

Applications of Tree Automata (in memoriam of FERENC GÉCSEG)

Andreas Maletti

Institute for Natural Language Processing
Universität Stuttgart, Germany

`maletti@ims.uni-stuttgart.de`

Szeged — October 9, 2015

Tree Automata

Literature



THATCHER & WRIGHT (1968):

Generalized Finite Automata Theory with an Application to a Decision Problem of Second-Order Logic.

Math. Systems Theory 2(1): 57–81



ROUNDS (1970):

Mappings and Grammars on Trees.

Mathematical Systems Theory 4(3): 257–287



GÉCSEG (1977):

Universal Algebras and Tree Automata.

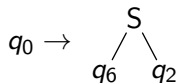
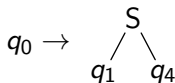
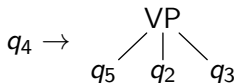
Proc. FCT: 98–112

Tree Automaton

tree automaton: (Q, Σ, I, R)

- finite set Q of states
- finite set Σ of terminals
- initial states $I \subseteq Q$
- finite set R of rules of the form $q \rightarrow \sigma(q_1, \dots, q_k)$
($\sigma \in \Sigma, k \geq 0, q, q_1, \dots, q_k \in Q$)

example rules:



derivation semantics: $\xi \Rightarrow \zeta$ for sentential forms $\xi, \zeta \in T_\Sigma(Q)$
if there exist leaf position w in ξ and rule $q \rightarrow r$ in R

$$\xi = \xi[q]_w$$

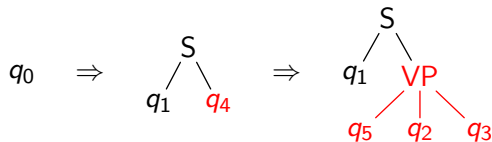
$$\zeta = \xi[r]_w$$

recognized tree language:

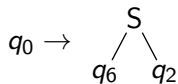
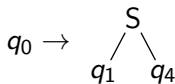
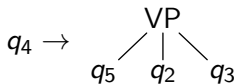
$$\{t \in T_\Sigma \mid \exists q_0 \in I: q_0 \Rightarrow^* t\}$$

Tree Automaton

example derivation:



example rules:



Standard references



GÉCSEG & STEINBY (1984):

Tree Automata.

Akadémiai Kiadó

2nd edition (2015): `arXiv:1509.06233`



GÉCSEG & STEINBY (1997):

Tree Languages.

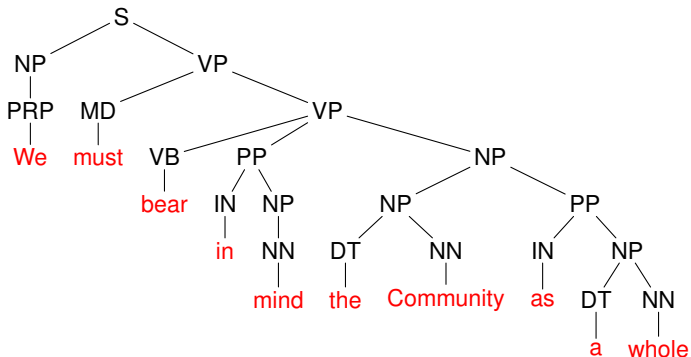
In *Handbook of Formal Languages 3*: 1–68, Springer

Parsing

Parsing

parsing: determining the syntactic structure of a sentence

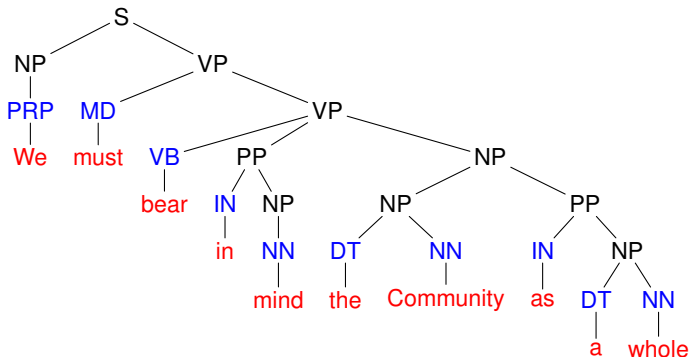
example: We must bear in mind the Community as a whole



