

# Characterizations of subregular tree languages

Andreas Maletti

Institute of Computer Science, Universität Leipzig, Germany

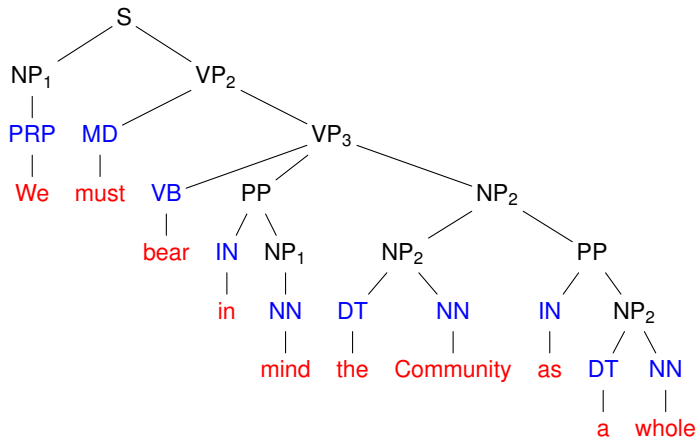
`andreas.maletti@uni-leipzig.de`



MPI, Leipzig — November 29, 2018

# Constituent Syntax Tree

Syntax tree for **We must bear in mind the Community as a whole**



# Constituent Syntax Tree

## Definition (Tree)

For sets  $\Sigma$  and  $V$  and  $\text{rk}: \Sigma \rightarrow \mathbb{N}$ , let  $T_{(\Sigma, \text{rk})}(V)$  be the least set  $T$  s.t.

- 1  $V \subseteq T$
- 2  $\sigma(t_1, \dots, t_{\text{rk}(\sigma)}) \in T$  for all  $\sigma \in \Sigma$  and  $t_1, \dots, t_{\text{rk}(\sigma)} \in T$

# Constituent Syntax Tree

## Definition (Tree)

For sets  $\Sigma$  and  $V$  and  $\text{rk}: \Sigma \rightarrow \mathbb{N}$ , let  $T_{(\Sigma, \text{rk})}(V)$  be the least set  $T$  s.t.

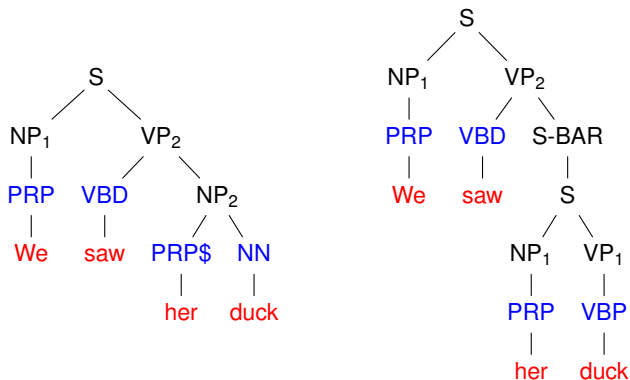
①  $V \subseteq T$

②  $\sigma(t_1, \dots, t_{\text{rk}(\sigma)}) \in T$  for all  $\sigma \in \Sigma$  and  $t_1, \dots, t_{\text{rk}(\sigma)} \in T$

- 2nd item: **top concatenation**
- 'rk' often implicit (we often write  $T_{\Sigma}(V)$  instead of  $T_{(\Sigma, \text{rk})}(V)$ )
- **( $\Sigma$ -)tree language** = set  $L \subseteq T_{\Sigma}(\emptyset)$  of trees

# Constituent Syntax Trees

Syntax tree is not unique (weights are used for disambiguation)



# Tree Languages

## Representations

- enumeration

# Tree Languages

## Representations

- enumeration
- local tree languages
- tree substitution languages
- **regular tree languages**

# Tree Languages

## Representations

- enumeration
- local tree languages
- tree substitution languages
- **regular tree languages**

### Definition (Regular tree language [Brainerd 1984])

$L \subseteq T_{\Sigma}(\emptyset)$  **regular** iff  $\exists$  congruence  $\cong$  (top-concatenation) on  $T_{\Sigma}(\emptyset)$  s.t.

