

# Parsing Beyond Context-Free Grammars: TAG: Grammar induction

Laura Kallmeyer

Heinrich-Heine-Universität Düsseldorf

slides partly done by Tatiana Bladier

Wintersemester 2021

# Overview

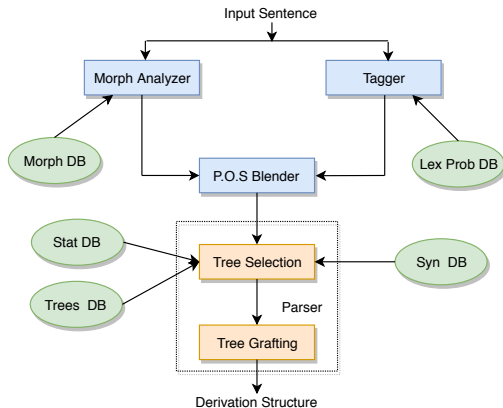
- 1 Hand-crafted TAGs
- 2 Annotated text corpora
- 3 TAG induction: Top-down approach
- 4 TAG induction: Bottom-up approach

# TAG induction: motivation

- A wide-coverage TAG for a natural language contains several thousands of trees and takes a lot of effort to build
- Examples: hand-crafted TAGs for English (XTAG project [G<sup>+</sup>98]) or French (FTAG)
- Growing amount of large-scale constituency treebanks for different languages (English, German, French, Spanish, Vietnamese, Korean etc.)
- Idea: automatic induction of linguistically plausible TAGs from syntactically annotated corpora
- Orientation on already existing hand-crafted LTAGs

# XTAG: LTAG implementation for English

- large coverage LTAG for English
- implemented (2001) by the XTAG Research Group at the University of Pennsylvania
- architecture [G<sup>+</sup>98]:



# XTAG: facts

- Morphological Analyzer and MorphDB
  - $\approx 317000$  inflected items derived from over 90000 stems
  - returns: root form + POS + inflectional information
- POS Tagger and Lex Prob DB
  - Wall-Street Journal trained trigram tagger
  - extended to output N-best POS sequences
- Syntactic DB
  - more than 30.000 entries, each consisting of:
    - ★ uninflected form of the word + POS
    - ★ list of trees or tree-families associated with the word
    - ★ list of feature equations that capture lexical idiosyncrasies

[G<sup>+</sup>98]

# XTAG: facts

- Tree DB
  - 1004 unanchored trees divided into 53 families and 221 individual trees
- Tree Selection
  - for each word, tree templates are chosen from the Tree Database and the anchor position is filled with the word
- Tree Grafting
  - once a particular lexicalized tree set is chosen for a sentence, parsing is done  $\implies$  output: derived (parse) tree and derivation tree

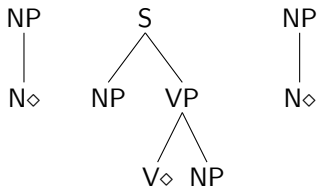
[G<sup>+</sup>98]

# XTAG: Lexical anchoring

Example

(1) John likes apples

templates:



lemmas

John

like

apple

words

John

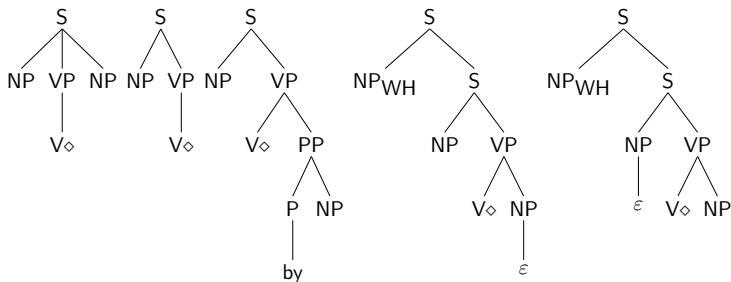
likes

apples

# XTAG: tree families

- XTAG *tree families*: grouping together (unanchored) trees that belong to the same subcategorization frame

Example: Family  $T_{n \times 0V_{n \times 1}}$  for transitive verbs (selected trees):



