters and the label of the mother
 - a tree is accepted if it halts in a final state at the root
- unranked case
 - additional word automaton to recognize the sequence of labels of the daughters
- variations
 - extension to “hedges” (Murata, Brüggemann-Klein, Wood)
 - tree-walking automata navigate through the tree
 - caterpillars (Brüggemann-Klein, Herrmann, Wood)
 - pushdown tree automata (Neumann/Seidl)



Query Automata (Neven/Schwentick)

- two-way tree automata + a selection function
 - partition state-alphabet symbol pairs in U and D
 - always works on a cut of the tree
 - down transition: a node is replaced by its children
 - up transition: all children of a node are replaced by the node
 - ⇒ new states for all nodes in the cut
- query: the set of nodes selected by the two-way tree automaton via the selection function
- Expressivity: unary MSO patterns
- Extension to the unranked case possible, but complicated





Intermediate Conclusion

Fundamental Problem

Scylla (Expressivity)

VS.

Charybdis (Undecidability)



Fundamental Problem II

classical

semistructured

RA

“Automata”

FO

MSO

SQL

Tree Description language
with regular path descriptions





Secondary Edges

Cross-serial dependencies

*(... weil) der Karl die Maria dem Peter den Hans
schwimmen lehren helfen liess*

German: Dijck-Sprache—CF

*(... omdat) Karel Marie Piet Jan liet helpen leren
zwemmen*

Dutch: $a^n b^n$ —CF

(... wil) mer de maa em chind lönd hälffe schwüme



Cross-serial dependencies

*(... weil) der Karl die Maria dem Peter den Hans
schwimmen lehren helfen liess*

German: Dijck-Sprache—CF

*(... omdat) Karel Marie Piet Jan liet helpen leren
zwemmen*

Dutch: $a^n b^n$ —CF

(... wil) mer de maa em chind lönd hälffe schwüme



Cross-serial dependencies

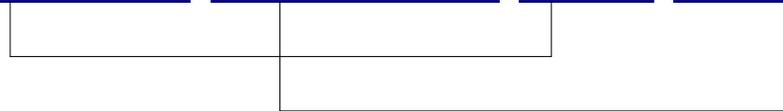
*(... weil) der Karl die Maria dem Peter den Hans
schwimmen lehren helfen liess*

German: Dijck-Sprache—CF

*(... omdat) Karel Marie Piet Jan liet helpen leren
zwemmen*

Dutch: $a^n b^n$ —CF

(... wil) mer de maa em chind lönd hälffe schwüme



Cross-serial dependencies

*(... weil) der Karl die Maria dem Peter den Hans
schwimmen lehren helfen liess*

German: Dijck-Sprache—CF

*(... omdat) Karel Marie Piet Jan liet helpen leren
zwemmen*

Dutch: $a^n b^n$ —CF

(... wil) mer de maa em chind lönd hälffe schwüme



Swiss-German: $a^n b^m c^n d^m$ —Non-CF



Regular Query of context-sensitive relations



Overview: Two-Step Approach

- Assumption: there exists a description of the context-sensitive phenomena (e.g., via a context-free tree grammar)
- database query
 - MSO query into the treebank
 - ⇒ set of possible answers/trees
 - lifting of the grammar and the possible answer trees
 - ⇒ tree automaton representing the grammar and a set of lifted candidate trees
 - the automaton filters the lifted candidate trees
 - ⇒ set of lifted answer trees
- reconstruction of the intended trees via MSO-transduction or macro tree transducer
 - ⇒ set of intended answer trees



Context-Free Tree Grammars

A context-free tree grammar Γ is a 5-tuple $\langle \Sigma, \mathbf{F}, S, \mathbf{X}, \mathbf{P} \rangle$, with

- Σ, \mathbf{F} ranked alphabets of *inoperatives* and *operatives*;
- $S \in \mathbf{F}$ start symbol;
- \mathbf{X} a set of variables; and
- \mathbf{P} a set of productions.



Context-Free Tree Grammars

The productions $p \in \mathbf{P}$ are of the form

$$F(x_1, \dots, x_n) \longrightarrow t$$

with $F \in \mathbf{F}$, $x_1, \dots, x_n \in \mathbf{X}$, and t is a term over $\Sigma \cup \mathbf{F} \cup \{x_1, \dots, x_n\}$.

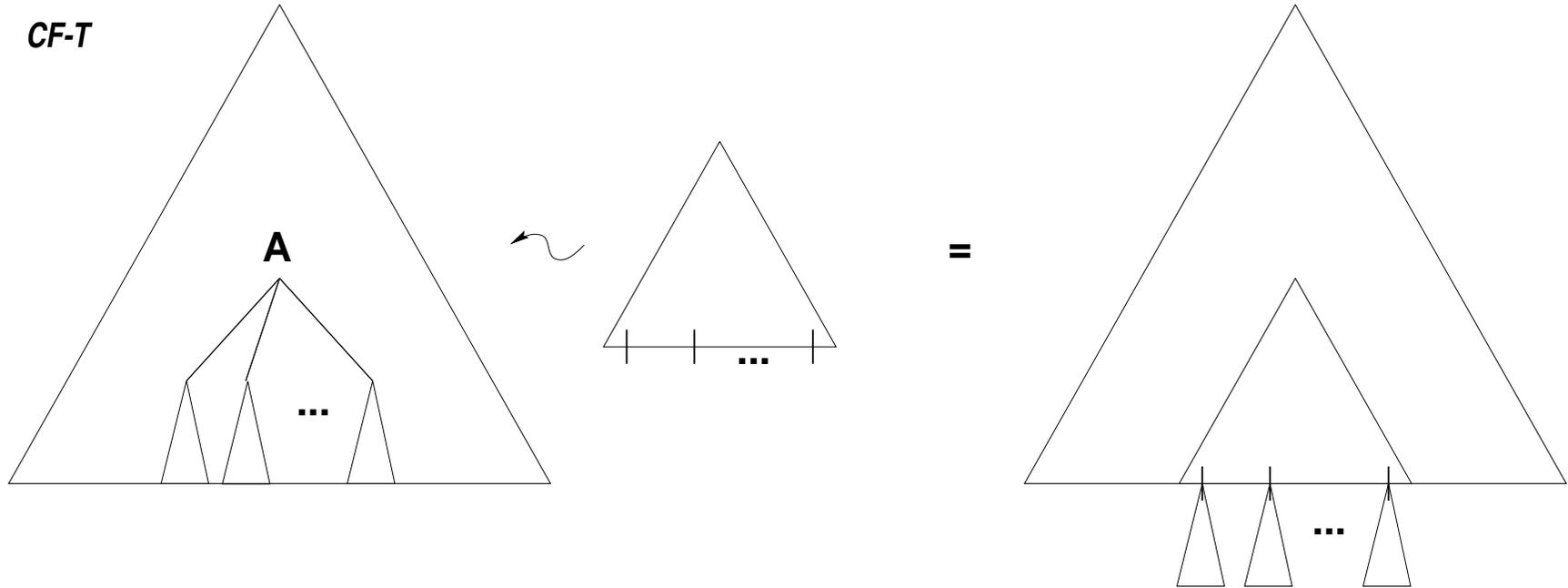
An application of a rule $F(x_1, \dots, x_n) \longrightarrow t$ “rewrites” a subtree rooted in F as the tree t with its respective variables substituted by F 's daughters.

Tree grammars with $\mathbf{F}_n = \emptyset, n \neq 0$, are called *regular*.

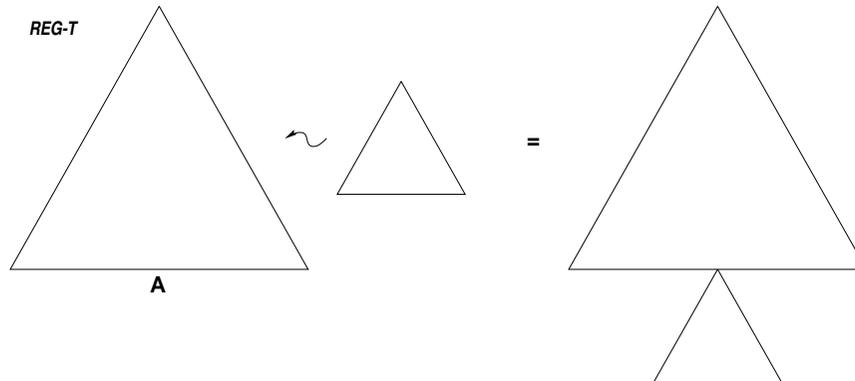


Context-Free Tree Grammars

CF-T



REG-T



A CFTG for $a^n b^m c^n d^m$: \mathcal{G}

$$\Sigma_0 = \{\varepsilon, a, b, c, d\}$$

$$\Sigma_2 = \{\bullet\}$$

$$\mathbf{F}_0 = \{S\}$$

$$\mathbf{F}_4 = \{F\}$$

$$S = S$$

$$\mathbf{X} = \{x_1, x_2, x_3, x_4\}$$

$$\mathbf{P} = \left\{ \begin{array}{l} S \longrightarrow \varepsilon \\ S \longrightarrow F(a, \varepsilon, c, \varepsilon) \\ S \longrightarrow F(\varepsilon, b, \varepsilon, d) \\ F(x_1, x_2, x_3, x_4) \longrightarrow F(\bullet(a, x_1), x_2, \bullet(c, x_3), x_4) \\ F(x_1, x_2, x_3, x_4) \longrightarrow F(x_1, \bullet(b, x_2), x_3, \bullet(d, x_4)) \\ F(x_1, x_2, x_3, x_4) \longrightarrow \bullet(\bullet(\bullet(x_1, x_2), x_3), x_4) \end{array} \right.$$