- 1  $\cong$  has finite index (finitely many equiv. classes)
- 2  $\cong$  saturates  $L$ ; i.e.  $L = \bigcup_{t \in L} [t]_{\cong}$



# Regular Tree Languages

Examples for  $\Sigma = \{\sigma/2, \delta/2, \alpha/0\}$ :

- 2 equivalence classes ( $L$  and  $T_\Sigma(\emptyset) \setminus L$ )

$$L = \{t \in T_\Sigma(\emptyset) \mid t \text{ contains odd number of } \alpha\}$$

- 3 equivalence classes (“no  $\sigma$ ”, “some  $\sigma$ , but legal”, illegal)

$$L' = \{t \in T_\Sigma(\emptyset) \mid \sigma \text{ never below } \delta\}$$

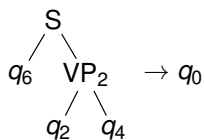
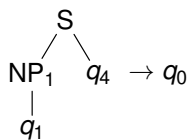
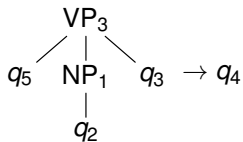
# Regular Tree Languages

## Definition (Regular tree grammar [Brainerd, 1969])

**Regular tree grammar**  $G = (Q, \Sigma, I, P)$

- alphabet  $Q$  of nonterminals and initial nonterminals  $I \subseteq Q$
- alphabet of terminals  $\Sigma$
- finite set of productions  $P \subseteq T_{\Sigma}(Q) \times Q$   
(we write  $r \rightarrow q$  for productions  $(r, q)$ )

## Example productions

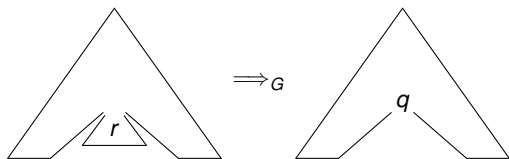


# Regular Tree Languages

## Derivation semantics and recognized tree language

Regular tree grammar  $G = (Q, \Sigma, I, P)$

- for each production  $r \rightarrow q \in P$



- generated tree language

$$L(G) = \{t \in T_{\Sigma}(\emptyset) \mid \exists q \in I: t \Rightarrow_G^* q\}$$

# Regular Tree Languages

Recall 3 equivalence classes (“no  $\sigma$ ”, “some  $\sigma$ , but legal”, illegal)

$$L' = \{t \in T_{\Sigma}(\emptyset) \mid \sigma \text{ never below } \delta\}$$

$$\mathcal{C}_1 = [\alpha]$$

$$\mathcal{C}_2 = [\sigma(\alpha, \alpha)]$$

$$\mathcal{C}_3 = [\delta(\sigma(\alpha, \alpha), \alpha)]$$

# Regular Tree Languages

Recall 3 equivalence classes (“no  $\sigma$ ”, “some  $\sigma$ , but legal”, illegal)

$$L' = \{t \in T_{\Sigma}(\emptyset) \mid \sigma \text{ never below } \delta\}$$

$$\mathcal{C}_1 = [\alpha]$$

$$\mathcal{C}_2 = [\sigma(\alpha, \alpha)]$$

$$\mathcal{C}_3 = [\delta(\sigma(\alpha, \alpha), \alpha)]$$

Productions with nonterminals  $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$

$$\alpha \rightarrow \mathcal{C}_1 \quad \delta(\mathcal{C}_1, \mathcal{C}_1) \rightarrow \mathcal{C}_1$$

$$\sigma(\mathcal{C}_1, \mathcal{C}_1) \rightarrow \mathcal{C}_2 \quad \sigma(\mathcal{C}_1, \mathcal{C}_2) \rightarrow \mathcal{C}_2 \quad \sigma(\mathcal{C}_2, \mathcal{C}_1) \rightarrow \mathcal{C}_2 \quad \sigma(\mathcal{C}_2, \mathcal{C}_2) \rightarrow \mathcal{C}_2$$

$$\delta(\mathcal{C}_1, \mathcal{C}_2) \rightarrow \mathcal{C}_3 \quad \delta(\mathcal{C}_1, \mathcal{C}_3) \rightarrow \mathcal{C}_3 \quad \delta(\mathcal{C}_2, \mathcal{C}_1) \rightarrow \mathcal{C}_3 \quad \delta(\mathcal{C}_2, \mathcal{C}_2) \rightarrow \mathcal{C}_3$$