◇ marks the position of the lexical anchor



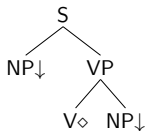
# XTAG: tree families

- tree families for verb classes:
  - transitive verbs, e.g. *buy* [<sub>NP</sub>]  
*Joe bought a book.*
  - ditransitive verbs, e.g. *buy* [<sub>NP</sub> <sub>NP</sub>]  
*Joe bought Rose a book.*
  - ditransitive verbs with PP, e.g. *put* [<sub>NP</sub> <sub>PP</sub>]  
*Joe put the book on the shelf.*
  - verbs with sentential complement, e.g. *tell* [<sub>NP</sub> <sub>S</sub>]  
*Joe told Rose that Tom lives in Amsterdam.*
  - light verbs, e.g. 'make comment'  
*Joe made a comment on my book.*
  - ...

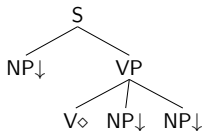
# XTAG: tree families

- transitive / ditransitive / light verbs

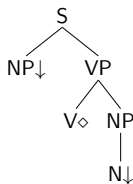
$\alpha NX0VNX1$



$\alpha NX0VNX2NX1$



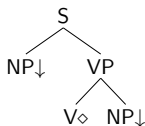
$\alpha NX0LVN1$



- X◇: lexical anchor
- lexical insertion from the lexicon

# XTAG: tree families

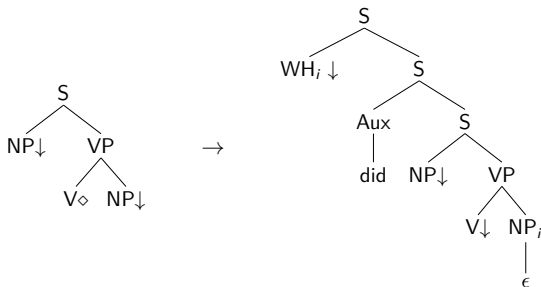
- tree family of *buy* [ $\_NP$ ] for
  - John bought a book.
 the 'base tree' ( $\alpha NX0VNX1$ )



- *Who bought a book? / What does John buy?* (wh-extraction)
  - *A book was bought by John.* (passivization)
  - *A book, Joe bought.* (topicalization)
- elementary trees derived from the base tree by transformation
- deriving *tree templates*

# XTAG: tree families

- elementary trees derived from the base tree by transformation
- e.g. base tree  $\rightarrow$  object extraction (wh-extraction)



# XTAG: tree families

- transform base trees while changing the properties of particular words → derives a different tree family
- for example: ergative verbs as *melt*  
*The sun melted the ice.* → *The ice melted.*
- XTAG covers the following phenomena:
  - wh-extraction, topicalization
  - relative clauses, raising and small clause constructions
  - clausal adjuncts and imperatives
  - it-clefts, wh-clefts
  - auxiliaries, copula, infinitives, gerunds, passives, adjuncts
  - PRO constructions, noun-noun modifications, determiner sequences
  - genitives, negation, noun-verb contractions

# Syntactically annotated text corpora: treebanks

- Hand-annotated (i.e. parsed and manually corrected) corpora are used for supervised and semi-supervised training of algorithms
- Constituency based corpora: Penn Treebank (English), NeGra German Corpus, TiGer treebank, French Treebank, Korean Treebank etc.
- Different formats: bracketed sentences, XML, NeGra

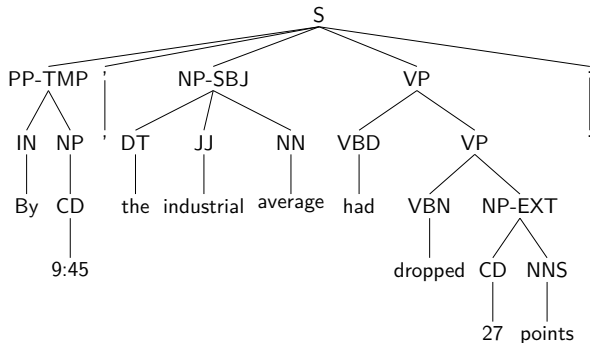
# Penn Treebank

- American English, syntactically annotated corpus, constituency treebank
- Created between 1989 and 1996
- Several genres, including IBM computer manuals, nursing notes, Wall Street Journal (WSJ) articles, and transcribed telephone conversations.
- The WSJ sections constitute the most used corpus for evaluating syntactic parsing.
- Standard split of the WSJ sections: train (sec. 0–18, 38.219 sentences, 912.344 tokens), dev (sec. 19–21, 5.527 sentences, 131.768 tokens), and test (sec. 22–24, 5.462 sentences, 129.654 tokens)

