

# CS626: NLP, Speech and the Web

Pushpak Bhattacharyya  
CSE Dept.,  
IIT Bombay

Lecture 17, 18: Probabilistic parsing;  
parser comparisons

24<sup>th</sup> and 27<sup>th</sup> September, 2012

*(16<sup>th</sup> lecture was on UNL by Avishek)*

# Example of Sentence labeling: Parsing

[S<sub>1</sub> [S [S [VP [V<sub>B</sub> Come] [NP [NNP July]]]]  
[,']  
[CC and]  
[S [NP [DT the] [JJ IIT] [NN campus]]  
[VP [AUX is]  
[ADJP [JJ abuzz]  
[PP [IN with]  
[NP [ADJP [JJ new] [CC and] [V<sub>BG</sub> returning]]  
[NNS students]]]]]]]  
[.]]]

# Noisy Channel Modeling

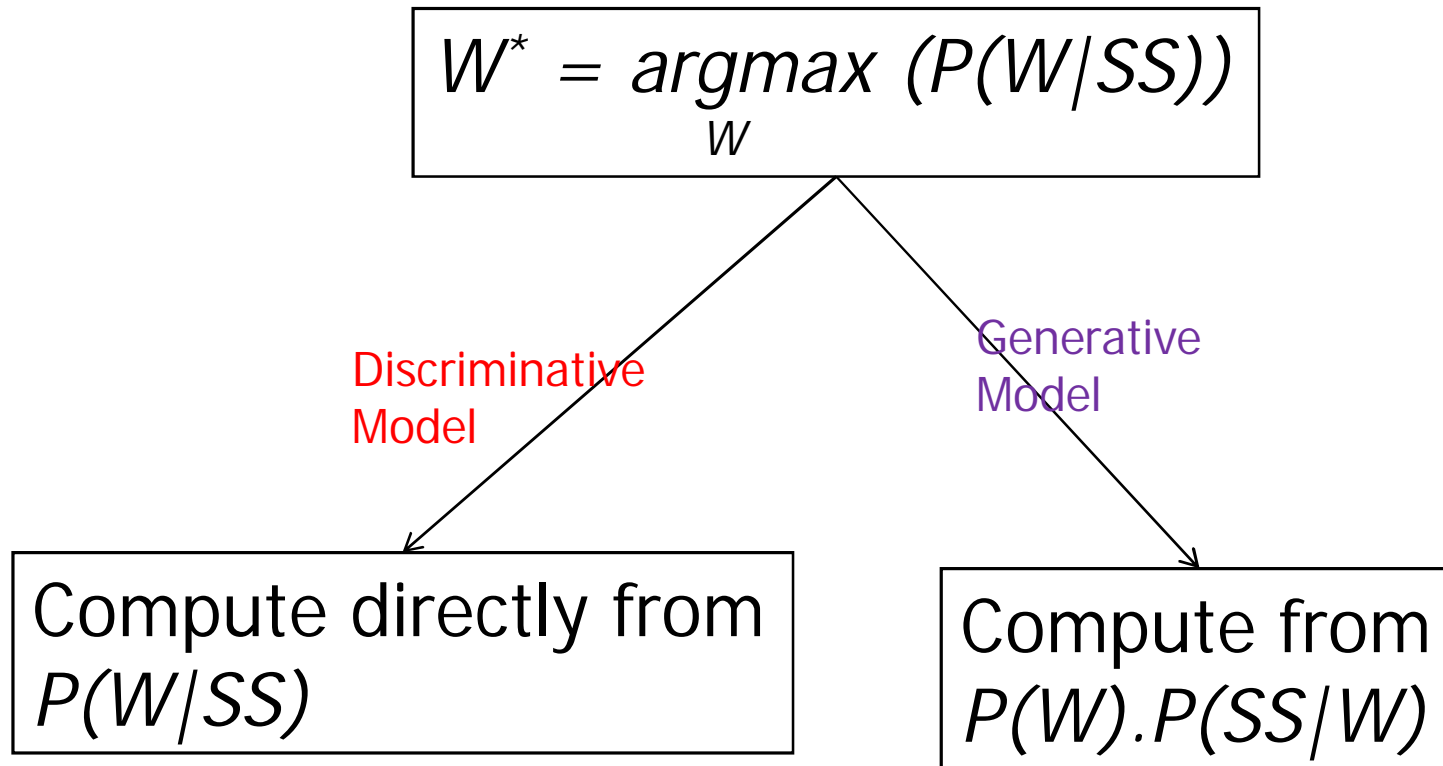


$$\begin{aligned} T^* &= \underset{T}{\operatorname{argmax}} [P(T|S)] \\ &= \underset{T}{\operatorname{argmax}} [P(T).P(S|T)] \\ &= \underset{T}{\operatorname{argmax}} [P(T)], \text{ since given the parse the} \\ &\quad \text{sentence is completely} \\ &\quad \text{determined and } P(S|T)=1 \end{aligned}$$