$$\delta(\mathcal{C}_2, \mathcal{C}_3) \rightarrow \mathcal{C}_3 \quad \delta(\mathcal{C}_3, \mathcal{C}_1) \rightarrow \mathcal{C}_3 \quad \delta(\mathcal{C}_3, \mathcal{C}_2) \rightarrow \mathcal{C}_3 \quad \delta(\mathcal{C}_3, \mathcal{C}_3) \rightarrow \mathcal{C}_3$$

$$\sigma(\mathcal{C}_1, \mathcal{C}_3) \rightarrow \mathcal{C}_3 \quad \sigma(\mathcal{C}_2, \mathcal{C}_3) \rightarrow \mathcal{C}_3 \quad \sigma(\mathcal{C}_3, \mathcal{C}_1) \rightarrow \mathcal{C}_3 \quad \sigma(\mathcal{C}_3, \mathcal{C}_2) \rightarrow \mathcal{C}_3$$

$$\sigma(\mathcal{C}_3, \mathcal{C}_3) \rightarrow \mathcal{C}_3$$

# Regular Tree Languages

## Properties

- ✓ simple
- ✓ most expressive class we consider
- ✗ ambiguity, (several explanations for a generated tree)  
but can be removed
- ✓ closed under all Boolean operations  
(union/intersection/complement: ✓/✓/✓)
- ✓ all relevant properties decidable (emptiness, inclusion, ...)

## Characterizations

- finite index congruences
- regular tree grammars
- (deterministic) tree automata
- regular tree expressions
- second-order logic formulas
- ...

# Tree Languages

## Representations

- ~~enumerate trees~~
- local tree languages
- tree substitution languages
- regular tree languages



# Tree Languages

## Representations

- enumerate trees
- local tree languages
- tree substitution languages
- regular tree languages

Definition (Local tree grammar [Gécseg, Steinby 1984])

Local tree grammar = finite set of legal branchings  
(together with a set of root labels)

$$G = (\Sigma, I, P) \text{ with } I \subseteq \Sigma \text{ and } P \subseteq \bigcup_{k \in \mathbb{N}} \text{rk}^{-1}(k) \times \Sigma^k$$

# Local Tree Languages

Example (with root label S)

$$S \rightarrow NP_1 VP_2$$
$$NP_2 \rightarrow NP_2 PP$$
$$MD \rightarrow \text{must}$$
$$VP_2 \rightarrow MD VP_3$$
$$VP_3 \rightarrow VB PP NP_2$$

...

# Local Tree Languages

Example (with root label S)

$S \rightarrow NP_1 VP_2$

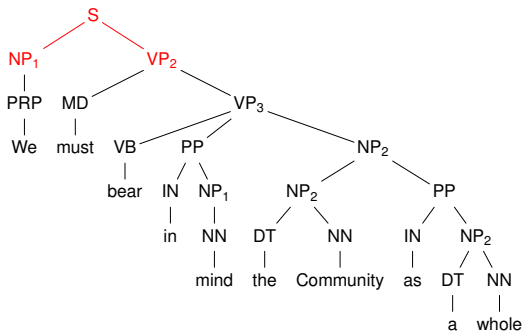
$NP_2 \rightarrow NP_2 PP$

$MD \rightarrow \text{must}$

$VP_2 \rightarrow MD VP_3$

$VP_3 \rightarrow VB PP NP_2$

...



# Local Tree Languages

Example (with root label S)

$S \rightarrow NP_1 VP_2$

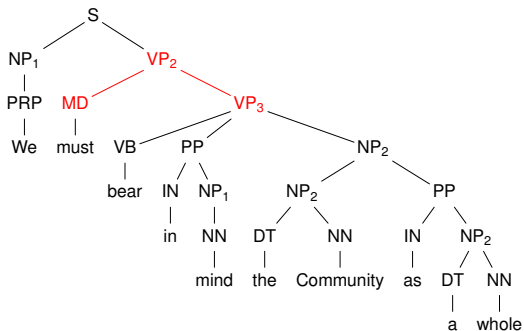
$NP_2 \rightarrow NP_2 PP$

$MD \rightarrow \text{must}$

$VP_2 \rightarrow MD VP_3$

$VP_3 \rightarrow VB PP NP_2$

...