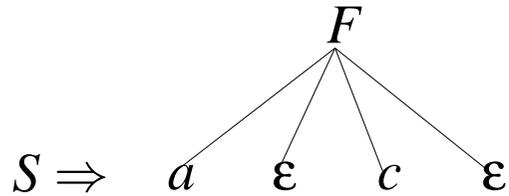


An example derivation of *aabccd*

S



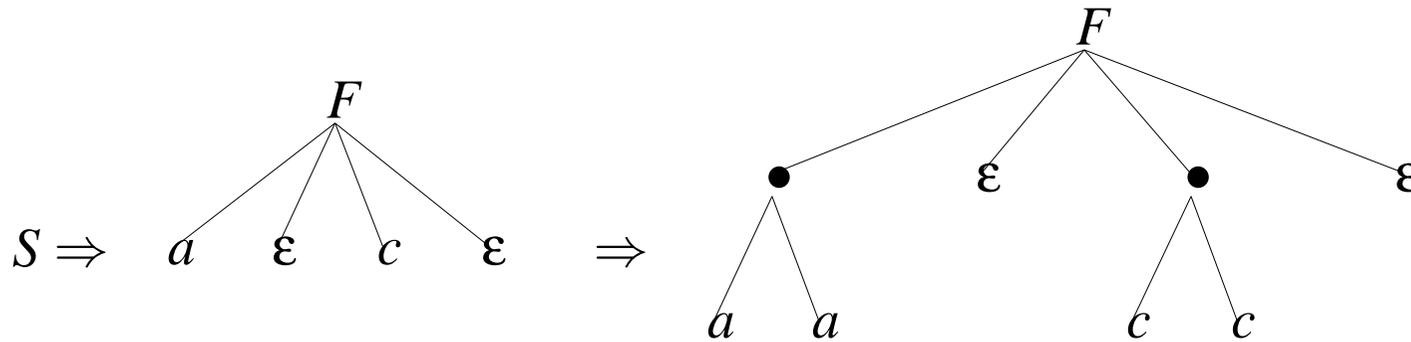
An example derivation of $aabcccd$



$$S \longrightarrow F(a, \varepsilon, c, \varepsilon)$$



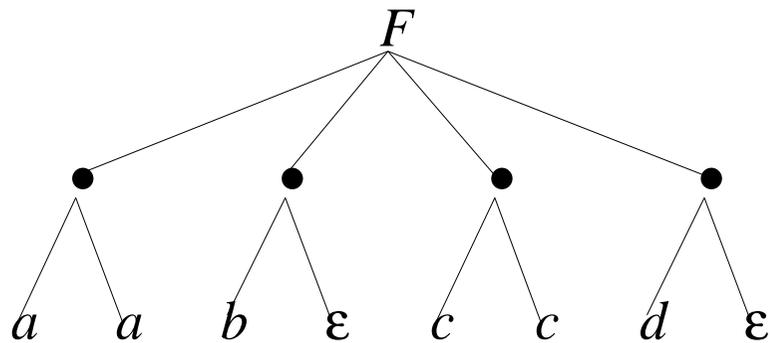
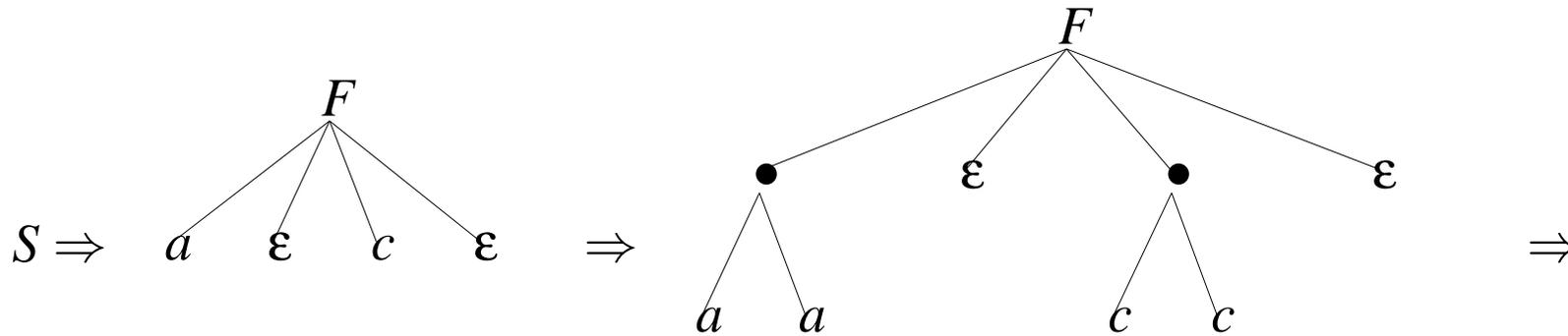
An example derivation of $aabccd$



$$F(x_1, x_2, x_3, x_4) \longrightarrow F(\bullet(a, x_1), x_2, \bullet(c, x_3), x_4)$$



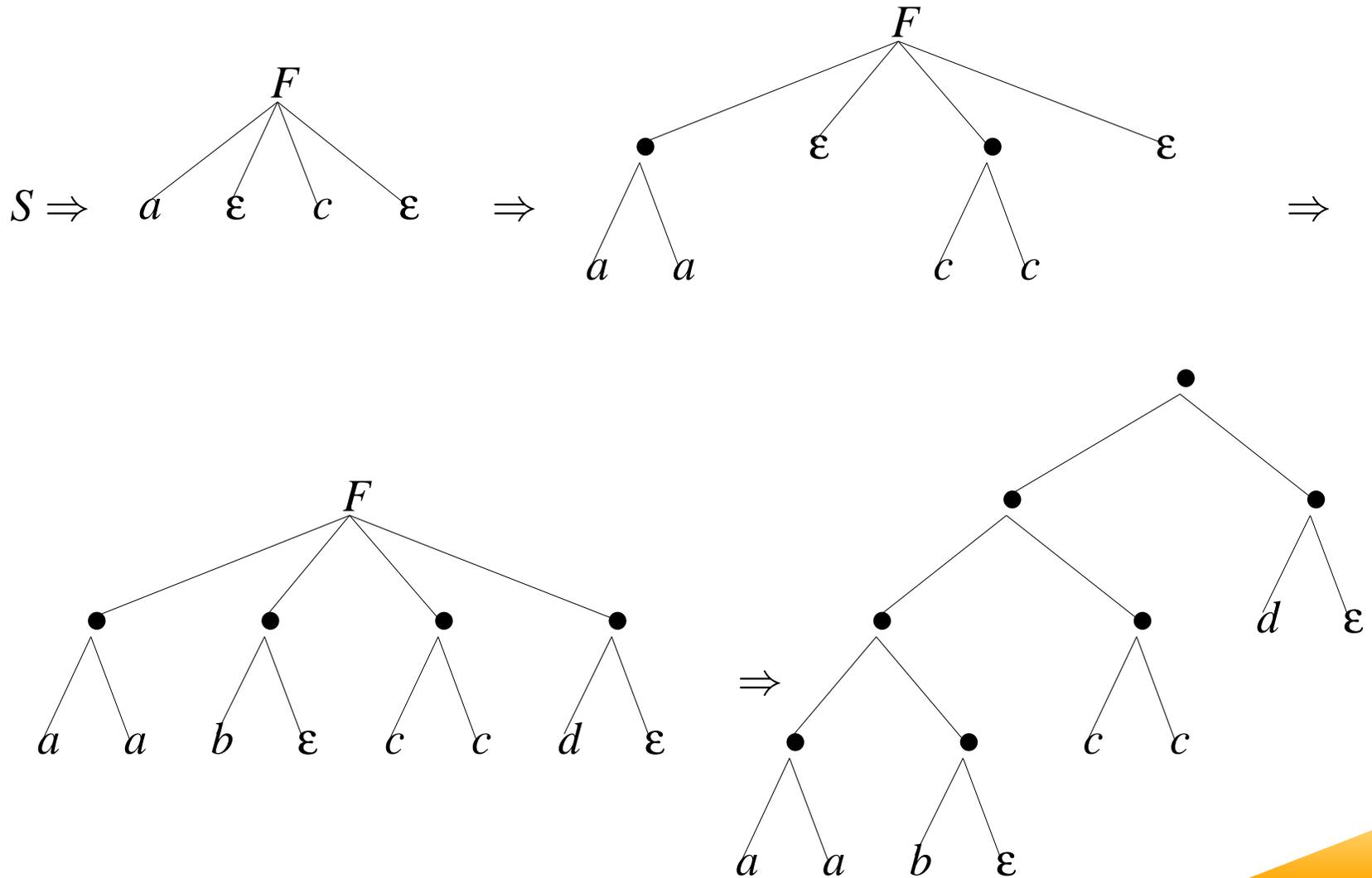
An example derivation of $aabccd$



$$F(x_1, x_2, x_3, x_4) \longrightarrow F(x_1, \bullet(b, x_2), x_3, \bullet(d, x_4))$$



An example derivation of $aabccd$



$$F(x_1, x_2, x_3, x_4) \longrightarrow \bullet(\bullet(\bullet(x_1, x_2), x_3), x_4)$$





MSO as a Query Language

MSO Queries

$$Q = \{t \models \Phi_r\}$$

$r \stackrel{def}{\iff}$ regular tree expression

$\Phi_r \stackrel{def}{\iff}$ translation of r into an MSO formula

Example:

$$(NP_{akk} \cup NP_{dat})^* \cdot _ \cdot (V_{akk} \cup V_{dat})^*$$



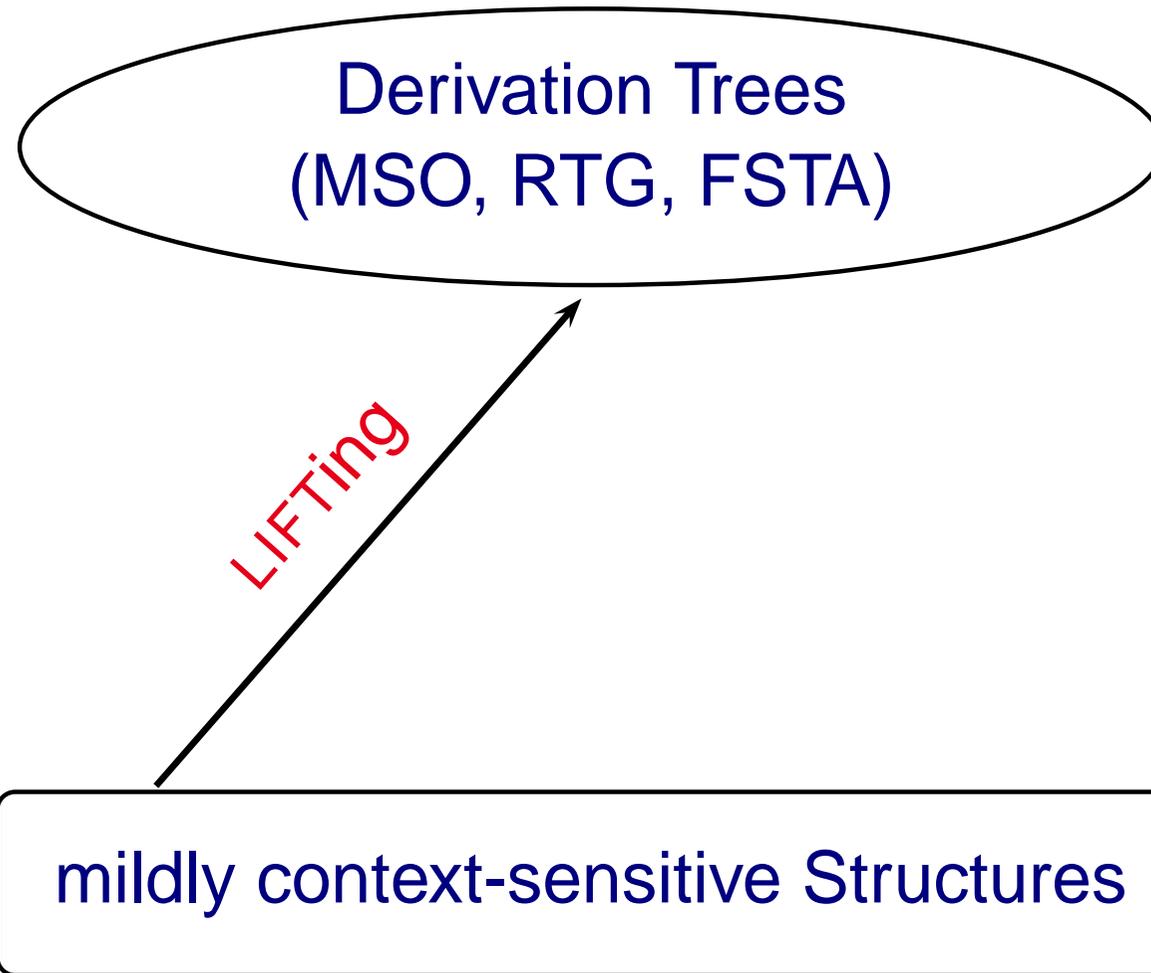
Overview Two-Steps

Derivation Trees
(MSO, RTG, FSTA)