# Corpus

- A collection of text called *corpus*, is used for collecting various language data
- With annotation: more information, but manual labor intensive
- Practice: *label automatically; correct manually*
- The famous *Brown Corpus* contains 1 million tagged words.
- **Switchboard**: very famous corpora 2400 conversations, 543 speakers, many US dialects, annotated with orthography and phonetics

# Discriminative vs. Generative Model



# Language Models

- N-grams: sequence of n consecutive words/characters
- Probabilistic / Stochastic Context Free Grammars:
  - Simple probabilistic models capable of handling recursion
  - A CFG with probabilities attached to rules
  - Rule probabilities → how likely is it that a particular rewrite rule is used?

# PCFGs

- Why PCFGs?
  - Intuitive probabilistic models for tree-structured languages
  - Algorithms are extensions of HMM algorithms
  - Better than the n-gram model for language modeling.

# Data driven parsing is tied up with what is seen in the training data

- “*Stand right walk left.*” Often a message on the airport escalators.



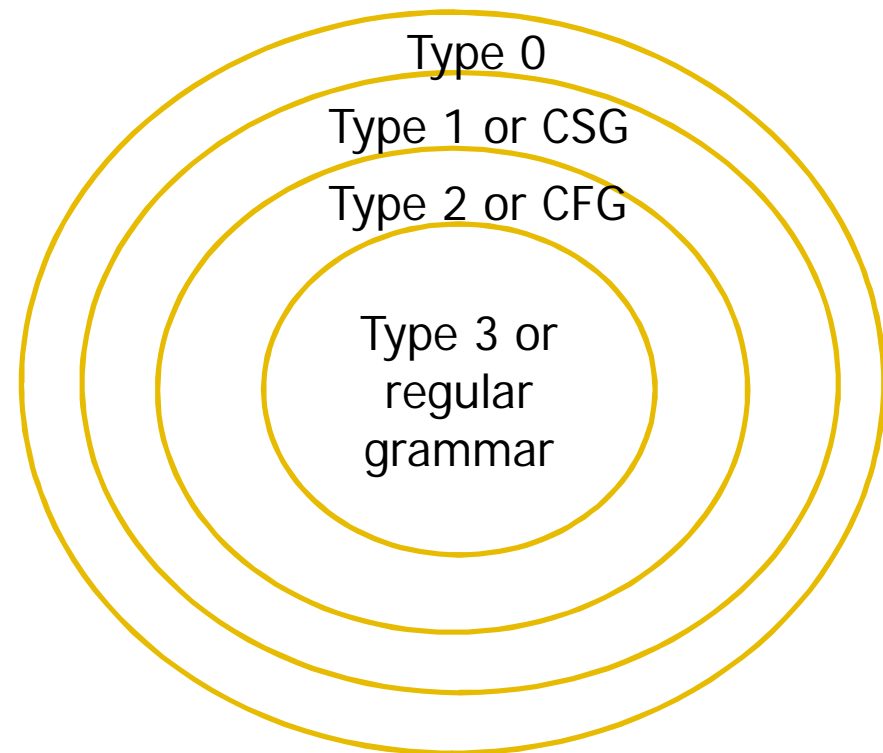
- “Walk left” can come from :- “*After climbing , we realized the walk left was still considerable.*”
- “Stand right” can come from: - “*The stand right though it seemed at that time, does not stand the scrutiny now.*”
- Then the given sentence “Stand right walk left” will never be parsed correctly.



# Chomski Hierarchy

Properties of the hierarchy:-

- *The containment is proper.*
- *Grammar is not learnable even for the lowest type from positive examples alone.*



$Type\ 0 \supset Type\ 1 \supset Type\ 2 \supset Type\ 3$

# Rationalism vs Empiricism

- Scientific paper : - “*The pendulum swings back*” by Ken Church.
- Every 20 years there is a shift from rationalism to empiricism.
- AI approaches are of two types:-
  - Rational (Theory and model driven e.g. Brill Tagger; humans decide the templates)
  - Empirical (Data driven e.g. HMM)
- Is it possible to classify probabilistic CFG into any of these categories?
  - Very difficult to answer.
  - Anything that comes purely from brain can be put into the category of “rationalism”.
- There is no pure rationalism or pure empiricism