# Local Tree Languages

Example (with root label S)

$S \rightarrow NP_1 VP_2$

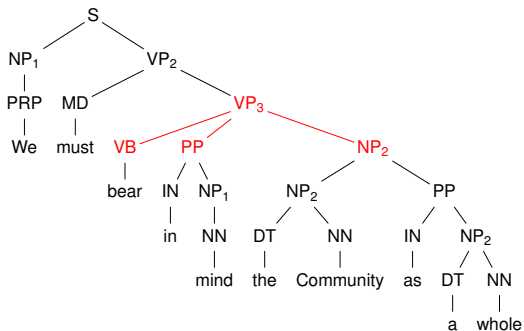
$NP_2 \rightarrow NP_2 PP$

$MD \rightarrow \text{must}$

$VP_2 \rightarrow MD VP_3$

$VP_3 \rightarrow VB PP NP_2$

...



# Local Tree Languages

Example (with root label S)

$S \rightarrow NP_1 VP_2$

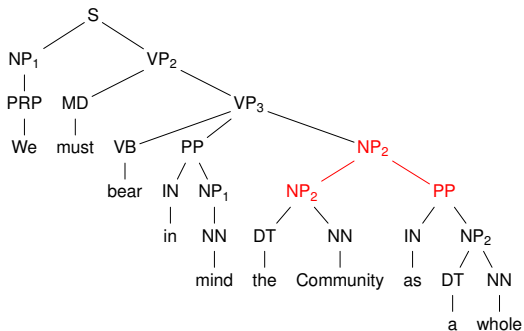
$NP_2 \rightarrow NP_2 PP$

$MD \rightarrow \text{must}$

$VP_2 \rightarrow MD VP_3$

$VP_3 \rightarrow VB PP NP_2$

...



# Local Tree Languages

Example (with root label S)

$S \rightarrow NP_1 VP_2$

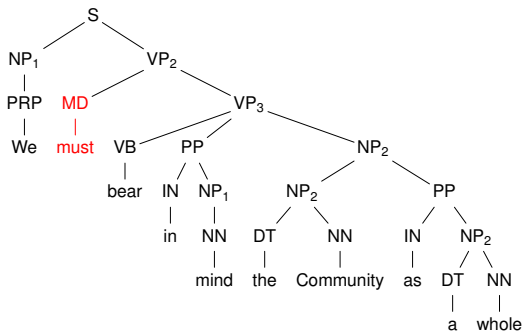
$NP_2 \rightarrow NP_2 PP$

$MD \rightarrow \text{must}$

$VP_2 \rightarrow MD VP_3$

$VP_3 \rightarrow VB PP NP_2$

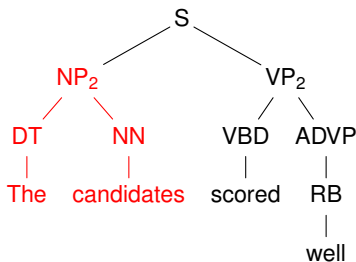
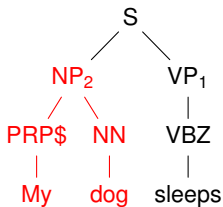
...



# Local Tree Languages

not closed under union

- these singletons are local



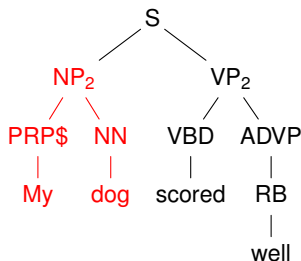
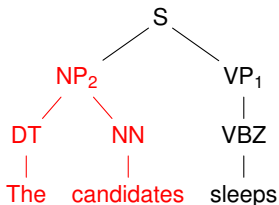
- but their union cannot be local



# Local Tree Languages

not closed under union

- these singletons are local

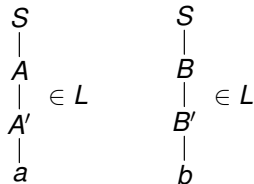


- but their union cannot be local  
(as we also generate these trees — overgeneralization)

# Local Tree Languages

not closed under complement

- this tree language  $L$  is local

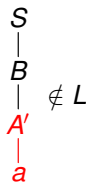
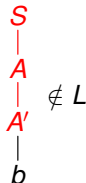
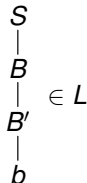
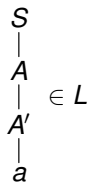


- but its complement cannot be local

# Local Tree Languages

not closed under complement

- this tree language  $L$  is local



- but its complement cannot be local  
(as we also generate these trees — overgeneralization)