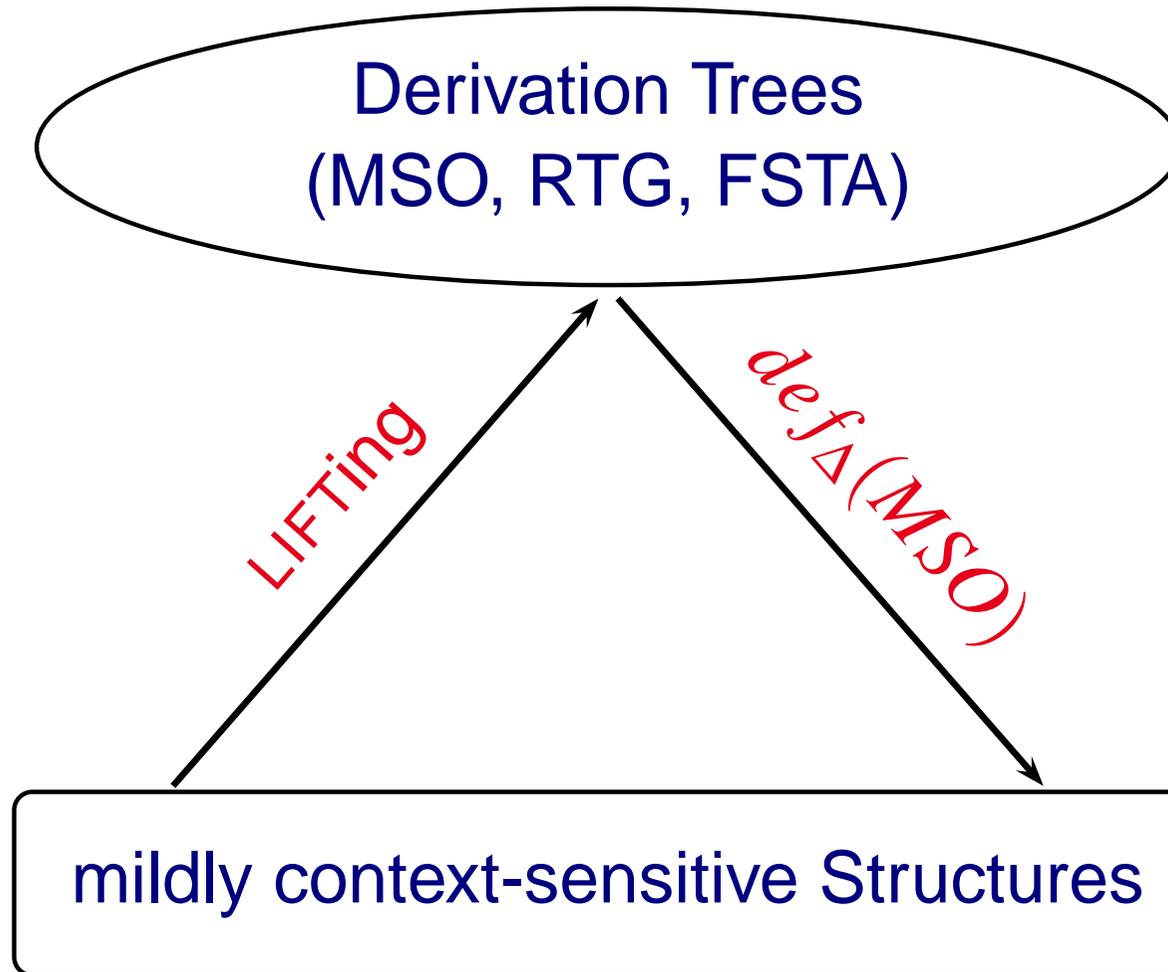
mildly context-sensitive Structures



Overview Two-Steps



Overview Two-Steps





Lifting

Lifting: Intuition

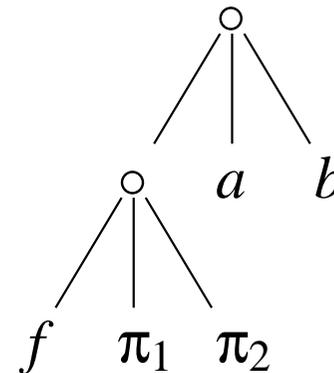
- Is a simple primitive recursive function.
- Makes control structure visible.
- All function symbols become constants.
- Resulting grammar is *regular*.
- Therefore expressible by an MSO-formula, and
- there is a tree automaton for the lifted grammar.



Lifting: Intuition

- Is a simple primitive recursive function.
- Makes control structure visible.
- All function symbols become constants.
- Resulting grammar is *regular*.
- Therefore expressible by an MSO-formula, and
- there is a tree automaton for the lifted grammar.

$$f(a, b) \rightsquigarrow (f \circ (\pi_1, \pi_2)) \circ (a, b) \rightsquigarrow$$



Lifting

For $k \geq 0$, $\text{LIFT}_{\Sigma}^k : T(\Sigma, \mathbf{X}_k) \rightarrow T(\Sigma^L, k)$ is defined as follows:

$$\text{LIFT}_{\Sigma}^k(x_i) = \pi_i^k$$

$$\text{LIFT}_{\Sigma}^k(f) = c_{(0,k)}(f') \text{ f\"ur } f \in \Sigma_0$$

$$\text{LIFT}_{\Sigma}^k(f(t_1, \dots, t_n)) =$$

$$c_{(n,k)}(f', \text{LIFT}_{\Sigma}^k(t_1), \dots, \text{LIFT}_{\Sigma}^k(t_n))$$

$$\text{f\"ur } f \in \Sigma_n, n \geq 1$$



The lifted Regular Tree Grammar: \mathcal{G}'

$$\Sigma_0^L = \{\varepsilon, a', b', c', d', S'\}$$

$$\Sigma_2^L = \{\bullet'\}$$

$$\Sigma_4^L = \{F'\} \cup \{\pi_1, \pi_2, \pi_3, \pi_4\}$$

$$\Sigma_{n,k}^L = \{c\} \text{ for simplicity}$$

$$\mathbf{F}_0 = \{S', F'\}$$

$$S = S'$$

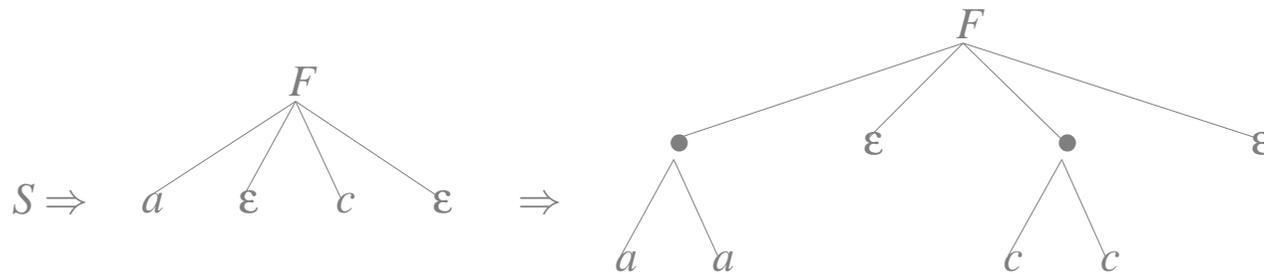
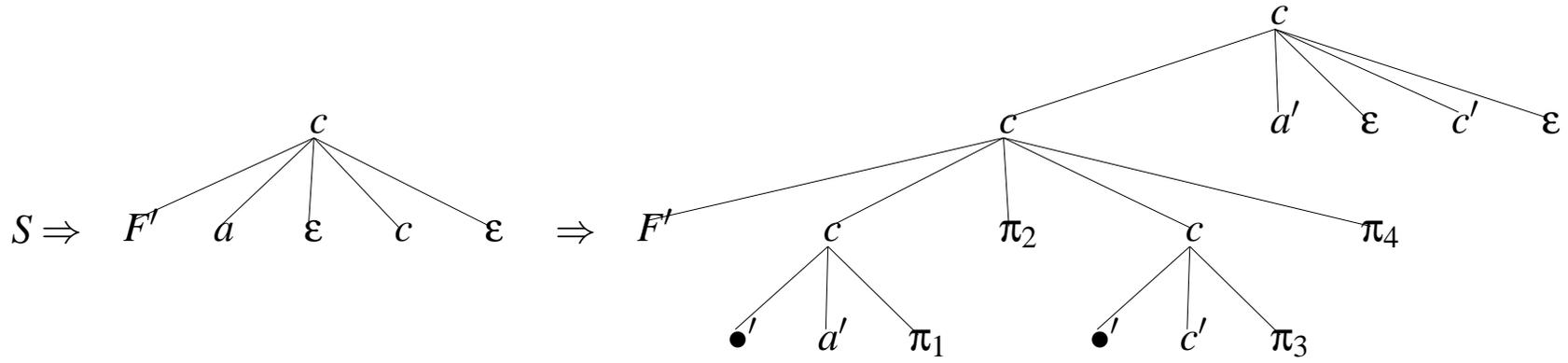
$$\mathbf{X} = \emptyset$$

$$\mathbf{P} = \left\{ \begin{array}{l} S' \longrightarrow \varepsilon \\ S' \longrightarrow c(F', a', \varepsilon, c', \varepsilon) \\ S' \longrightarrow c(F', \varepsilon, b', \varepsilon, d') \\ F' \longrightarrow c(F', c(\bullet', a', \pi_1), \pi_2, c(\bullet', c', \pi_3), \pi_4) \\ F' \longrightarrow c(F', \pi_1, c(\bullet', b', \pi_2), \pi_3, c(\bullet', d', \pi_4)) \\ F' \longrightarrow c(\bullet', c(\bullet', c(\bullet', \pi_1, \pi_2), \pi_3), \pi_4) \end{array} \right\}$$

$$F(x_1, x_2, x_3, x_4) \longrightarrow F(\bullet(a, x_1), x_2, \bullet(c, x_3), x_4)$$

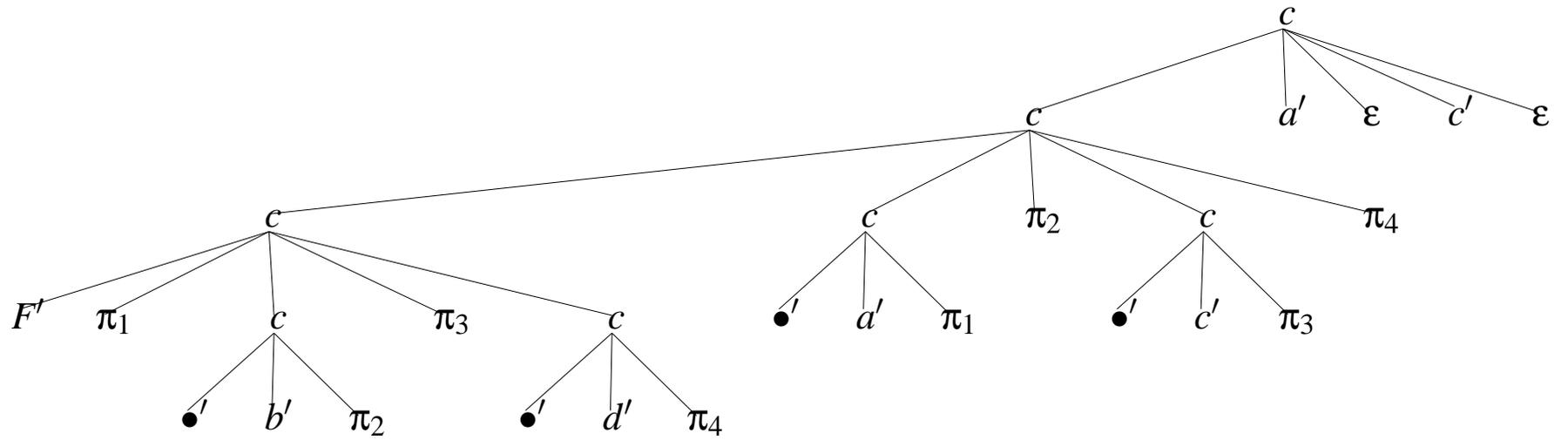


A sample derivation of \mathcal{G}'

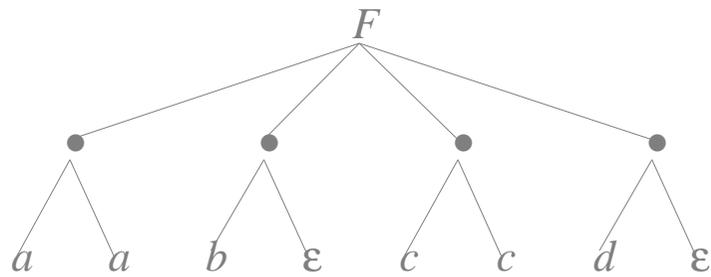


A sample derivation of \mathcal{G}'

