

Grammar Implementation with TAG

Motivation for TAG

Timm Lichte

HHU Düsseldorf

SS 2011

13.04.2010

1. Why CFG is not enough
2. Tree Substitution Grammars
3. Tree Adjoining Grammars
 - Adjunction and substitution
 - Adjunction constraints

Why CFG is not enough

... for treating natural language:

1. only atomic non-terminals
2. only weak lexicalization (lexicalization challenge)
3. expressive power is too low (expressivity challenge)

Why CFG is not enough (1) - Atomic non-terminals

$S \rightarrow NP VP$ $NP \rightarrow John$ $NP \rightarrow Mary$
 $VP \rightarrow V$ $VP \rightarrow V NP$ $V \rightarrow sleeps$ $V \rightarrow likes$

Possible derivation:

$S \Rightarrow NP VP \Rightarrow John VP \Rightarrow John V \Rightarrow John sleeps$

$S \xRightarrow{*} John likes Mary$

$S \xRightarrow{*} John sleeps Mary$

How to treat subcategorization frames, number agreement, and case marking?

- (1) a. Kim depends on Sandy.
 *Kim depends Sandy.
 *Kim depends.
 b. *The children depends on Sandy.
 c. Kim depends on her/*she.

Why CFG is not enough (1)

How to treat subcategorization frames, number agreement, and case marking?

⇒ encode the necessary information into the non-terminal symbols

$S \rightarrow NP_{3sg/nom} VP_{3sg/itr}$	$S \rightarrow NP_{3sg/nom} VP_{3sg/tr}$
$VP_{3sg/itr} \rightarrow V_{3sg/itr}$	$VP_{3sg/tr} \rightarrow V_{3sg/tr} NP_{3sg/acc}$
$NP_{3sg/nom} \rightarrow John$	$NP_{3sg/acc} \rightarrow Mary$
$V_{3sg/itr} \rightarrow sleeps$	$V_{3sg/tr} \rightarrow likes$

$S \xrightarrow{*} John\ likes\ Mary$

$S \xrightarrow{*} John\ sleeps$

Drawback: Every possible combination of subcategorization frame, number agreement, and case marking necessitates its own rule (let alone the number of non-terminal symbols).

Why CFG is not enough (1)

Example from German: $NP \rightarrow D N$ (determiner noun pairs)

Müller(2007) presents a CFG with 48 non-terminal symbols and 24 rules!

$NP_{3sg/nom} \rightarrow D_{fem/sg/nom} N_{fem/sg/nom}$
 $NP_{3sg/nom} \rightarrow D_{masc/sg/nom} N_{masc/sg/nom}$
 $NP_{3sg/nom} \rightarrow D_{neu/sg/nom} N_{neu/sg/nom}$
 $NP_{3pl/nom} \rightarrow D_{fem/pl/nom} N_{fem/pl/nom}$
 $NP_{3pl/nom} \rightarrow D_{masc/pl/nom} N_{masc/pl/nom}$
 $NP_{3pl/nom} \rightarrow D_{neu/pl/nom} N_{neu/pl/nom}$
...

\Rightarrow grammar writing is tedious and error prone

\Rightarrow generalizations are hardly expressible

Remedy: feature structures instead of atomic non-terminal symbols, unification, underspecification

Lexicalization

In a lexicalized grammar, each element of the grammar contains at least one lexical item (terminal symbol).

$G_1: S \rightarrow SS, S \rightarrow a$

$G_2: S \rightarrow aS, S \rightarrow a$

- **Computationally interesting:** the number of analyses for a sentence is finite (if the grammar is finite of course).
- **Linguistically interesting:** each lexical item comes with the possibility of certain partial syntactic constructions, therefore one would like to associate it to a set of substructures.

Why CFG is not enough (2)

Lexicalizing a CFG

- Greibach normal form: $A \rightarrow aB_1...B_k$ ($k \geq 0$)
- weak lexicalization: string language is preserved
- strong lexicalization: tree structure is preserved

Question: can CFGs be lexicalized such that the set of trees remains the same (strong lexicalization)?

Answer: No. Only weak lexicalization (same string language).