# Local Tree Languages

## Properties

- ✓ simple
- ✓ no ambiguity (unique explanation for each recognized tree)
- ✗ not closed under Boolean operations  
(union/intersection/complement: ✗/✓/✗)
- ✗ not closed under (non-injective) relabelings
- ✓ locality of a regular tree language decidable

# Local Tree Languages

$$\begin{array}{ccccccc} \alpha \rightarrow \mathcal{C}_1 & \delta(\mathcal{C}_1, \mathcal{C}_1) \rightarrow \mathcal{C}_1 & & & & & \text{(irrelevant productions omitted)} \\ \sigma(\mathcal{C}_1, \mathcal{C}_1) \rightarrow \mathcal{C}_2 & \sigma(\mathcal{C}_1, \mathcal{C}_2) \rightarrow \mathcal{C}_2 & \sigma(\mathcal{C}_2, \mathcal{C}_1) \rightarrow \mathcal{C}_2 & \sigma(\mathcal{C}_2, \mathcal{C}_2) \rightarrow \mathcal{C}_2 & & & \end{array}$$

# Local Tree Languages

$$\begin{array}{llll} \alpha \rightarrow \mathcal{C}_1 & \delta(\mathcal{C}_1, \mathcal{C}_1) \rightarrow \mathcal{C}_1 & & \text{(irrelevant productions omitted)} \\ \sigma(\mathcal{C}_1, \mathcal{C}_1) \rightarrow \mathcal{C}_2 & \sigma(\mathcal{C}_1, \mathcal{C}_2) \rightarrow \mathcal{C}_2 & \sigma(\mathcal{C}_2, \mathcal{C}_1) \rightarrow \mathcal{C}_2 & \sigma(\mathcal{C}_2, \mathcal{C}_2) \rightarrow \mathcal{C}_2 \end{array}$$

- 1 extract all possible branches and root labels

$$\left\{ \begin{array}{l} \delta \rightarrow \alpha \alpha, \delta \rightarrow \alpha \delta, \delta \rightarrow \delta \alpha, \delta \rightarrow \delta \delta, \\ \sigma \rightarrow \alpha \alpha, \sigma \rightarrow \alpha \delta, \sigma \rightarrow \delta \alpha, \sigma \rightarrow \delta \delta, \\ \sigma \rightarrow \alpha \sigma, \sigma \rightarrow \delta \sigma, \sigma \rightarrow \sigma \alpha, \sigma \rightarrow \sigma \delta, \sigma \rightarrow \sigma \sigma \end{array} \right\}$$

# Local Tree Languages

$$\begin{array}{l} \alpha \rightarrow \mathcal{C}_1 \quad \delta(\mathcal{C}_1, \mathcal{C}_1) \rightarrow \mathcal{C}_1 \quad (\text{irrelevant productions omitted}) \\ \sigma(\mathcal{C}_1, \mathcal{C}_1) \rightarrow \mathcal{C}_2 \quad \sigma(\mathcal{C}_1, \mathcal{C}_2) \rightarrow \mathcal{C}_2 \quad \sigma(\mathcal{C}_2, \mathcal{C}_1) \rightarrow \mathcal{C}_2 \quad \sigma(\mathcal{C}_2, \mathcal{C}_2) \rightarrow \mathcal{C}_2 \end{array}$$

- 1 extract all possible branches and root labels

$$\begin{aligned} & \{ \delta \rightarrow \alpha \alpha, \delta \rightarrow \alpha \delta, \delta \rightarrow \delta \alpha, \delta \rightarrow \delta \delta, \\ & \quad \sigma \rightarrow \alpha \alpha, \sigma \rightarrow \alpha \delta, \sigma \rightarrow \delta \alpha, \sigma \rightarrow \delta \delta, \\ & \quad \sigma \rightarrow \alpha \sigma, \sigma \rightarrow \delta \sigma, \sigma \rightarrow \sigma \alpha, \sigma \rightarrow \sigma \delta, \sigma \rightarrow \sigma \sigma \} \end{aligned}$$

- 2 check whether this local tree grammar  $G$  overgeneralizes  
(check whether  $L(G) \subseteq L$ )

# Local Tree Languages

## Characterizations

- local tree grammars
- parse trees of context-free grammars
- (not much available, but seems well understood)