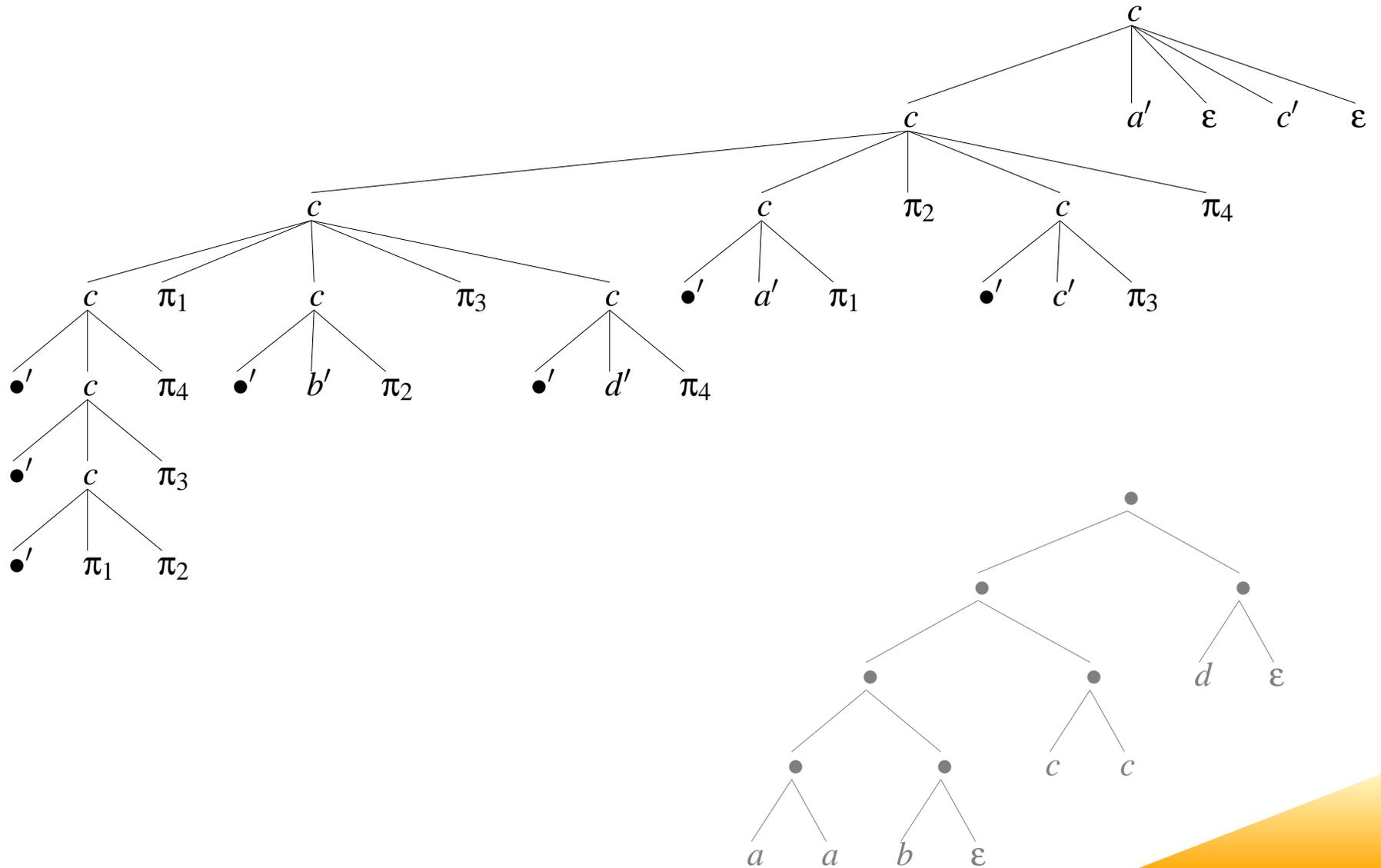
\Rightarrow



\Rightarrow



A sample derivation of \mathcal{G}'



Lifted Query

$$Q^L = \{t^L \models \varphi_r^L \wedge \psi\}$$

$$\varphi_r^L \stackrel{def}{\iff} \text{“Lifted” } \varphi_r$$

$$\psi \stackrel{def}{\iff} \text{“Lifted” specification of} \\ \text{cross-serial dependencies}$$

But: the resulting trees are not in the format of the treebank and for humans/linguists unreadable.





Reconstruction

Goal of the Reconstruction

- The intended tree is contained in the lifted tree.
- Homomorphism:

$$h(f') = f(x_1, \dots, x_n) \text{ for } f \in \Sigma_n$$

$$h(\pi_i^n) = x_i$$

$$h(c_{(n,k)}(t, t_1, \dots, t_n)) = h(t)[h(t_1), \dots, h(t_n)]$$

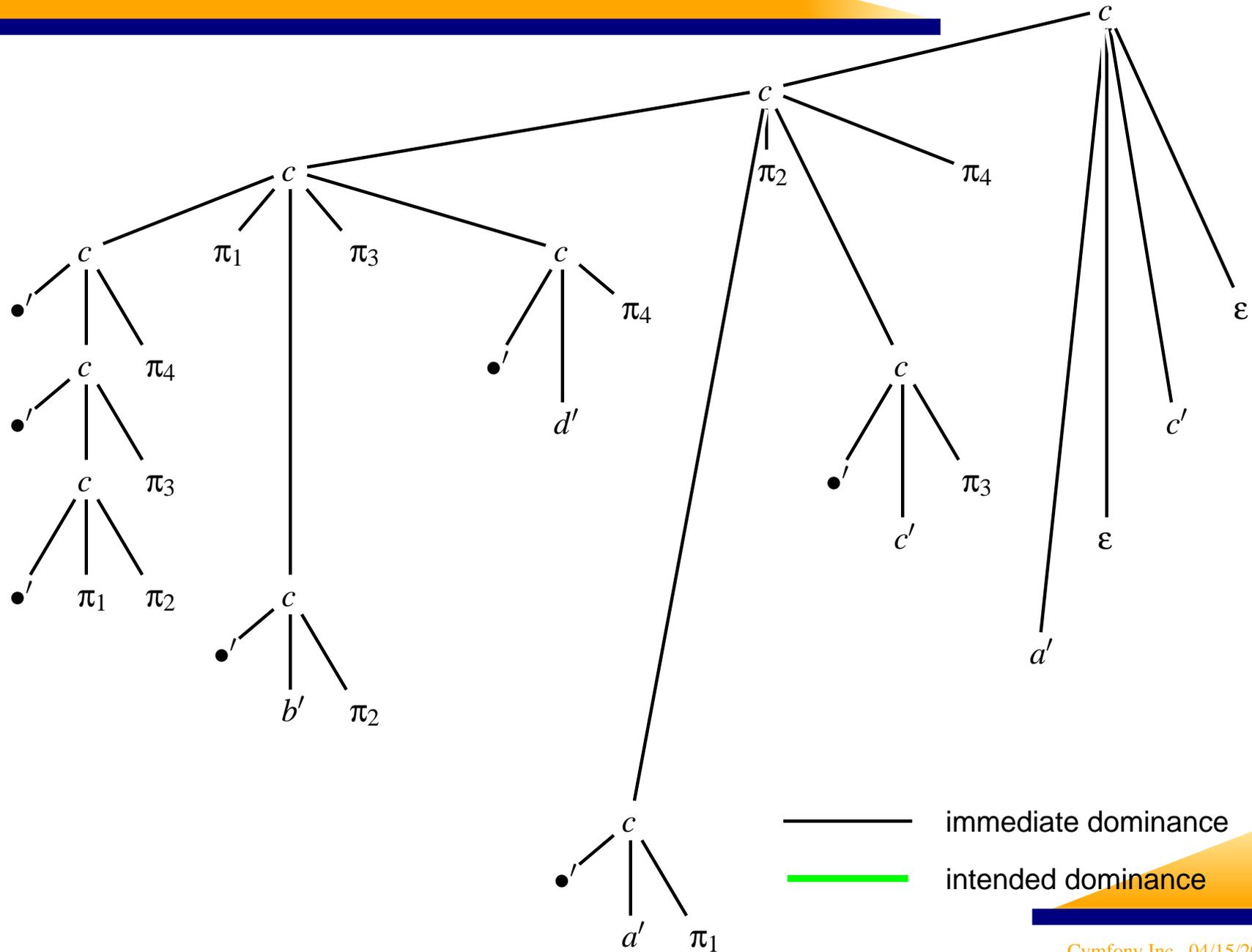
- Definition of a set of intended relations:

$$R^I = \{ \blacktriangleleft^* \text{ (dominance), } \blacktriangleleft, \blacktriangleleft^+, \triangleleft \text{ (precedence), } \dots \}$$

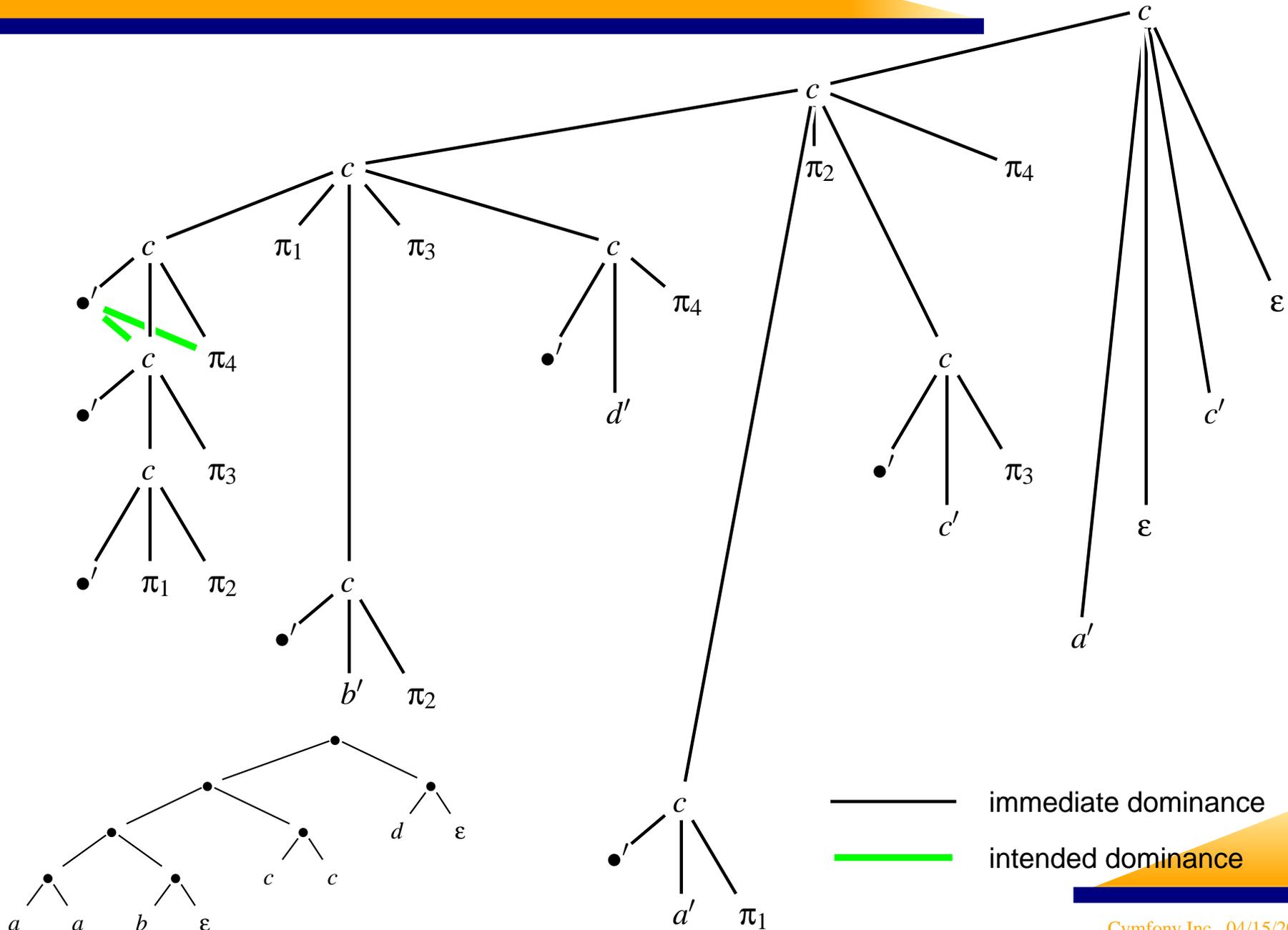
- ⇒ MSO-Transduction or Macro-Tree Transducer