$G_1: S \rightarrow SS, S \rightarrow a$

$G_2: S \rightarrow aS, S \rightarrow a$

G_1 cannot be strongly lexicalized with some finite CFG, e.g. G_2 .

Why CFG is not enough (3) - Low expressive power

Question: Are CFGs powerful enough to describe all natural language phenomena?

Answer: No.

Example: **cross-serial dependencies** in Dutch and in Swiss German

(1)

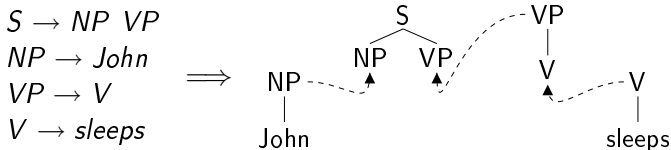
... dat Wim Jan Marie de kinderen zag helpen leren zwemmen
... that Wim Jan Marie the children saw help teach swim
'... that Wim saw Jan help Marie teach the children to swim'

A formalism that can generate cross-serial dependencies must be able to generate the copy language $\{ww \mid w \in \{a, b\}^*\}$.

But: The **copy language** is **not context-free**.

Tree Substitution Grammar (TSG)

- a tree rewriting version of CFG
 - weakly equivalent generative capacity
- A CFG-production corresponds to a TSG-tree with the LHS as root and the RHS as daughters.
- Applying a CFG-production corresponds to **substituting** a non-terminal leaf for a new tree.



Tree Substitution Grammar (2)

A **Tree Substitution Grammar (TSG)** is a triple $G = \langle N, T, I \rangle$ such that

- T and N are disjoint alphabets, the terminals and nonterminals, and
- I is a finite set of **initial** trees.

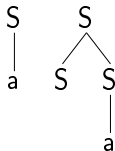
The trees can be combined into larger trees by **substitution**.

The **tree language** of a TSG is the set of trees generated in this way that do not contain any remaining non-terminal leaves.

Tree Substitution Grammar (3)

Some important facts:

- TSG is weakly equivalent to CFG (same string language).
- TSG is not powerful enough to describe cross-serial dependencies.
- It is not possible to find a strongly equivalent (same trees) lexicalized TSG for each CFG.

$$S \rightarrow SS$$
$$S \rightarrow a$$


⇒ Solution: **adjunction operation** and **adjunction constraints**!

Tree Adjoining Grammar (TAG)

TAG = TSG + adjunction + adjunction constraints

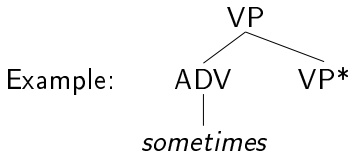
- The definition of TAG goes back to Joshi et al. (1975).
- TAG is among the most frequently used grammar formalisms in computational linguistics.
- TAG is interesting both for its computational properties (**mildly context-sensitivity**) and for its linguistic applications.
- There are large coverage TAG grammars for English (XTAG, Philadelphia) and French (FTAG, Paris).

Tree Adjoining Grammar - Adjunction (1)

Rewriting operations:

- **substitution**: replacing a leaf with a new tree.
- **adjunction**: replacing an internal node with a new tree.

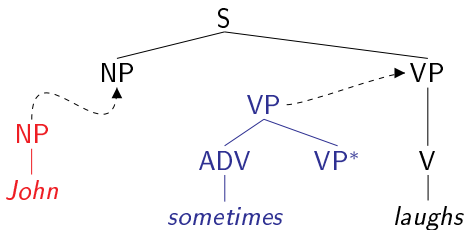
Trees that may adjoin are called **auxiliary trees** and have a special leaf, the **footnode** (marked by *). After adjunction, the subtree below the target node appears below the footnode.



The root node and the footnode are required to carry the same label. The path from the root node to the footnode is called the **spine**.