# Laws of machine learning

- “Learning from vacuum (zero knowledge) is impossible”.
- “Inductive bias decides what to learn”.
  - *In case of POS tagger , bias is that we assume the tagger to be a HMM model.*

# Formal Definition of PCFG

- A PCFG consists of
  - A set of terminals  $\{w_k\}$ ,  $k = 1, \dots, V$   
 $\{w_k\} = \{ \text{child, teddy, bear, played...} \}$
  - A set of non-terminals  $\{N^i\}$ ,  $i = 1, \dots, n$   
 $\{N_i\} = \{ \text{NP, VP, DT...} \}$
  - A designated start symbol  $N^1$
  - A set of rules  $\{N^i \rightarrow \zeta^j\}$ , where  $\zeta^j$  is a sequence of terminals & non-terminals  
 $\text{NP} \rightarrow \text{DT NN}$
  - A corresponding set of rule probabilities

# Rule Probabilities

- Rule probabilities are such that

$$\forall i \sum_j P(N^i \rightarrow \zeta^j) = 1$$

*E.g.*,  $P(\text{NP} \rightarrow \text{DT NN}) = 0.2$

$$P(\text{NP} \rightarrow \text{NN}) = 0.5$$

$$P(\text{NP} \rightarrow \text{NP PP}) = 0.3$$

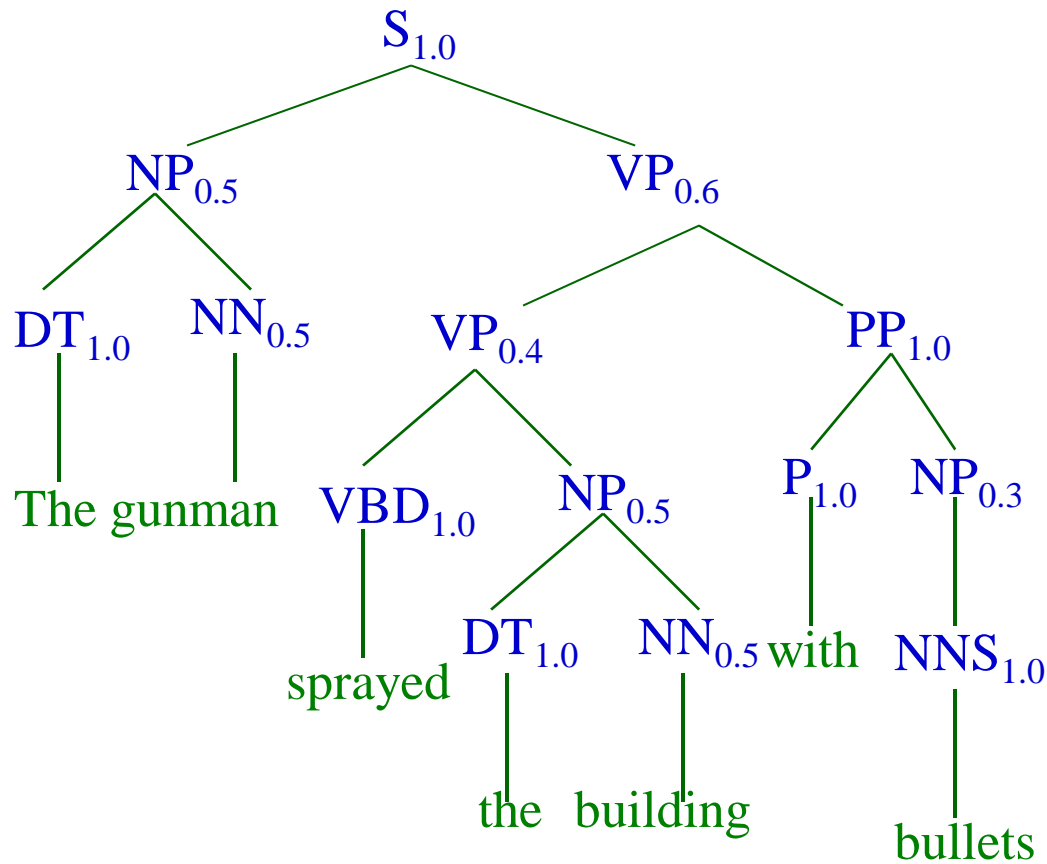
- $P(\text{NP} \rightarrow \text{DT NN}) = 0.2$ 
  - Means 20 % of the training data parses use the rule  $\text{NP} \rightarrow \text{DT NN}$