[MSM93]

# Penn Treebank

Sample tree:



Corresponding bracketed notation:

(S (PP-TMP (IN By) (NP (CD 9:45))) ( , ) (NP-SBJ (DT the) (JJ industrial) (NN average))) (VP (VBD had) (VP (VBN dropped) (NP-EXT (CD 27) (NNS points)))) ( . ))



# TAG induction: idea

- Sentences in treebanks are seen as derived TAG trees
- These derived TAG trees are cut in smaller pieces which correspond to elementary trees in TAG
- Treebanks provide additional information which can be used as features
- While extracting elementary trees we can also learn probabilities about trees (for example, statistics about frequencies of elementary trees)
- Probabilistic TAGs
- **Reversibility**

[XPJ00]

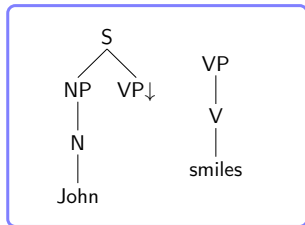
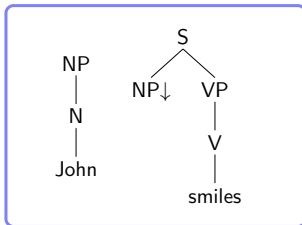
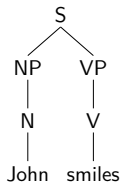
## TAG induction: idea

- Annotated trees can be cut to elementary trees in different ways.
- The choice of how the elementary trees should look like depends on the grammar theory. If there is a hand-crafted TAG, it can be used as orientation.
- There are certain linguistic principles to follow:
  - Dependencies (including long-distance dependencies) should be localized in the same elementary tree (i.e. head and its complements are included into elementary tree)
  - Recursion is factored into separate auxiliary trees (i.e. modifiers are excluded)
  - No elementary tree is anchored only by an empty element.

# TAG induction: idea

(2) John smiles

Two options for cutting the tree for (2) into elementary trees:



[XPJ00]

TAG induction from the PTB is usually inspired by the choices in the XTAG grammar.

# LTAG extraction algorithm

Conversion of a treebank tree into elementary trees:

- ① Mark children in the tree as being head, complement or modifier using heuristic percolation tables
- ② Decompose tree into elementary trees
  - elementary trees contain the corresponding lexical anchor and they represent a particular syntactic construction with slots for its complements
  - the corresponding unanchored elementary trees are called supertags

Further step: create derivation trees based on the node addresses.

[XP06]

# Percolation tables for heads and modifiers

- Percolation tables are used to mark the roles of the children in the tree
  - ★ Head child
  - ★ Syntactically obligatory child (a complement)
  - ★ Modifier
- Needed to extract linguistically plausible LTAGs
- Deterministic extraction procedure (only one LTAG is possible)
- Language dependent and treebank dependent
- A head can have several arguments, and its projections can be modified by other components.
  - ★ e.g. a verbal phrase (VP) can have a verb (V) as its head and nominal phrase (NP) as its argument
  - ★ a VP can be modified by adverbial phrases (APs), prepositional phrases (PPs) etc.