# Tree Languages

## Representations

- ~~enumerate trees~~
- local tree languages
- **tree substitution languages**
- regular tree languages

# Tree Languages

## Representations

- enumerate trees
- local tree languages
- tree substitution languages
- regular tree languages

Definition (Tree substitution grammar [Joshi, Schabes 1997])

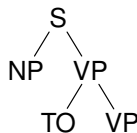
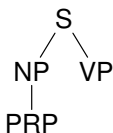
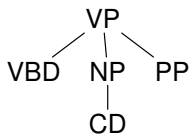
Tree substitution grammar = finite set of legal fragments  
(together with a set of root labels)

$G = (\Sigma, I, P)$  with  $I \subseteq \Sigma$  and finite  $P \subseteq T_{\Sigma}(\Sigma)$

# Tree Substitution Languages

## Typical fragments

[Post 2011]



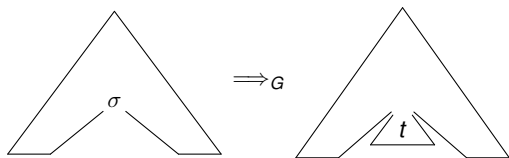
## Derivation step     $\xi \Rightarrow_G \zeta$

- $\xi = c[\text{root}(t)]$  and  $\zeta = c[t]$  for some context  $c$  and fragment  $t \in P$

# Tree Substitution Languages

Tree substitution grammar  $G = (\Sigma, l, P)$

- for each fragment  $t \in P$  with root label  $\sigma$



- generated tree language

$$L(G) = \{t \in T_{\Sigma}(\emptyset) \mid \exists \sigma \in l: \sigma \Rightarrow_G^* t\}$$

# Tree Substitution Languages

## Fragments

$S(NP_1(\text{PRP}), VP_2)$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{PRP}(\text{We})$

$\text{MD}(\text{must})$

## Derivation

S

# Tree Substitution Languages

## Fragments

$S(NP_1(PRP), VP_2)$

$VP_2(MD, VP_3(VB, PP, NP_2))$

PRP(We)

MD(must)

## Derivation

S

# Tree Substitution Languages

## Fragments

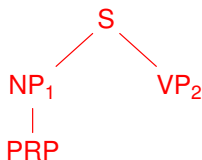
$S(NP_1(PRP), VP_2)$

$VP_2(MD, VP_3(VB, PP, NP_2))$

PRP(We)

MD(must)

## Derivation



# Tree Substitution Languages

## Fragments

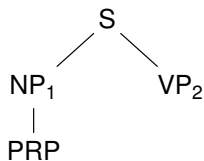
$S(NP_1(\text{PRP}), VP_2)$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{PRP}(\text{We})$

$\text{MD}(\text{must})$

## Derivation





# Tree Substitution Languages

## Fragments

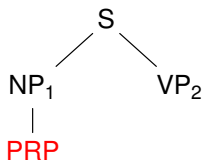
$S(NP_1(PRP), VP_2)$

$VP_2(MD, VP_3(VB, PP, NP_2))$

PRP(We)

MD(must)

## Derivation



# Tree Substitution Languages

## Fragments

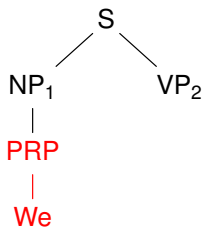
$S(NP_1(\text{PRP}), VP_2)$

$\text{PRP}(\text{We})$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{MD}(\text{must})$

## Derivation



# Tree Substitution Languages

## Fragments

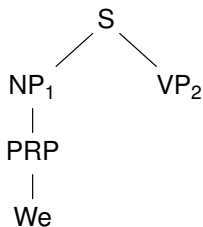
$S(NP_1(\text{PRP}), VP_2)$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{PRP}(\text{We})$

$\text{MD}(\text{must})$

## Derivation



# Tree Substitution Languages

## Fragments

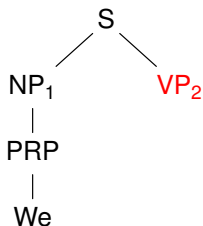
$S(NP_1(\text{PRP}), VP_2)$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{PRP}(\text{We})$