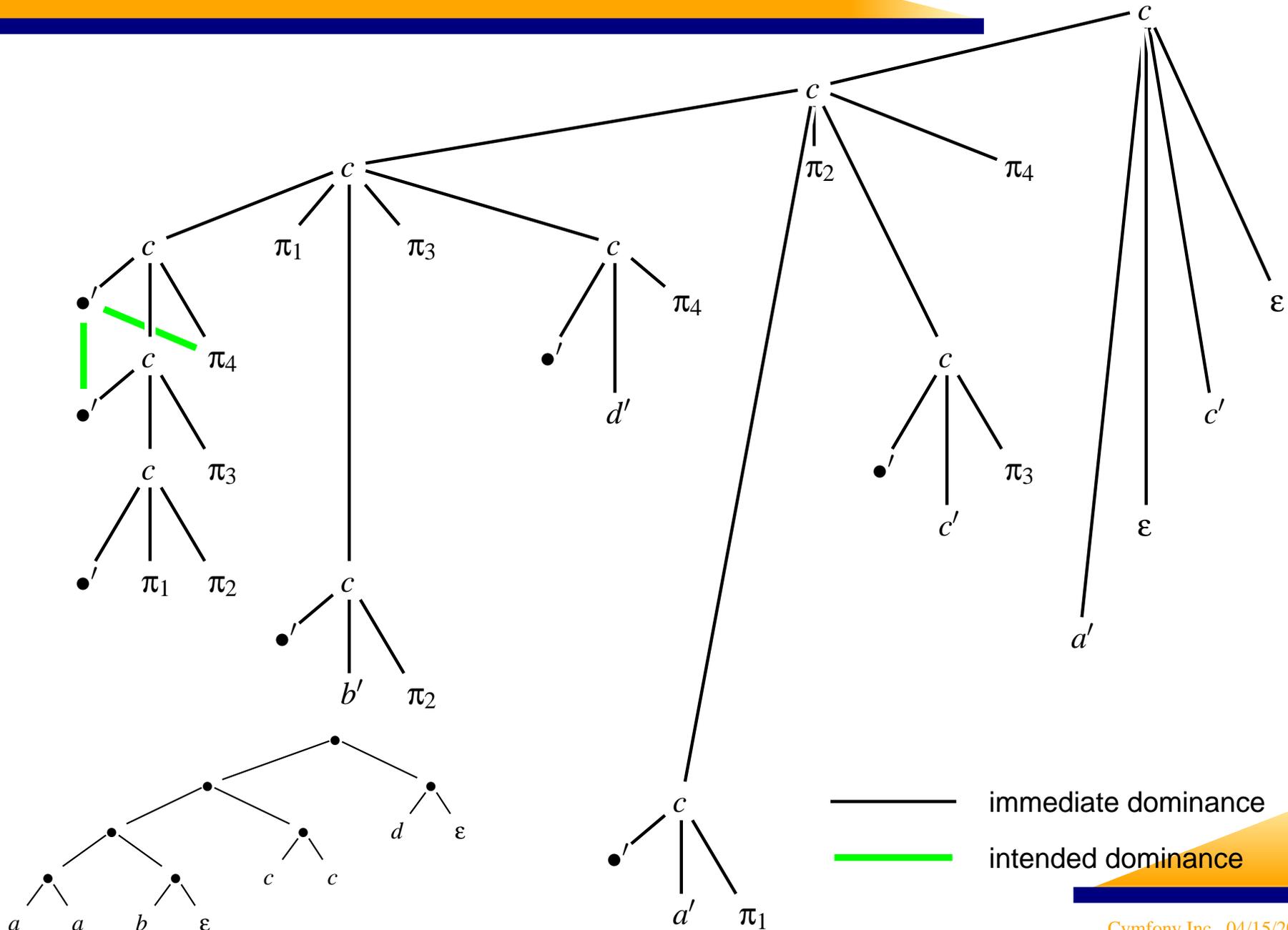
Reconstruction of the Intended Structures



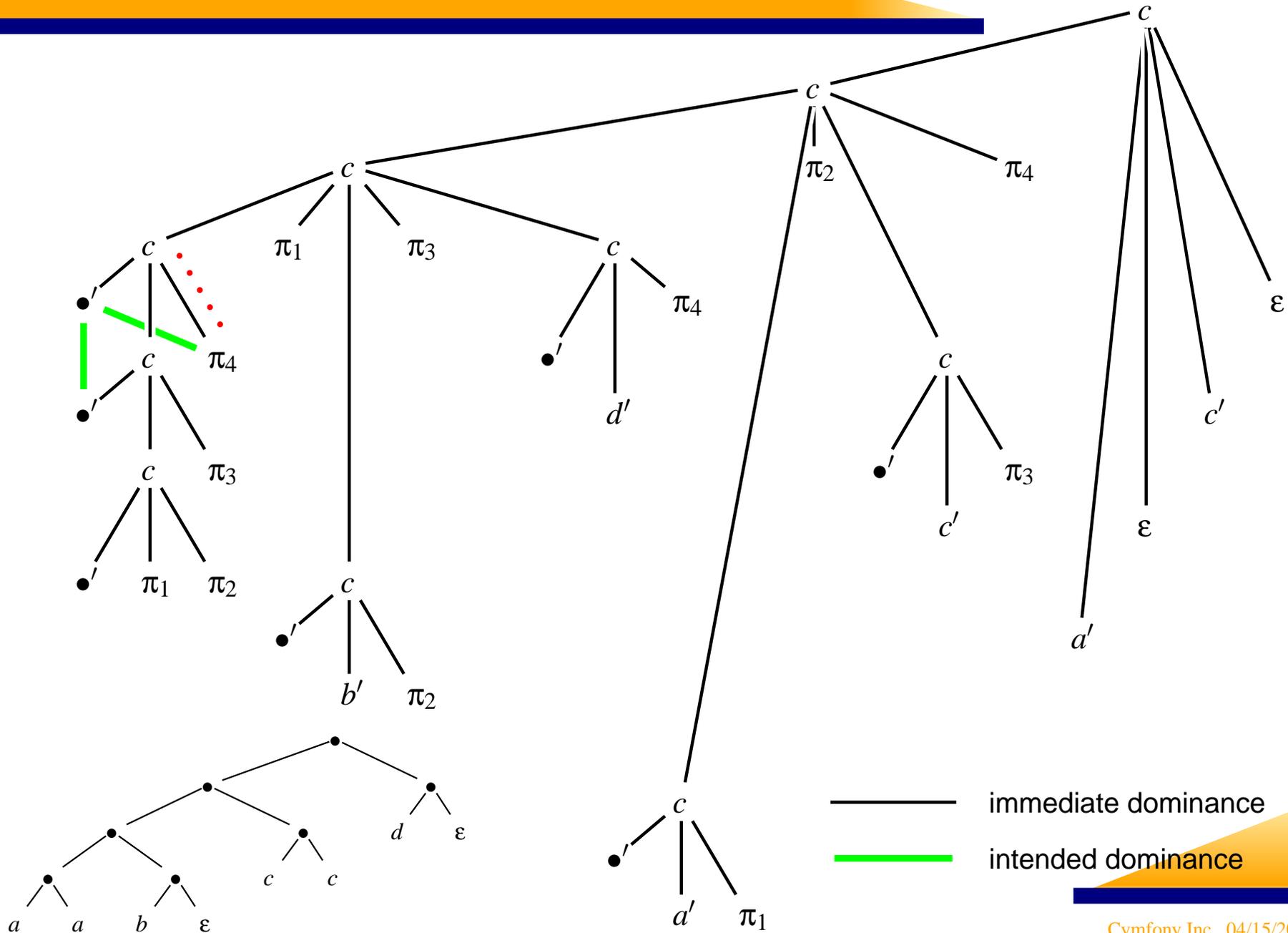
Reconstruction of the Intended Structures



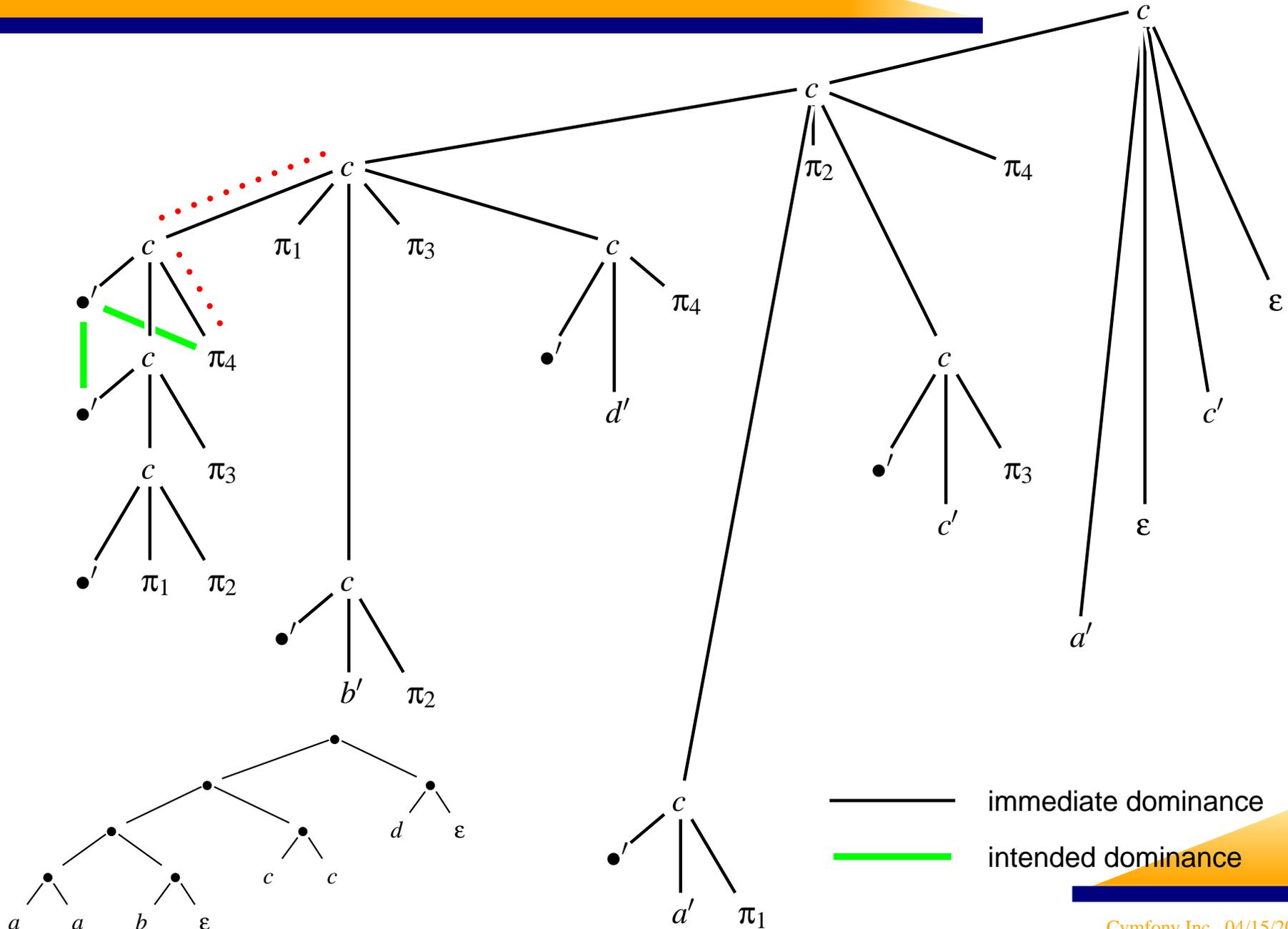
Reconstruction of the Intended Structures