[XP06]

## Percolation tables for heads and modifiers

- Head rules for PTB (WSJ) based on Collins (1999 p. 240) (the list is not exhaustive)

parent node	search direction	head candidates
adjp	left-to-right	nns qp nn \$ advp jj vbn vbg adjp jjr np jjs dt
conjp	right-to-left	cc rb in
nac	left-to-right	nn nns nnp nnps np nac ex \$ cd qp prp vbg jj jjs jjr
vp	left-to-right	to vbd vbn md vbz vb vbg vbp vp adjp nn nns np

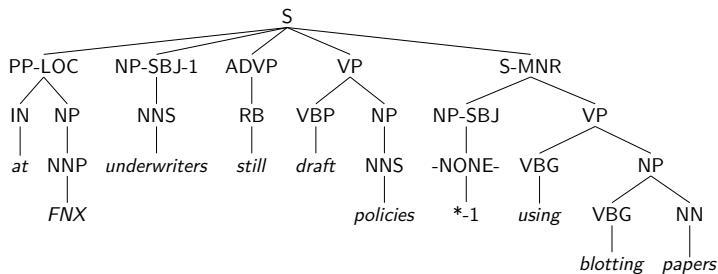
If no head can be found, the leftmost child is chosen as a default.

- Modifier rules for PTB (WSJ) (excerpt)

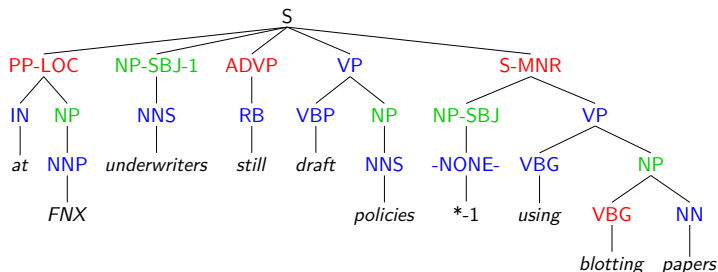
parent node	modifier nodes
s	pp pp-loc adv advp s-mnr
np	cc vbg jj jjs adjp nn np
vp	np jj jjs adjp pp-loc s-mnr

- All other nodes (non-leaves!) are marked as complements

# Step 1a): Head and modifier identification



# Step 1a): Head and modifier identification



Marking heads, modifiers and complements.



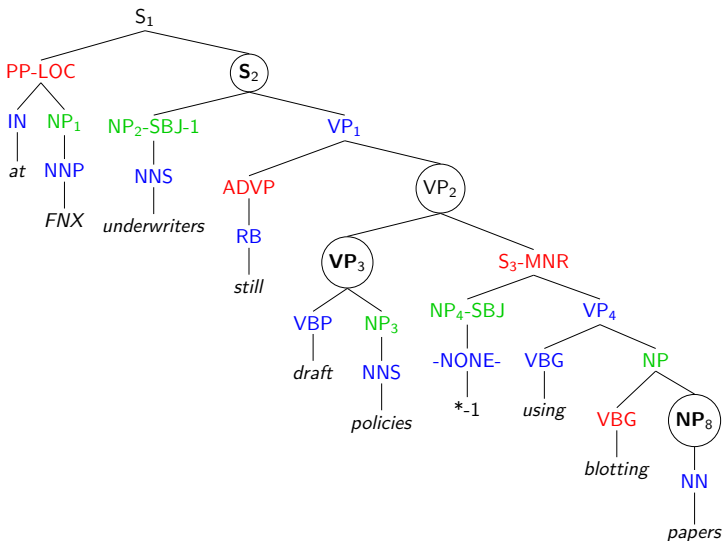
## Step 1b): Bracketing around modifiers

For each subtree  $v \rightarrow v_1 \dots v_n$  in the original tree, we add further bracketings whenever we encounter a modifier.