Parsing

parsing: determining the syntactic structure of a sentence

example: We must bear in mind the Community as a whole



POS-tag: part-of-speech tag, “class” of a word

short history:

- before 1990

- ▶ hand-crafted CFG (refined via POS tags)
- ▶ corrections and selection by human annotators

- 1990s

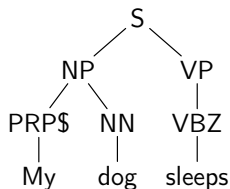
- ▶ PENN tree bank (1M words)
- ▶ weighted CFG (statistic methods)
- ▶ WALL STREET JOURNAL tree bank (30M words)

- since 2000

- ▶ weighted tree automata

popular parsers (for English):

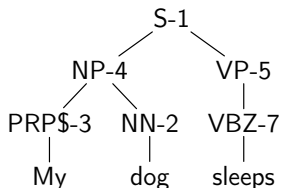
- **COLLINS parser:** CFG with manual subcategorization
[**COLLINS**, 1999]
- **STANFORD parser:** CFG with manual subcategorization
[**KLEIN, MANNING**, 2003]
- **BLLIP parser:** CFG with manual subcategorization and reranking
[**CHARNIAK, JOHNSON**, 2005]
- **BERKELEY parser:** CFG with automatic subcategorization
[**PETROV, KLEIN**, 2007]



tags:

- official tags often conservative
 - ▶ **English:** \approx 50 tags
 - ▶ **German:** \gg 200 tags

ADJA-Sup-Dat-Sg-Fem



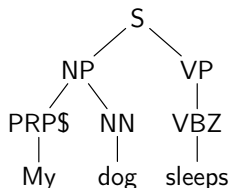
tags:

- official tags often conservative

- ▶ **English:** \approx 50 tags
- ▶ **German:** \gg 200 tags

ADJA-Sup-Dat-Sg-Fem

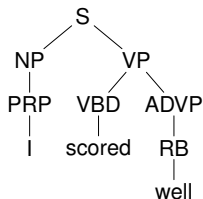
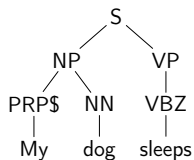
- modern parsers use refined tags \rightarrow **subcategorization**



tags:

- official tags often conservative
 - ▶ **English:** \approx 50 tags
 - ▶ **German:** \gg 200 tags
- modern parsers use refined tags \rightarrow **subcategorization**
- but return parse over official tags \rightarrow **relabeling**

ADJA-Sup-Dat-Sg-Fem

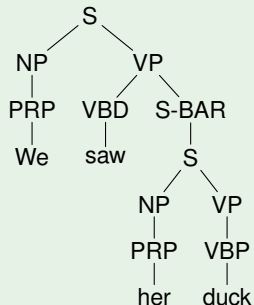
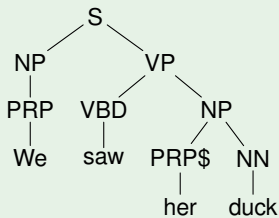


read off CFG productions:

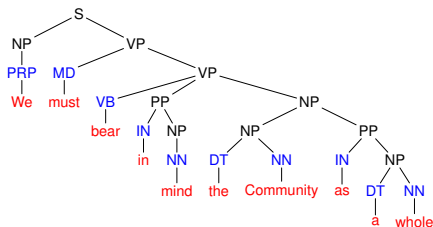
$S \rightarrow NP VP$
 $PRP\$ \rightarrow My$
 $VP \rightarrow VBZ$
 $NP \rightarrow PRP$
 $VP \rightarrow VBD ADVP$
 $ADVP \rightarrow RB$