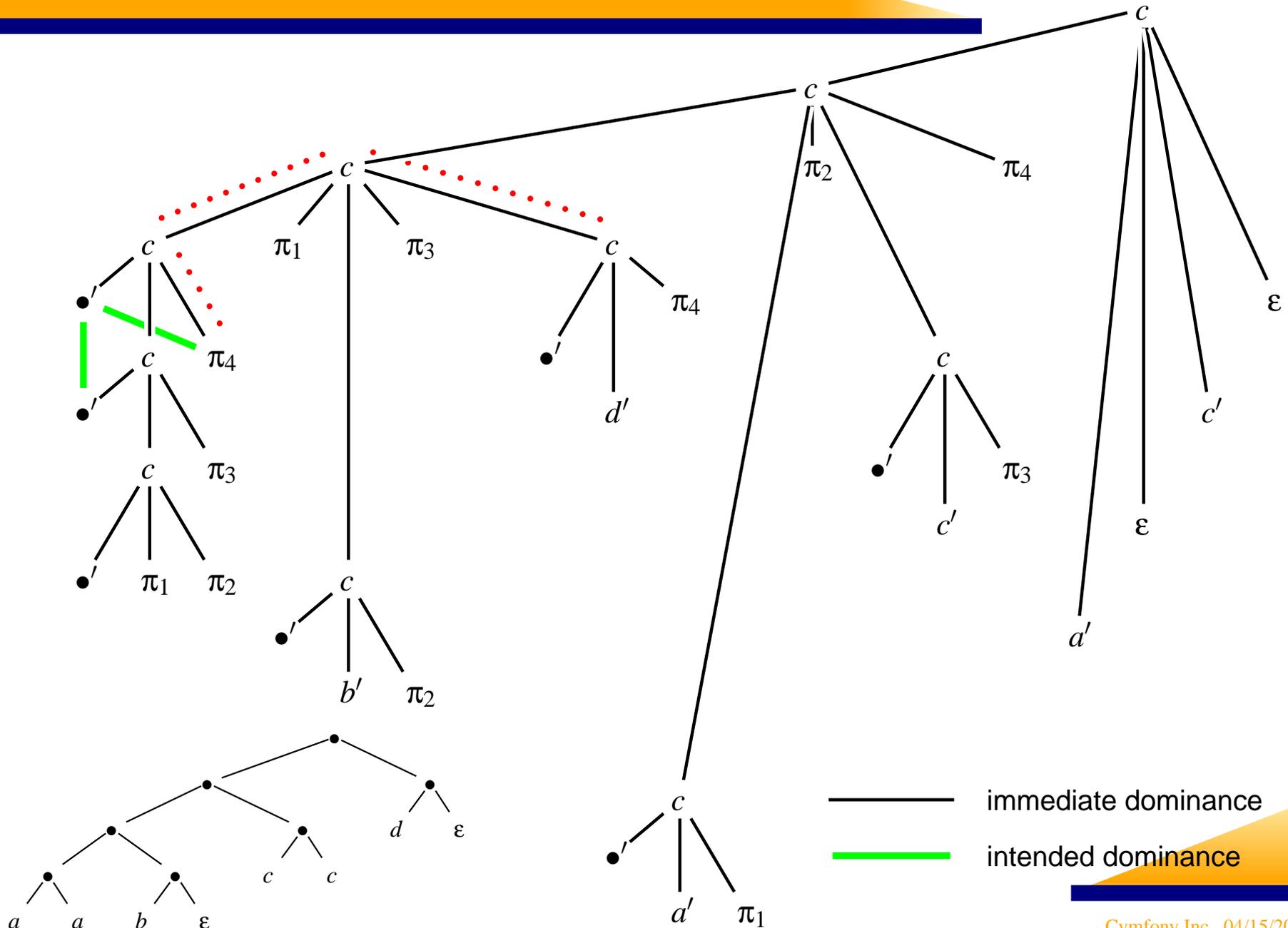
Reconstruction of the Intended Structures



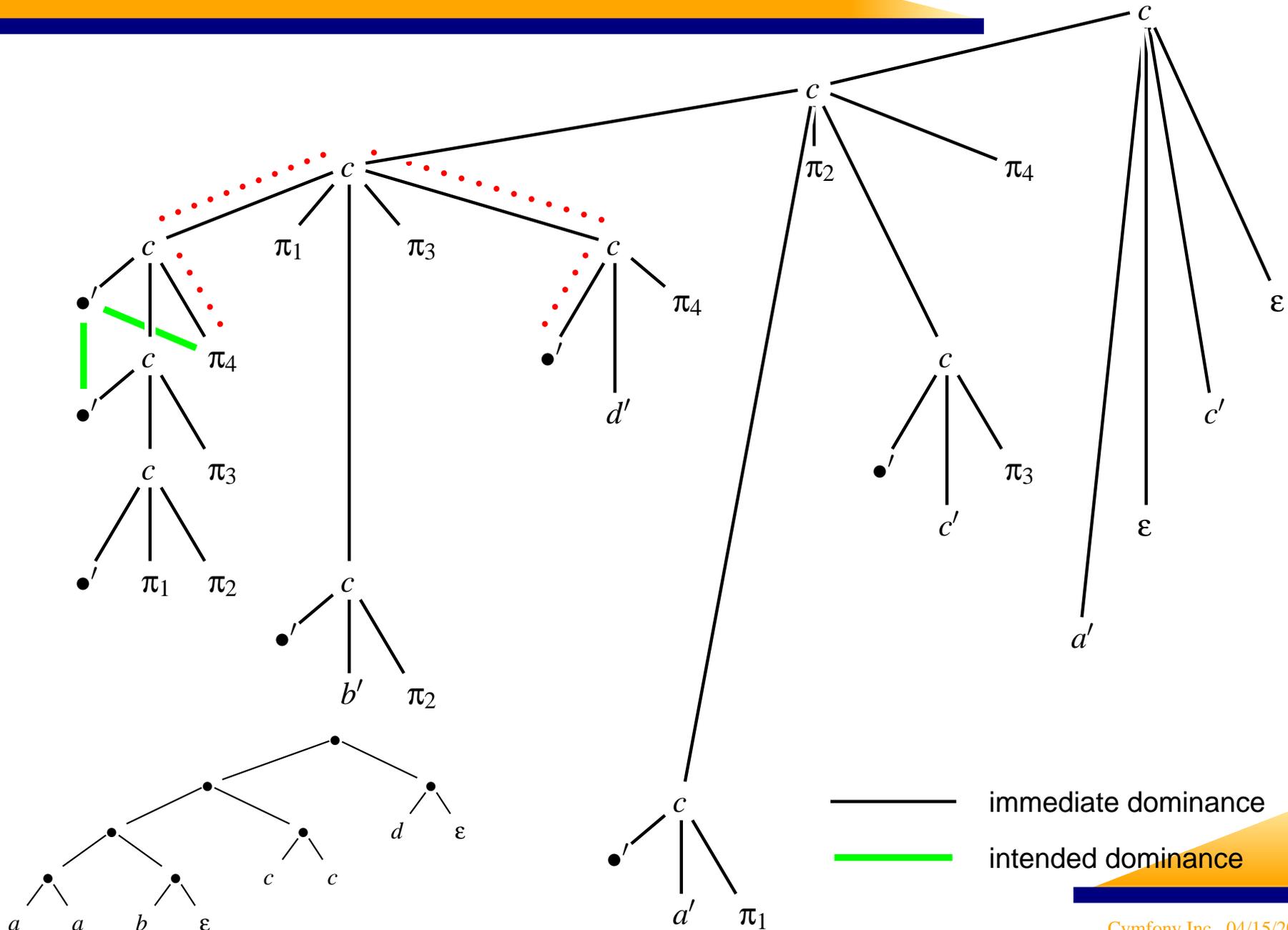
Reconstruction of the Intended Structures



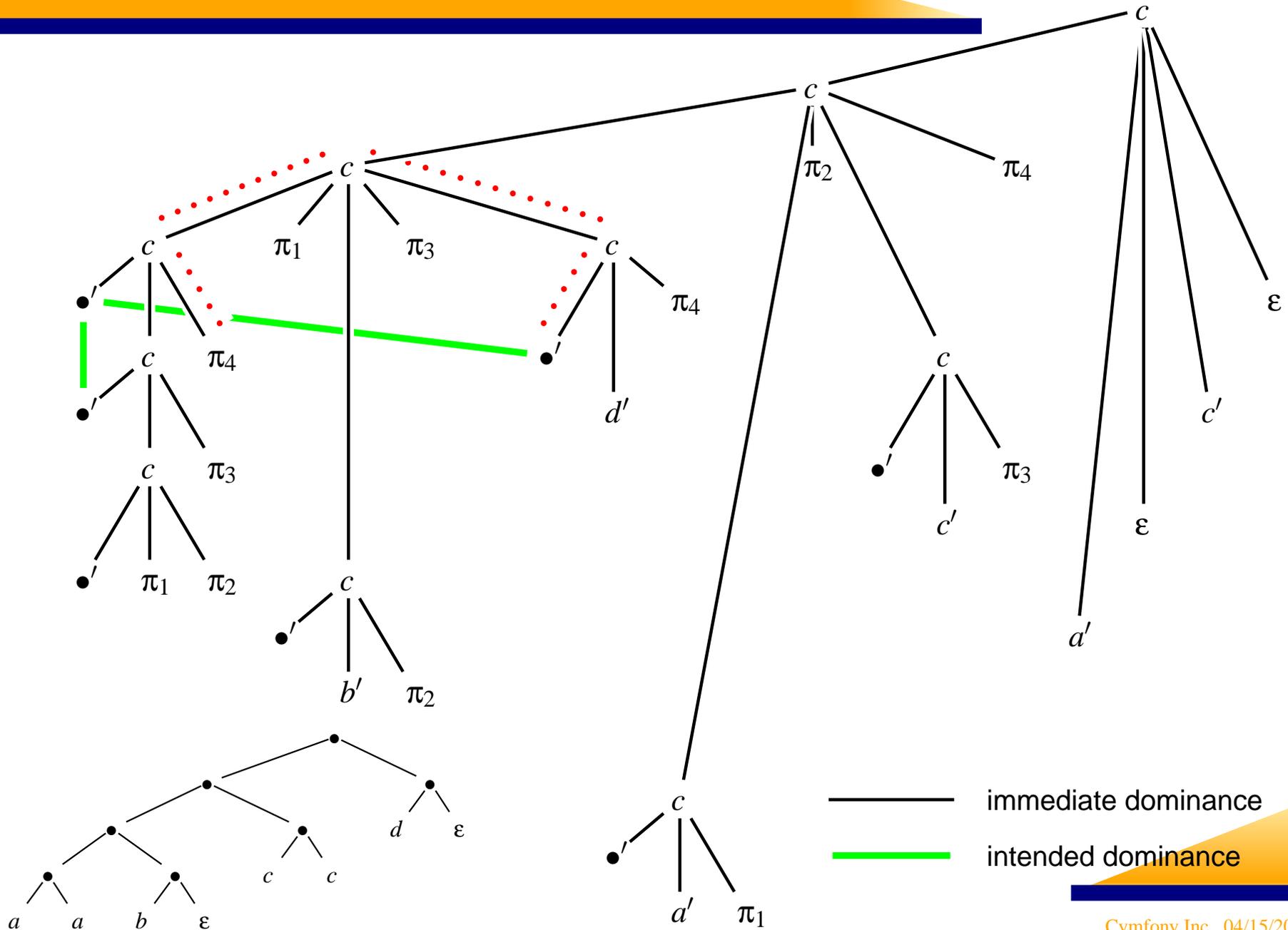
Reconstruction of the Intended Structures