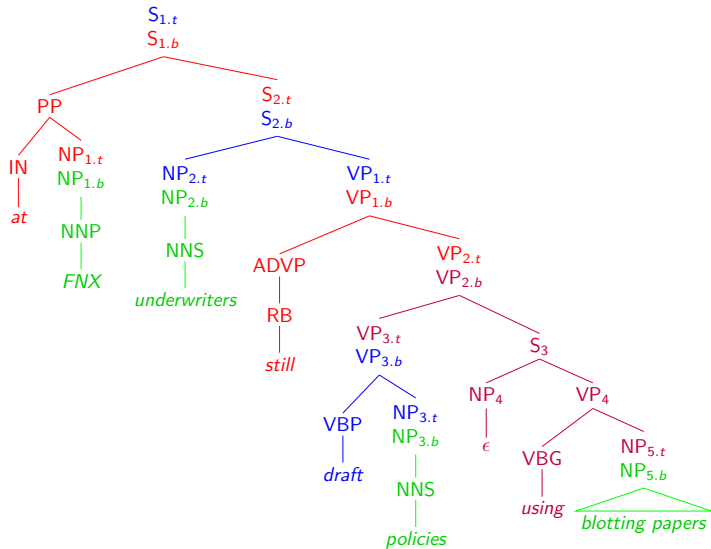
- ① (i) For each subtree  $v \rightarrow v_1 \dots v_n$ :  
if  $v_m$  is the leftmost modifier,  $v_h$  is the head, and  $1 < m < h$ : take a new node  $v'$ , replace our subtree with  $v \rightarrow v_1 \dots v_{m-1} v'$  ( $v'$  is a head) and add the subtree  $v' \rightarrow v_m \dots v_n$ .  $I(v') = I(v_h)$ .
  - (ii) For each subtree  $v \rightarrow v_1 \dots v_n$ :  
if  $v_1$  is the leftmost modifier and either  $n > 2$  or  $n = 2$  and  $I(v) \neq I(v_n)$ : take a new node  $v'$ , replace our subtree with  $v \rightarrow v_1 v'$  ( $v'$  is a head) and add the subtree  $v' \rightarrow v_2 \dots v_n$ .  $I(v') = I(v)$ .
- ② Repeat 1. until it cannot be applied any longer.

Then do the same (mirrored) for the modifiers to the right of the head.

# Step 1b): Bracketing around modifiers

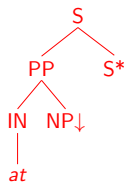


## Step 2: Extracting elementary trees



## Step 2: Extracting elementary trees

#1



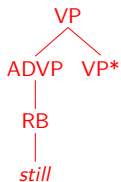
#2



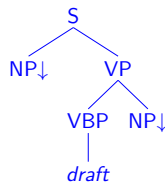
#3



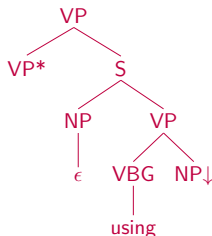
#4



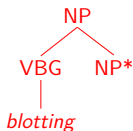
#5



#7



#8



#9



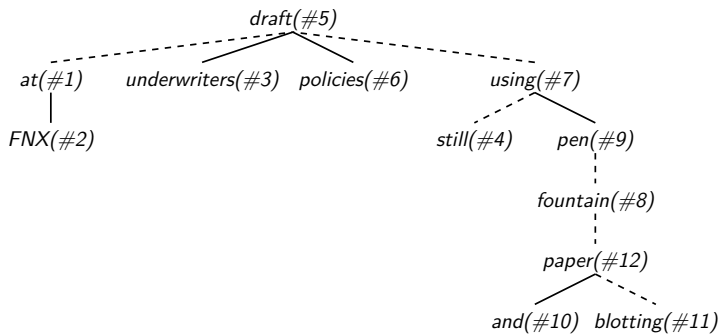
#6



## Step 3: filter out linguistically implausible trees

- Linguistically implausible trees are :
  - trees with erroneous labels (e.g. an adjective erroneously marked as a preposition) ,
  - trees with an erroneously marked head child (e.g. an adjective marked as a head of an NP)
  - other trees which do not meet the criteria of the design for extracted TAG
- Filtering out implausible trees is not a trivial task
- One possibility to rule out implausible trees is to first decompose a template into a set of sub-templates then mark the template as plausible if and only if every sub-template is plausible.
- Another possibility is to manually delete the trees from the output (if possible)