$NP \rightarrow PRP\$ NN$
 $NN \rightarrow dog$
 $VBZ \rightarrow sleeps$
 $PRP \rightarrow I$
 $VBD \rightarrow scored$
 $RB \rightarrow well$

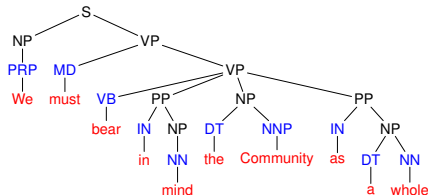
Weights for disambiguation



BERKELEY parser [Reference]:



BLLIP parser:



State-of-the-art models

- CFG with subcategorization (CFG_{sub})
- tree substitution grammars with subcategorization (TSG_{sub})
[SHINDO et al., 2012]
- (both as expressive as weighted tree automata)
- other models

State-of-the-art models

- CFG with subcategorization (CFG_{sub})
- tree substitution grammars with subcategorization (TSG_{sub})
[SHINDO et al., 2012]
- (both as expressive as weighted tree automata)
- other models

Experiment [SHINDO et al., 2012]

grammar model	F_1
wCFG	72.6
wTSG [COHN et al., 2010]	84.7
wCFG _{sub} [PETROV, 2010]	91.8
wTSG _{sub} [SHINDO et al., 2012]	92.4

grammar with subcategorization:

- a grammar G generating $L(G) \subseteq T_{\Sigma}(W)$
- a (total) mapping $\rho: \Sigma \rightarrow \Delta$

(grammar with relabeling)

(subcategorized trees)

(functional relabeling)

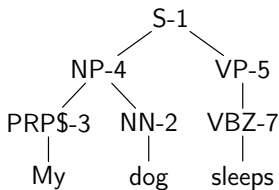
grammar with subcategorization:

- a grammar G generating $L(G) \subseteq T_{\Sigma}(W)$
- a (total) mapping $\rho: \Sigma \rightarrow \Delta$

(grammar with relabeling)

(subcategorized trees)

(functional relabeling)



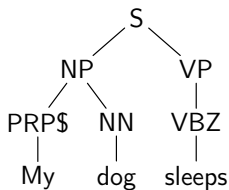
grammar with subcategorization:

- a grammar G generating $L(G) \subseteq T_{\Sigma}(W)$
- a (total) mapping $\rho: \Sigma \rightarrow \Delta$

(grammar with relabeling)

(subcategorized trees)

(functional relabeling)



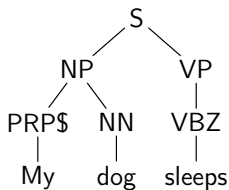
grammar with subcategorization:

- a grammar G generating $L(G) \subseteq T_{\Sigma}(W)$
- a (total) mapping $\rho: \Sigma \rightarrow \Delta$

(grammar with relabeling)

(subcategorized trees)

(functional relabeling)