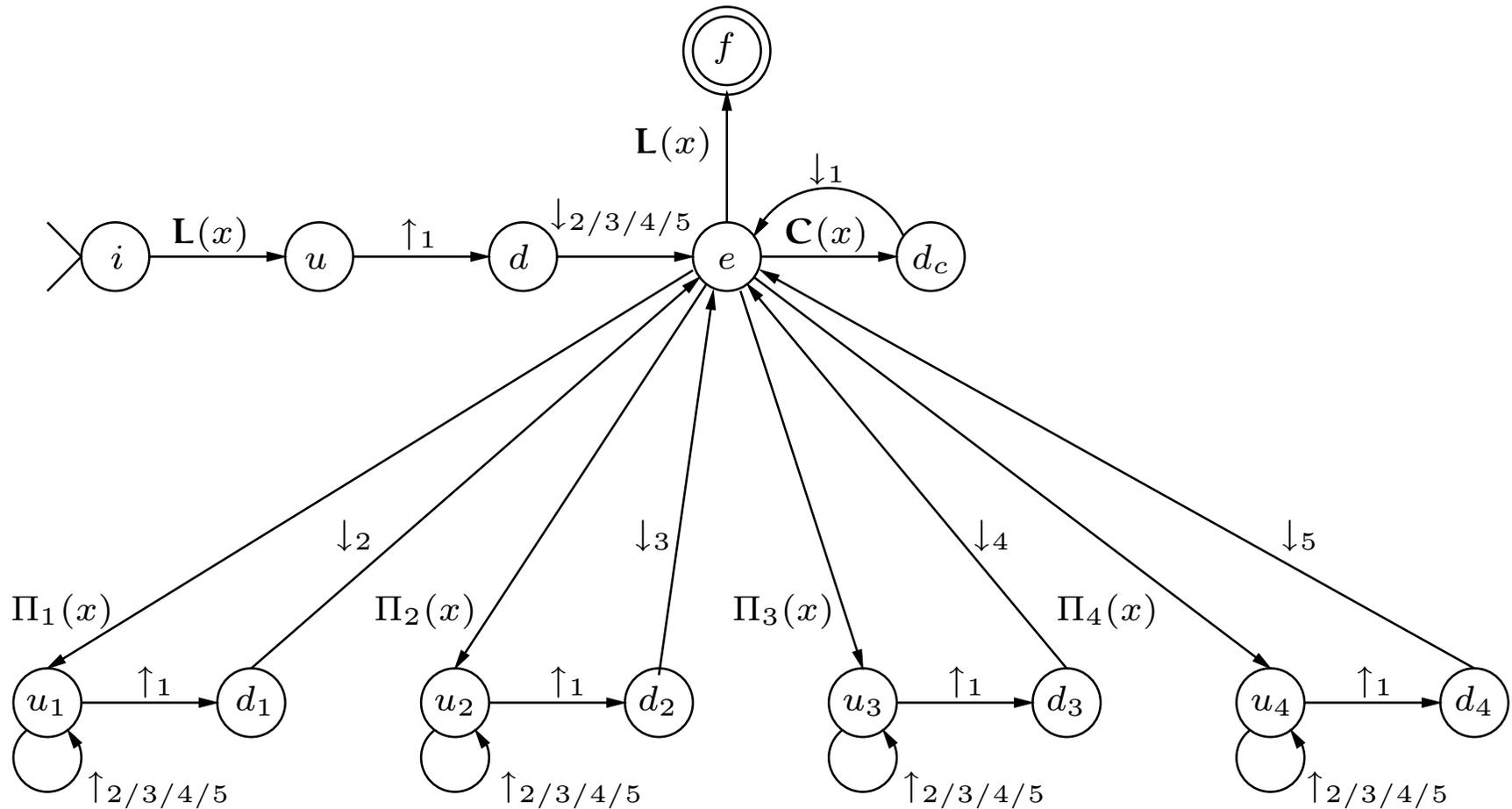
Reconstruction of the Intended Structures



Reconstruction of the Intended Structures



The tree-walking automaton \mathcal{A}



Regular Tree-Node Relations

For any tree t , such a tree-walking automaton \mathfrak{A} computes a node relation

$$R_t(\mathfrak{A}) = \{(x, y) \mid (x, q_i) \xRightarrow{*} (y, q_f) \text{ for} \\ \text{some } q_i \in I \text{ and some } q_f \in F\}$$

where for all states $q_i, q_j \in Q$ and nodes x, y in t

$$(x, q_i) \Rightarrow (y, q_j) \text{ iff } \exists d \in \Delta : (q_i, d, q_j) \in \delta \\ \text{and } y \text{ is reachable from } x \text{ in } t \text{ via } d.$$

If all the tests $\phi(x)$ of \mathfrak{A} are *MSO* definable, \mathfrak{A} specifies a *regular tree-node relation*, which is itself *MSO* definable.



MSO definable Transduction

$$\mathcal{R} \rightsquigarrow \mathcal{S}$$

$$\Delta = (\chi, \xi, (\theta_s)_{s \in \mathcal{S}})$$

χ the domain of the transduction

ξ the resulting domain of \mathcal{S}

θ_s the new relations

Definition by Courcelle 1997

Basic Idea Rabin 1965



Intended Structures via MSO transduction

$$(\varphi, \psi, (\theta_s)_{s \in S})$$

$$S = \{\triangleleft, \triangleleft^*, \triangleleft^+, \triangleleft, \dots\}$$

$$\chi = \varphi_r^L \wedge \psi$$

ξ = nodes with a “linguistic” label

$$\theta_{\triangleleft} = \text{trans}_{W_{\triangleleft}}(x, y)$$

$$\theta_{\triangleleft^*} = (\forall X)[\triangleleft\text{-closed}(X) \wedge x \in X \rightarrow y \in X]$$

$$\theta_{\triangleleft^+} = x \triangleleft^* y \vee x \not\approx y$$

$$\theta_{\triangleleft} = \text{trans}_{W_{\triangleleft}}(x, y)$$

$$\theta_{\text{label}} = \text{taken from } R$$



Top-Down Tree Transducer

A top-down tree transducer (TDTT) is a tuple $T = \langle Q, \Sigma, \Omega, q_0, P \rangle$ with states Q , ranked alphabets Σ and Ω (input and output), initial state q_0 and a finite set of productions P of the form

$$q(\sigma(x_1, \dots, x_n)) \longrightarrow t$$

where $n \geq 0$, $\sigma \in \Sigma_n$ and $t \in T(\Omega \cup \Sigma(X) \cup Q)$.

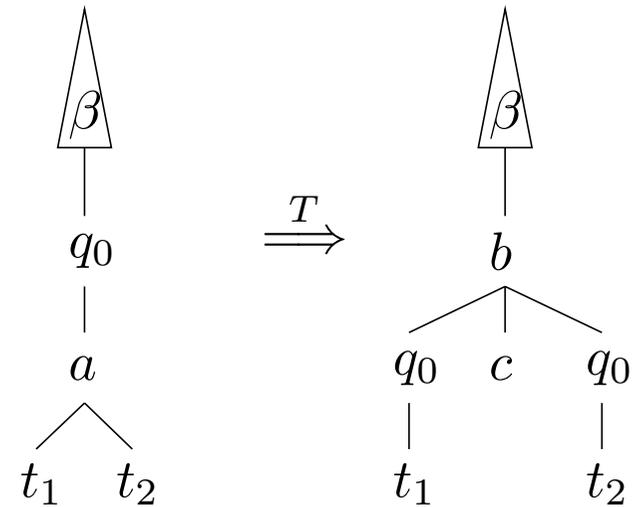
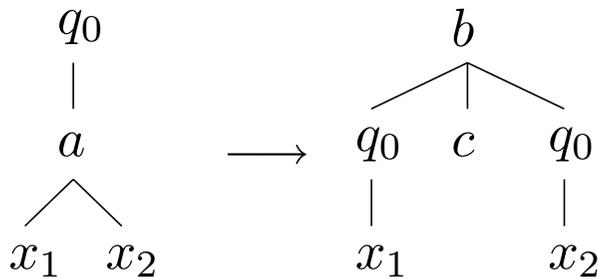
The transition relation (\xrightarrow{T}) is defined as usual. The transduction realized by a top-down tree transducer T is then defined to be $\{(t_1, t_2) \in T(\Sigma) \times T(\Omega) \mid q_0(t_1) \xrightarrow{T}^* t_2\}$.



Top-Down Tree Transducer

$$q_0(a(x_1, x_2)) \longrightarrow b(q_0(x_1), c, q_0(x_2))$$

$$q_0(p) \longrightarrow q$$



Macro Tree Transducer

A macro tree transducer (MTT) is a five-tuple $M = \langle Q, \Sigma, \Omega, q_0, P \rangle$ with Q a ranked alphabet of states, ranked alphabets Σ and Ω (input and output), initial state q_0 of rank 1, and a finite set of productions P of the form

$$q(\sigma(x_1, \dots, x_n), y_1, \dots, y_m) \longrightarrow t$$

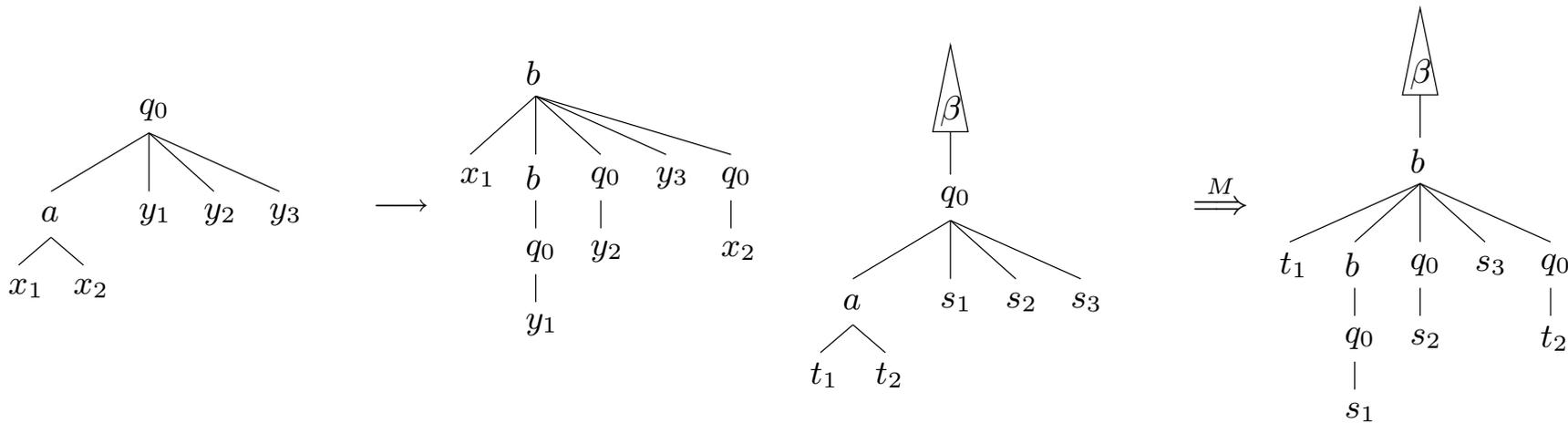
where $n, m \geq 0$, $q \in Q_{m+1}$, $\sigma \in \Sigma_n$ and $t \in RHS(\Sigma, \Omega, n, m)$.

The productions $p \in P$ of M are used as term rewriting rules in the usual way. The transition relation of M is denoted by \xrightarrow{M} . The transduction realized by M is the function $\{(t_1, t_2) \in T(\Sigma) \times T(\Omega) \mid (q_0, t_1) \xrightarrow{M}^* t_2\}$.



Macro Tree Transducer

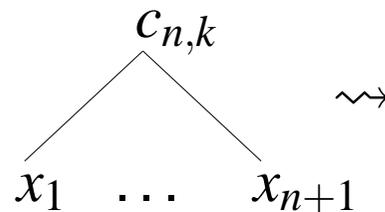
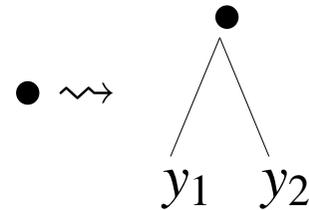
$$q_0(a(x_1, x_2), y_1, y_2, y_3) \longrightarrow b(x_1, b(q_0(y_1)), q_0(y_2), y_3, q_0(x_2))$$



Constructing an MTT: Intuition

$$\sigma \rightsquigarrow \sigma$$

$$\pi_i \rightsquigarrow y_i$$

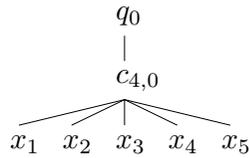


$$q_k(\underbrace{c_{n,k}(x_1, \dots, x_{n+1}), y_1, \dots, y_k}_{k+1}) \longrightarrow$$

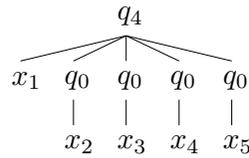
$$\underbrace{q_n(x_1, q_k(x_2, y_1, \dots, y_k), \dots, q_k(x_{n+1}, y_1, \dots, y_k))}_{n+1}$$



The Example MTT for \mathcal{G}'