## Step 4: Creating derivation trees (optional)



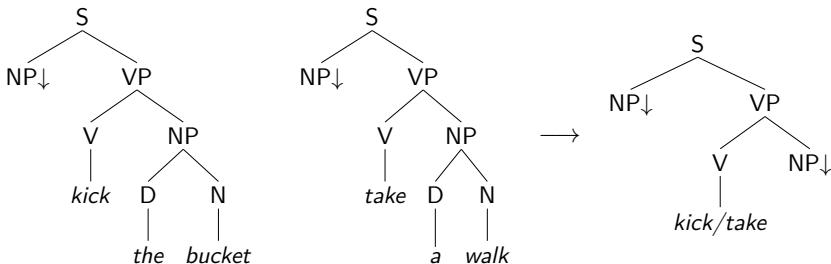
## Some considerations

- It might be useful to reduce the set of the part-of-speech tags and functional tags in the treebank in order to reduce the number of elementary trees in the extracted LTAG

	tags in the PTB	merged tags
adjectives	JJ/JJR/JJS	A
adverbs	RB/RBR/RBS/WRB	Adv
determiners	DT/PDT/WDT/PRP\$/WP\$	D
nouns	CD/NN/NNS/NNP/NNPS/PRP/WP/EX/\$/#	N
verbs	V/MD/VB/VBP/VBZ/VBN/VBD/VBG/TO	V
clauses	S/SQ/SBAR/SBARQ/SINV	S
noun phrases	NAC/NP/NX/QP/WHNP	NP
adjective phrases	ADJP/WHADJP	AP
adverbial phrases	ADVP/WHADVP	AdvP
preposition phrases	PP/WHPP	PP

## Some considerations: Co-anchors

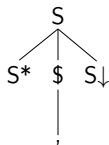
- Although TAG elementary trees typically have one anchor, there are situations where more than one anchor could be desirable (for example, idiomatic expressions or multi-word expressions)





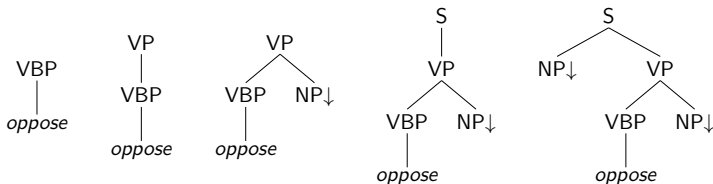
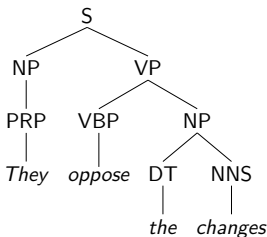
## Some considerations: Punctuation

- Paired punctuation marks (e.g. quotation marks) vs. not-paired
- Punctuation marks are sometimes left out in hand-crafted LTAGs
- Paired punctuation marks are frequently beyond the TAG domain of locality.
- One way of punctuation treatment is by marking each punctuation symbol as an adjunct.



# Bottom-up approach (1) [CBVS06]

- Tree-extraction procedure is applied recursively and bottom-up
- During its application, there are multiple elementary trees in various stages are generated. [CBVS06]



## Bottom-up approach (2) [CBVS06]

- One starts at the leaf node and constructs the elementary tree in the following manner:
  - ★ extending the trunk of the partial tree by one level
  - ★ adding complements as substitution nodes to the root of the resulting tree
- elementary trees are extracted in parallel
- one tree is selected to be the *head partial tree* based on the roots of the partial trees and their common parent.
- a partial tree whose root is labeled as adjunct is factored into a modifier auxiliary tree that adjoines onto the head partial tree.

[CBVS06]

# Bottom-up approach (algorithm) [CBVS06]

- 1 The procedure begins by associating an anchor (leave) and its corresponding preterminal with partial tree  $\gamma_0$ . A *partial tree* is an elementary tree in a possibly completed or still incomplete stage of generation.
- 2 The full elementary tree is constructed bottom up from  $\gamma_0$ , alternating between two steps: growing the trunk of the partial tree by one level, and adding substitution nodes as daughters of the root of the partial tree.
- 3 (for continuous steps, see next page)

# Bottom-up approach (algorithm) [CBVS06]

- ① Let  $\mu_1, \mu_2, \dots, \mu_n$  be the root nodes of sibling subtrees  $\gamma_1, \gamma_2, \dots, \gamma_n$  respectively. Let  $\mu_p$  be an immediately dominating node corresponding to these siblings' parent:
  - One partial tree  $\mu_i$  is selected to be the *head partial tree* based on the roots of the partial trees and their common parent  $\mu_p$
  - Roots of the remaining partial trees are marked as complements or adjuncts
  - The head partial tree is then extended by one level
  - Duplicates are made of complement nodes such that they become substitution nodes in the head partial tree.
  - A partial tree whose root is labeled as an adjunct is factored into a modifier auxiliary tree that adjoins onto the head partial tree
  - Repeat until reaching the partial tree whose root is labeled with a start symbol and which has marked substitution sites for the arguments of the head child.