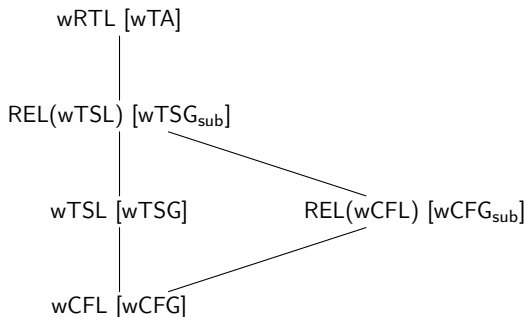
Semantics

$$L(G, \rho) = \rho(L(G)) = \{\rho(t) \mid t \in L(G)\}$$

Language class: **REL**(\mathcal{L}) for language class \mathcal{L}

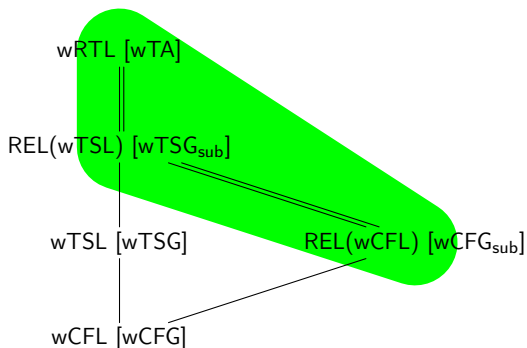
Theorem

$$\text{REL}(\text{wCFL}) = \text{REL}(\text{wTSL}) = \text{wRTL}$$



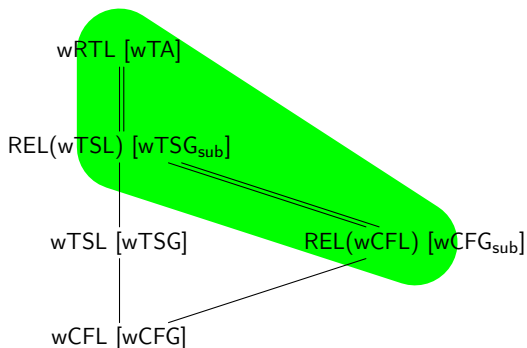
Theorem

$$\text{REL}(\text{wCFL}) = \text{REL}(\text{wTSL}) = \text{wRTL}$$



Theorem

$$\text{REL}(\text{wCFL}) = \text{REL}(\text{wTSL}) = \text{wRTL}$$



hence: **subcategorization** \approx **finite-state**

Comparison:

- 1 rule of subcategorized grammar:

$S-1 \rightarrow ADJP-2 S-1$

weight: 0.003545

with relabeling $\rho(S-1) = S, \dots$

Comparison:

- 1 rule of subcategorized grammar:

$$S-1 \rightarrow ADJP-2 S-1 \quad \text{weight: } 0.003545$$

with relabeling $\rho(S-1) = S, \dots$

- 2 corresponding rule of tree automaton

$$S-1 \rightarrow S(ADJP-2, S-1) \quad \text{weight: } 0.003545$$

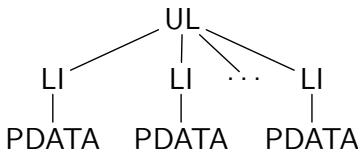
XML

History:

- XML = eXtensible Markup Language
- developed by W3C
- readable by humans and machines
- open standard; widely supported
- uncountable applications (MS Office, Apple iWork, etc.)

Example:

```
<UL>  
<LI> first item </LI>  
<LI> second item </LI>  
...  
<LI> last item </LI>  
</UL>
```



Similar development:





- 2000:
 - ▶ original DTD (Document Type Definition) specification
 - ▶ essentially a CFG
- 2005–2010:
 - ▶ more expressive specifications
 - ▶ DSD, XML-Schema, ISO Relax NG, etc.
 - ▶ all implementable by (unranked) tree automata
- since 2010:
 - ▶ DTD largely irrelevant

Similar development:

- 2000:
 - ▶ original DTD (Document Type Definition) specification
 - ▶ essentially a CFG
- 2005–2010:
 - ▶ more expressive specifications
 - ▶ DSD, XML-Schema, ISO Relax NG, etc.
 - ▶ all implementable by (unranked) tree automata
- since 2010:
 - ▶ DTD largely irrelevant

tree automata as basis for XML documents

Selected Literature

-  **KLEIN, MANNING:** *Accurate Unlexicalized Parsing*
Proc. ACL **2003**
-  **MURATA, LEE, MANI, KAWAGUCHI:** *Taxonomy of XML schema languages using formal language theory*
ACM Trans. Internet Technology **2005**
-  **PETROV, BARRETT, THIBAU AND KLEIN:** *Learning Accurate, Compact, and Interpretable Tree Annotation.* Proc. ACL **2006**
-  **SHINDO, MIYAO, FUJINO AND NAGATA:** *Bayesian Symbol-Refined Tree Substitution Grammars for Syntactic Parsing.*
Proc. ACL **2012**