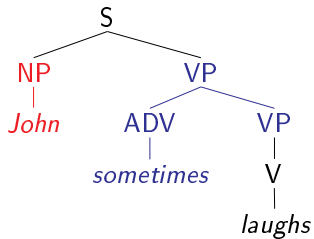
Tree Adjoining Grammar - Adjunction (2)

(2) John sometimes laughs



derived tree

laugh[1, *john*][2, *sometimes*]:



Tree Adjoining Grammar - Adjunction (3)

A **Tree Adjoining Grammar (TAG)** is a quadruple

$G = \langle N, T, I, A \rangle$ such that

- T and N are disjoint alphabets, the terminals and nonterminals,
- I is a finite set of **initial** trees, and
- A is a finite set of **auxiliary** trees.

The trees in $I \cup A$ are called **elementary** trees.

G is **lexicalized** iff each elementary tree has at least one leaf with a terminal label (LTAG).

Tree Adjoining Grammar - Adjunction (4)

A derivation starts with an initial tree.

In a final derived tree, all leaves must have terminal labels:

Let $G = \langle I, A, N, T \rangle$ be a TAG. Let γ and γ' be finite trees.

- $\gamma \Rightarrow \gamma'$ in G iff there is a node position p and an instance γ_0' of a tree (possibly derived from some) $\gamma_0 \in I \cup A$ such that $\gamma' = \gamma[p, \gamma_0']$.
 \Rightarrow^* is the reflexive transitive closure of \Rightarrow .
- The tree language of G is $L_T(G) := \{\gamma \mid \text{there is an } \alpha \in I \text{ such that } \alpha \xRightarrow{*} \gamma \text{ and all leaves in } \gamma \text{ have terminal labels}\}$.

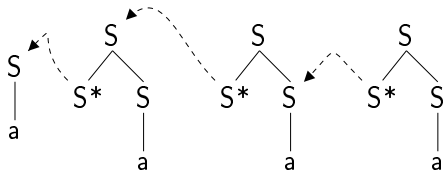
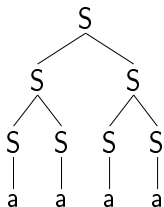
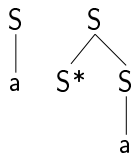
Tree Adjoining Grammar - Lexicalization challenge

LTAGs strongly lexicalize CFGs and TAGs.

Example:

$$\begin{aligned} S &\rightarrow SS \\ S &\rightarrow a \end{aligned}$$

is strongly equivalent with



Tree Adjoining Grammar - Adjunction constraints (1)

TAG as defined above are more powerful than CFG but they cannot generate the copy language.

In order to increase the expressive power, adjunction constraints are introduced that specify for each node

- 1 whether adjunction is mandatory and
- 2 which trees can be adjoined.

A **TAG with adjunction constraints** is a tuple $\langle N, T, I, A, O, C \rangle$ such that

- $\langle N, T, I, A \rangle$ is a TAG,
- $O : \{\mu \mid \mu \text{ is a node in a tree in } I \cup A\} \rightarrow \{1, 0\}$ is a function, and
- $C : \{\mu \mid \mu \text{ is a node in a tree in } I \cup A\} \rightarrow P(A)$ is a function.

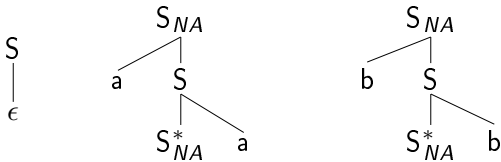
Three types of constraints are distinguished:

- **Obligatory Adjunction (OA):**
a node μ with $O(\mu) = 1$
- **Null Adjunction (NA):**
a node μ with $O(\mu) = 0$ and $C(\mu) = \emptyset$
- **Selective Adjunction (SA):**
a node μ with $O(\mu) = 0$ and $C(\mu) \neq \emptyset$ and $C(\mu) \neq A$

It is common practice to let the leaves carry the NA-constraint.

Tree Adjoining Grammar - Expressivity challenge

TAG for the copy language $\{ww \mid w \in \{a, b\}^*\}$:



Starting point: can we describe natural languages with CFGs?

- CFGs: string rewriting formalism, no strong lexicalization, no cross-serial dependencies.
- TSGs: tree rewriting formalism, no strong lexicalization, no cross-serial dependencies.
- TAG = TSG + adjunction + adjunction constraints
 - strong lexicalization
 - cross-serial dependencies