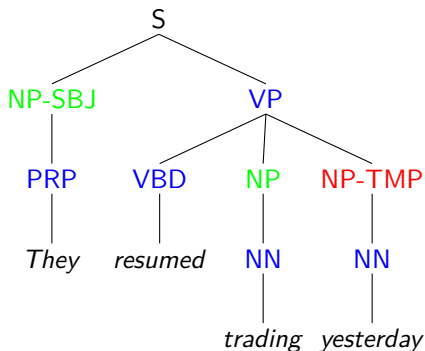
## Bottom-up approach (example)

Head, modifier and complement marking has to be done, just like in the preceding algorithm. Then start with the set of all subtrees consisting of leaves and their mothers.

Repeat until there are no more trees in the set:

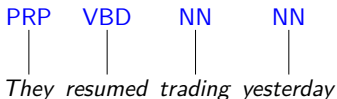
- For each tree obtained so far with root node  $v_r$ :
  - If  $v_r$  is a head:
    - add its mother node and all complements among  $v_r$ 's sisters to the tree;
  - If  $v_r$  is a complement or the root of the entire tree:
    - move the subtree to the final grammar.
  - If  $v_r$  is a modifier to the left (resp. right) of its head sister:
    - add new mother node and foot node sister to the right (resp. left) of  $v_r$ , both labeled with the label of the original mother of  $v_r$ .
    - Then move the tree to the final grammar.

# Bottom-up approach (example)

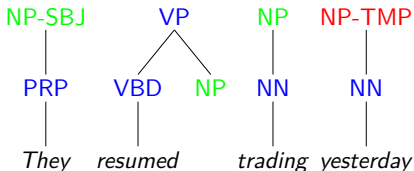


# Bottom-up approach (example)

- First round of the algorithm:



- Second round of the algorithm:



[CBVS06]

if root **head**:

add mother node  
and complements  
among sisters

if root **complement**:

elementary tree done

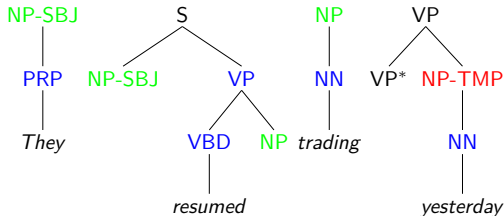
if root **modifier** with

mother node category X:  
add new mother node  
and foot node sister  
both labeled X



# Bottom-up approach (example)

- Third round of the algorithm:

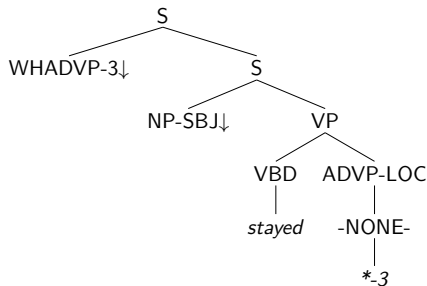
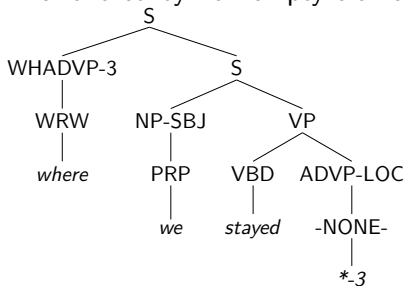


[CBVS06]

- if root **head**:
  - add mother node
  - and complements
  - among sisters
- if root **complement**:
  - elementary tree done
- if root **modifier** with
  - mother node category X:
  - add new mother node
  - and foot node sister
  - both labeled X