$\text{MD}(\text{must})$

## Derivation



# Tree Substitution Languages

## Fragments

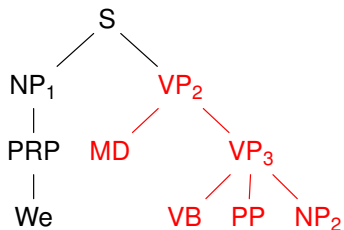
$S(NP_1(\text{PRP}), VP_2)$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{PRP}(\text{We})$

$\text{MD}(\text{must})$

## Derivation



# Tree Substitution Languages

## Fragments

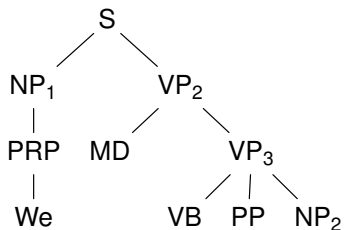
$S(NP_1(\text{PRP}), VP_2)$

$\text{PRP}(\text{We})$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{MD}(\text{must})$

## Derivation



# Tree Substitution Languages

## Fragments

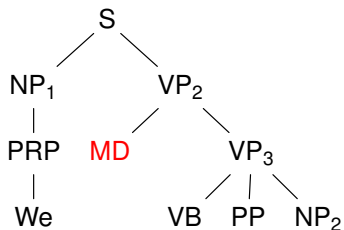
$S(NP_1(\text{PRP}), VP_2)$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{PRP}(\text{We})$

$\text{MD}(\text{must})$

## Derivation



# Tree Substitution Languages

## Fragments

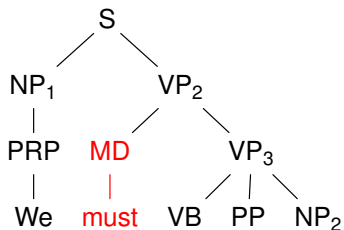
$S(NP_1(PRP), VP_2)$

$VP_2(MD, VP_3(VB, PP, NP_2))$

$PRP(We)$

$MD(must)$

## Derivation





# Tree Substitution Languages

## Fragments

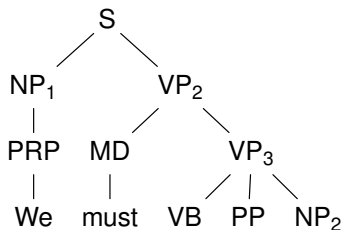
$S(NP_1(\text{PRP}), VP_2)$

$\text{PRP}(\text{We})$

$VP_2(\text{MD}, VP_3(\text{VB}, \text{PP}, NP_2))$

$\text{MD}(\text{must})$

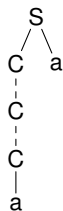
## Derivation



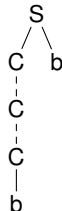
# Tree Substitution Languages

not closed under union

- these languages are tree substitution languages individually



$$L_1 = \{S(C^n(a), a) \mid n \in \mathbb{N}\}$$



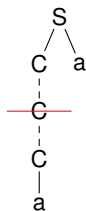
$$L_2 = \{S(C^n(b), b) \mid n \in \mathbb{N}\}$$

- but their union is not

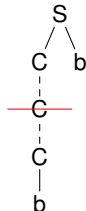
# Tree Substitution Languages

not closed under union

- these languages are tree substitution languages individually



$$L_1 = \{S(C^n(a), a) \mid n \in \mathbb{N}\}$$



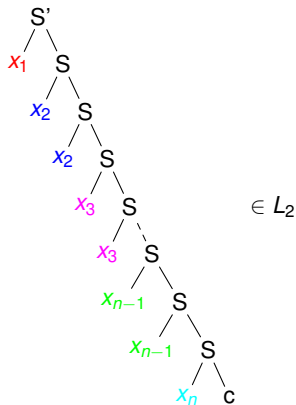
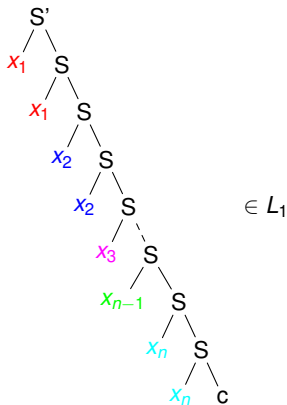
$$L_2 = \{S(C^n(b), b) \mid n \in \mathbb{N}\}$$

- but their union is not  
(exchange subtrees below the indicated cuts)