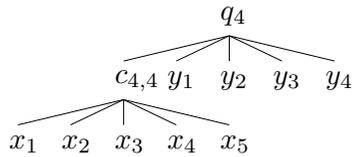


→

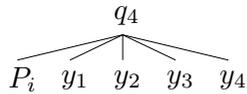
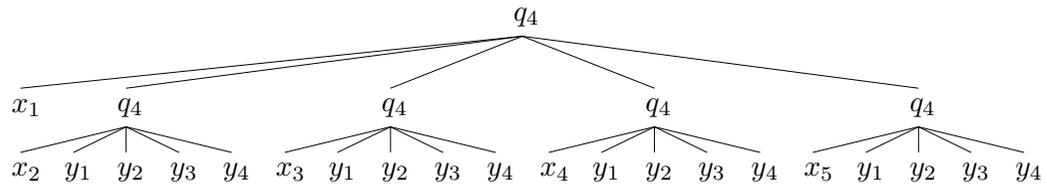


→ σ

for $\sigma \in \{a, b, c, d, \varepsilon\}$

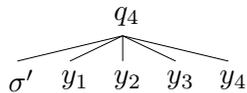


→



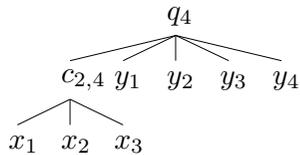
→ y_i

for $P_i = \pi_i$

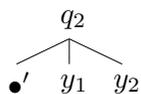
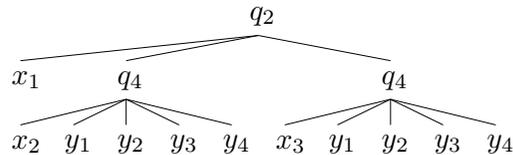


→ σ

for $\sigma \in \{a, b, c, d, \varepsilon\}$



→



→



Intended Query-Results

$$Q' = \{def_{\Delta}(t^L) \mid t^L \models \phi_r^L \wedge \psi\}$$

$def_{\Delta} \stackrel{def}{\iff}$ transduction from the domain of the lifted regular trees into original tree data base

$\Delta \stackrel{def}{\iff}$ MSO definable transduction





Summary and Outlook

Summary

- Querying context-sensitive relations with regular means.
- three “kinds” of trees
 - trees as elements of a free algebra
 - trees as elements of a free clone (Lawvere-Algebra)
 - trees as relational structures
- transduction is effectively realized via tree-walking automata or macro tree transducer



Outlook

- unranked trees
- implementation
- efficiency
 - MSO decidable on trees
 - MSO undecidable on graphs
 - MSO + one binary relation undecidable on trees
 - MSO \rightarrow tree automata hyper-exponential
 - “Guarded” MSO \rightarrow tree automata exponential
 - Lifting and transduction require linear time
 - Lifting of the candidate trees?



Overview

- Introduction and Motivation
- Semistructured Data
- Query Languages
- Intermediate Conclusion
- Secondary Edges
- Regular Queries of Context-Sensitive Relations
- MSO Queries
- Lifting
- Reconstruction
- Conclusion and Outlook





Appendix

Some Details and Definitions

MSO Logic on Trees (Rogers, 1994)

- K a set of individual constants
- P a set of predicate constants
- immediate dominance: \triangleleft
- reflexive dominance: \triangleleft^*
- proper dominance: \triangleleft^*
- proper precedence: \prec
- a set of *monadic* second-order variables (X, Y, \dots) ;
- a set of individual variables (x, y, \dots) ;
- boolean connectives: \neg, \wedge, \dots ;
- Quantifiers over individuals and sets: \exists, \forall ;
- equality and membership: \approx, \in .



Example formulas

$$X \subseteq Y \stackrel{def}{\iff} (\forall x)[x \in X \Rightarrow x \in Y]$$

$$\text{Sing}(X) \stackrel{def}{\iff}$$

% All subsets of X are equal to X or empty:

$$(\forall Y)[Y \subseteq X \Rightarrow (X \subseteq Y \vee (\forall Z)[Y \subseteq Z])] \wedge$$

% X is not the empty set:

$$(\exists Y)[X \not\subseteq Y]$$



Example formulas

$$\text{AC-Com}(x, y) \stackrel{\text{def}}{\iff}$$

% x c-commands y :

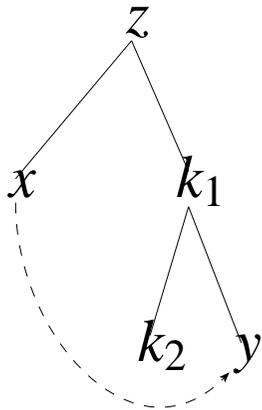
$$(\forall z)[z \triangleleft^+ x \Rightarrow z \triangleleft^+ y] \wedge \neg(x \triangleleft^* y) \wedge$$

% y does not c-command x :

$$\neg((\forall z)[z \triangleleft^+ y \Rightarrow z \triangleleft^+ x] \wedge \neg(y \triangleleft^* x)) \wedge$$

% x precedes y :

$$x \prec y$$



Example formulas

$LC(x, y) \stackrel{def}{\iff}$

$AC-Com(x, y) \wedge$

% There is no k_1 with property P :

$(\neg \exists k_1)[k_1 \in P \wedge$

% such that it intervenes between x and y :

$(\exists z)[z \triangleleft x \wedge z \triangleleft^+ k_1 \wedge k_1 \triangleleft^+ y]]$



Properties of MSO Logic on Trees

- Advantages:

- Concise Representation in a Powerful and Flexible Logic
- Descriptive Complexity Results
- Operational and Denotational Semantics
- Decidable
- Compilation: Formulas to Tree Automata

- Shortcomings

- Non-Elementary Complexity of the Decidability Algorithm:

$$|A| = 2^{2^{\cdot^{\cdot^{2^n}}}} \Big\}^m$$

- Descriptive Complexity: Only Context-Free String Languages



“Linguistic” Definitions

- A constraint is **satisfiable** iff it results in a non-empty automaton.
- A constraint is **valid** iff it results in the trivial non-empty automaton.
- A (maximally) binary branching **tree** is a subset T of the domain of MSO Logic on Trees. A **labeled tree** is a tuple $\langle T, F_1, \dots, F_k \rangle$.
- A **grammar** in this setting is a definition of a $k + 1$ -ary relation in MSO Logic on Trees designating all and only the well-formed labeled trees.
- The **recognition problem** is the determinization of the emptiness of the conjunction of the grammar formula/automaton with the input formula/automaton.
- The **parsing problem** is the conjunction of the grammar formula/automaton with the input formula/automaton.



Coding RTGs in MSO Logic (Thomas 1997)

- Define a tree automaton $\mathfrak{A}_{\mathcal{G}'}$ from the RTG \mathcal{G}'
- Code its behaviour in $\Phi_{\mathfrak{A}_{\mathcal{G}'}}$

$$(\exists X_0, \dots, X_m) \left[\bigwedge_{i \neq j} (\neg \exists y) [y \in X_i \wedge y \in X_j] \wedge \right.$$

$$(\forall x) [(\neg \exists y) [x \triangleleft y] \rightarrow x \in X_0] \wedge \quad \% \text{ Initial State}$$

$$\bigwedge_{1 \leq n \leq m} (\forall x_1, \dots, x_n, y) \left[\bigvee_{\substack{(i_1, \dots, i_n, \sigma, j) \in \alpha \\ 1 \leq k \leq n}} x_k \in X_{i_k} \wedge y \triangleleft x_k \wedge y \in X_j \wedge y \in P_\sigma \right]$$

$$\bigvee_{i \in F} (\exists x \forall y) [x \triangleleft^* y \wedge x \in X_i] \quad \% \text{ Root}$$



Lifting

Suppose that Σ is a ranked alphabet. The *derived* \mathbb{N} -sorted alphabet Σ^L is defined as follows:

For each $n \geq 0$,

$$\Sigma'_n = \{f' \mid f \in \Sigma_n\}$$

is a new set of symbols of sort n ;

back to Lifting



Lifting

Suppose that Σ is a ranked alphabet. The *derived* \mathbb{N} -sorted alphabet Σ^L is defined as follows:

for each $n \geq 1$ and each $i, 1 \leq i \leq n$,

$$\pi_i^n$$

is a new symbol, the i th projection symbol of sort n ;

back to Lifting



Lifting

Suppose that Σ is a ranked alphabet. The *derived* \mathbb{N} -sorted alphabet Σ^L is defined as follows:

for each $n \geq 0, k \geq 0$ the new symbol

$$c_{(n,k)}$$

is the (n, k) th composition symbol.

[back to Lifting](#)



Lifting

Suppose that Σ is a ranked alphabet. The *derived* \mathbb{N} -sorted alphabet Σ^L is defined as follows:

$$\Sigma_0^L = \Sigma'_0$$

$$\Sigma_n^L = \Sigma'_n \cup \{\pi_i^n \mid 1 \leq i \leq n\} \text{ for } n \geq 1$$

$$\Sigma_{n,k}^L = \{c_{(n,k)}\} \text{ for } n, k \geq 0$$

$$\Sigma_i^L = \emptyset \text{ otherwise}$$

[back to Lifting](#)



Model-Theoretic Interpretation

Basic Idea (Rabin 1965)

Obtaining a structure $\mathbb{B} = \langle B, Q \rangle$ from a structure $\mathbb{A} = \langle A, \mathcal{R} \rangle$ where \mathcal{R} and Q are families of relation symbols.

back to MSO transduction



Tree-Walking Automaton

A *tree-walking automaton (with tests)* (Bloem and Engelfriet 1997) over some ranked alphabet Σ is a finite automaton $\mathfrak{A} = (Q, \Delta, \delta, I, F)$ with states Q , directives Δ , transitions $\delta : Q \times \Delta Q$ and the initial and final states $I \subseteq Q$ and $F \subseteq Q$ which traverses a tree using three kinds of directives:

\uparrow_i “move up to the mother of the current node (if it has one and it is its i -th daughter)”,

\downarrow_i “move to the i -th daughter of the current node (if it exists)”,

$\phi(x)$ “verify that ϕ holds at the current node”.

[back](#) to the dominance automaton



Reflexive Closure in *MSO*

$$R\text{-closed}(X) \stackrel{\text{def}}{\iff} (\forall x, y)[x \in X \wedge R(x, y) \rightarrow y \in X]$$

$$R^*(x, y) \stackrel{\text{def}}{\iff} (\forall X)[R\text{-closed}(X) \wedge x \in X \rightarrow y \in X]$$

back to the MSO transduction



Walking language \Rightarrow MSO-formula

$$\mathit{trans}_{\emptyset}(x, y) = \perp$$

$$\mathit{trans}_{\downarrow_i}(x, y) = \mathit{edge}_i(x, y)$$

$$\mathit{trans}_{\uparrow_i}(x, y) = \mathit{edge}_i(y, x)$$

$$\mathit{trans}_{\phi(x)}(x, y) = \phi(x) \wedge x = y$$

$$\mathit{trans}_{W_1 \cup W_2}(x, y) = \mathit{trans}_{W_1}(x, y) \vee \mathit{trans}_{W_2}(x, y)$$

$$\mathit{trans}_{W_1 \cdot W_2}(x, y) = \exists z (\mathit{trans}_{W_1}(x, z) \wedge \mathit{trans}_{W_2}(z, y))$$

$$\mathit{trans}_{W^*}(x, y) = \mathit{trans}_W^*(x, y)$$

$$\mathit{trans}_W^*(x, y) = \forall X (\forall v, w (v \in X \wedge \mathit{trans}_W(v, w) \rightarrow w \in X) \\ \wedge x \in X \rightarrow y \in X)$$



Complex Right-Hand-Sides

Let Σ and Ω be ranked alphabets and $n, m \geq 0$. The set of right hand sides $RHS(\Sigma, \Omega, n, m)$ over Σ and Ω with n variables and m parameters is the smallest set $rhs \subseteq T(\Sigma \cup \Omega, X_n \cup Y_m)$ such that

1. $Y_m \subseteq rhs$
2. For $\omega \in \Omega_k$ with $k \geq 0$ and $\varphi_1, \dots, \varphi_k \in rhs$, $\omega(\varphi_1, \dots, \varphi_k) \in rhs$
3. For $q \in Q_{k+1}$ with $k \geq 0$, $x_i \in X_n$ and $\varphi_1, \dots, \varphi_k \in rhs$,
 $q(x_i, \varphi_1, \dots, \varphi_k) \in rhs$

back to MTTs