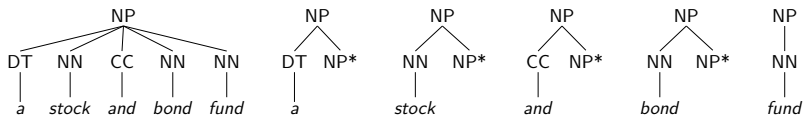
## Some considerations: Empty elements

- Empty elements are terminal nodes in bracketings that do not represent lexical items
- Present in Penn Treebank and in the hand-crafted TAGs → desirable to keep in the extracted LTAG?
- Modification of extraction algorithm can be made which grafts partial trees anchored by empty elements onto partial trees anchored by non-empty elements.



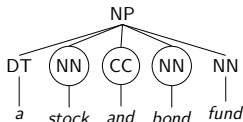
# Some considerations: Conjunctions

- According to the proposed algorithm, constructions with conjunctions are treated like this:

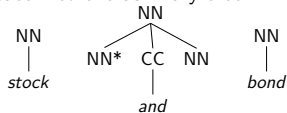


## Some considerations: Conjunctions

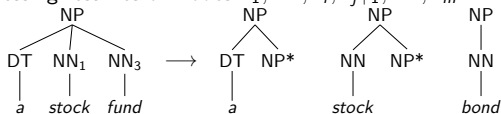
- Identify a sequence of daughters  $n_1, \dots, n_j$ , such that  $1 \leq i < j \leq m$  representing conjuncts and associated conjunctions.  $n_i$  is the leftmost conjunct,  $n_j$  the rightmost conjunct.



- From the daughters  $n_{i+1}, \dots, n_j$  factor the conjunction into its own auxiliary tree. Auxiliary tree is anchored with conjunction. It adjoins onto  $n_i$ . Conjuncts become initial trees that substitute into the auxiliary tree.



- Default processing resumes on nodes  $n_1, \dots, n_i, n_{j+1}, \dots, n_m$



# Extracted TAGs: statistics

Parameters	French LTAG [BvCSK18]	German LTAG [Kae12]	English LTAG [KFM <sup>+</sup> 17]
Elementary trees	5145	3426	4727
Elementary trees once	2693	1562	2165
POS tags	13 / 26	53	36
Sentences	21550	50000	44168
Avg. sentence length	31.34	17.71	appr. 20

# References I

- [BvCSK18] Tatiana Bladier, Andreas van Cranenburgh, Younes Samih, and Laura Kallmeyer. German and french neural supertagging experiments for Itag parsing (to appear). *ACL 2018 Student Research Workshop, Melbourne, Australia*, 2018.
- [CBVS06] John Chen, Srinivas Bangalore, and K Vijay-Shanker. Automated extraction of tree-adjoining grammars from treebanks. *Natural Language Engineering*, 12(3):251–299, 2006.
- [G<sup>+</sup>98] XTAG Research Group et al. A lexicalized tree adjoining grammar for english. *arXiv preprint cs/9809024*, 1998.
- [Kae12] Miriam Kaeshammer. *A German treebank and lexicon for tree-adjoining grammars*. PhD thesis, Master's thesis, Universitat des Saarlandes, Saarlandes, Germany, 2012.
- [KFM<sup>+</sup>17] Jungo Kasai, Robert Frank, Tom McCoy, Owen Rambow, and Alexis Nasr. Tag parsing with neural networks and vector representations of supertags. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1713–1723, 2017.
- [MSM93] Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993.
- [XP06] Fei Xia and Martha Palmer. From treebanks to tree-adjoining grammars. *Supertagging: using complex lexical descriptions in natural language processing*, pages 1–39, 2006.

# References II

- [XPJ00] Fei Xia, Martha Palmer, and Aravind Joshi.  
A uniform method of grammar extraction and its applications.  
In *Proceedings of the 2000 Joint SIGDAT conference on Empirical methods in natural language processing and very large corpora: held in conjunction with the 38th Annual Meeting of the Association for Computational Linguistics-Volume 13*, pages 53–62. Association for Computational Linguistics, 2000.