# Tree Substitution Languages

## not closed under intersection

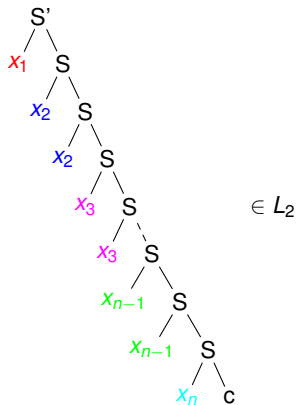
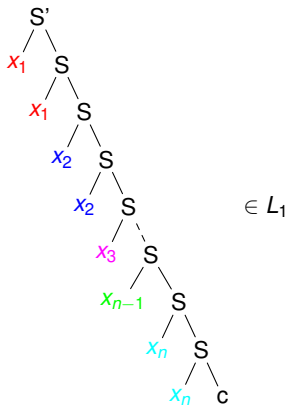
- these languages  $L_1$  and  $L_2$  are tree substitution languages individually for  $n \geq 1$  and arbitrary  $x_1, \dots, x_n \in \{a, b\}$



# Tree Substitution Languages

## not closed under intersection

- these languages  $L_1$  and  $L_2$  are tree substitution languages individually for  $n \geq 1$  and arbitrary  $x_1, \dots, x_n \in \{a, b\}$

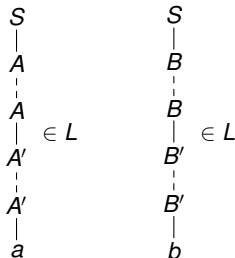


- but their intersection only contains trees with  $x_1 = x_2 = \dots = x_n$  and is not a tree substitution language

# Tree Substitution Languages

not closed under complement

- this language  $L$  is a tree substitution language

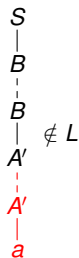
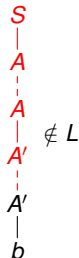
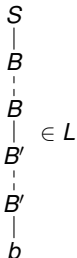
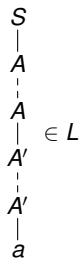


- but its complement is not

# Tree Substitution Languages

not closed under complement

- this language  $L$  is a tree substitution language



- but its complement is not  
(exchange as indicated in red)

# Tree Substitution Languages

## Properties

- ✓ simple
- ✓ contain all finite and co-finite tree languages
- ✗ ambiguity (several explanations for a generated tree)
- ✗ not closed under Boolean operations  
(union/intersection/complement: ✗/✗/✗)
- ✓ can express many finite-distance dependencies  
(extended domain of locality)



# Tree Substitution Languages

## Characterizations

- tree substitution grammars
- ???

(generally badly understood)

# Tree Substitution Languages

## Characterizations

- tree substitution grammars

- ???

(generally badly understood)

## Remark:

- several unions lead to additional power

# Tree Substitution Languages

## Open questions

- multiple intersections more expressive?
- which regular tree languages are tree substitution languages?
- relation to local tree languages?

# Tree Substitution Languages

## Open questions

- multiple intersections more expressive?
- which regular tree languages are tree substitution languages?
- relation to local tree languages?
- extension to weights
- application to parsing

# Tree Substitution Languages

## Open questions

- multiple intersections more expressive?
- which regular tree languages are tree substitution languages?
- relation to local tree languages?
- extension to weights
- application to parsing

Thank you for your attention!

# Tree Substitution Languages

Experiment

[Post, Gildea 2009]

grammar	size	Prec.	Recall	$F_1$
local	46k	75.37	70.05	72.61
“spinal” TSG	190k	80.30	78.10	79.18
“minimal subset” TSG	2,560k	76.40	78.29	77.33

(on WSJ Sect. 23)

# Tree Substitution Languages with Latent Variables

Experiment

[Shindo et al. 2012]

grammar	F1 score	
	$ w  \leq 40$	full
TSG [Post, Gildea, 2009]	82.6	
TSG [Cohn et al., 2010]	85.4	84.7
CFGlv [Collins, 1999]	88.6	88.2
CFGlv [Petrov, Klein, 2007]	90.6	90.1
CFGlv [Petrov, 2010]		91.8
TSGlv (single)	91.6	91.1
TSGlv (multiple)	92.9	92.4
Discriminative Parsers		
Carreras et al., 2008		91.1
Charniak, Johnson, 2005	92.0	91.4
Huang, 2008	92.3	91.7