# Probabilistic Context Free Grammars

- $S \rightarrow NP VP$  1.0
- $NP \rightarrow DT NN$  0.5
- $NP \rightarrow NNS$  0.3
- $NP \rightarrow NP PP$  0.2
- $PP \rightarrow P NP$  1.0
- $VP \rightarrow VP PP$  0.6
- $VP \rightarrow VBD NP$  0.4
- $DT \rightarrow the$  1.0
- $NN \rightarrow gunman$  0.5
- $NN \rightarrow building$  0.5
- $VBD \rightarrow sprayed$  1.0
- $NNS \rightarrow bullets$  1.0

# Example Parse $t_1$

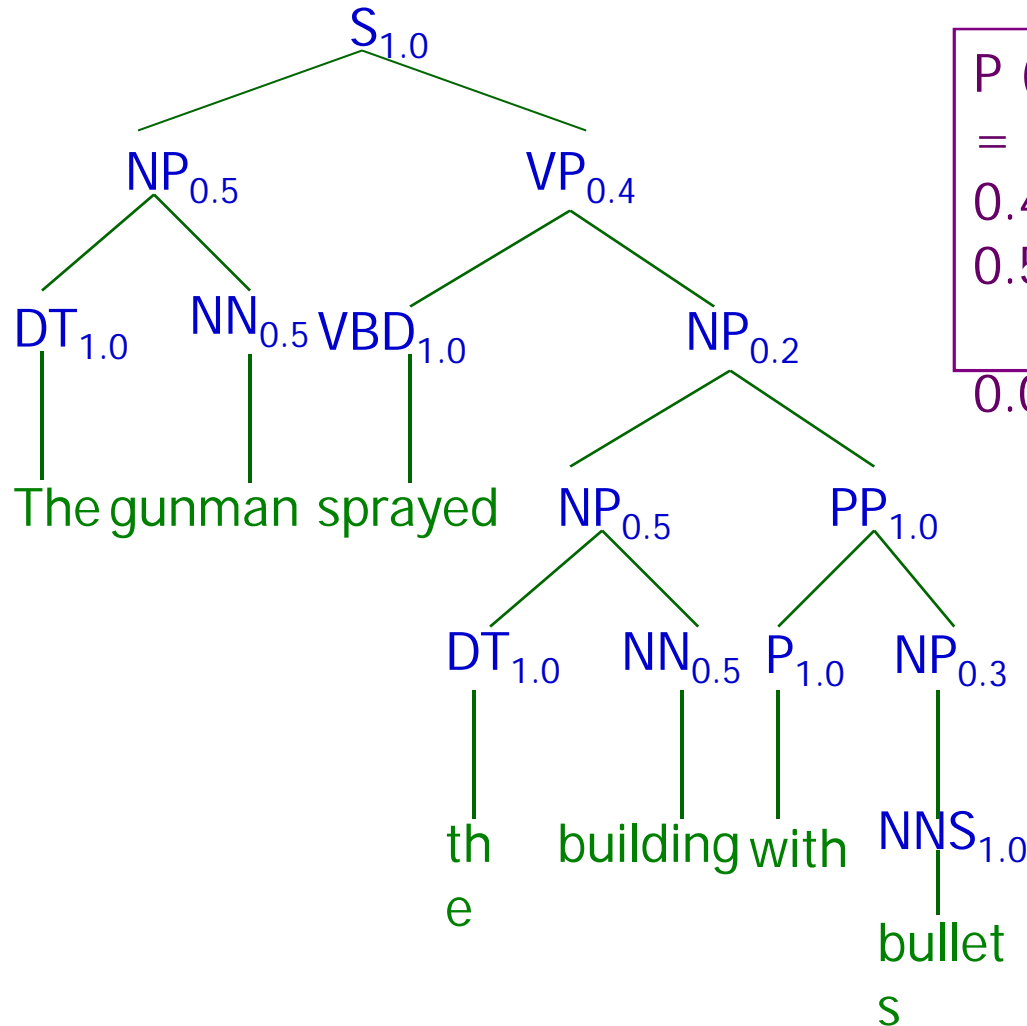
- The gunman sprayed the building with bullets.



$$\begin{aligned} P(t_1) &= 1.0 * \\ &0.5 * 1.0 * 0.5 * 0.6 * 0.4 * 1.0 \\ &* 0.5 * 1.0 * 0.5 * 1.0 * 1.0 * \\ &0.3 * 1.0 &= \\ &0.00225 \end{aligned}$$

# Another Parse $t_2$

- The gunman sprayed the building with bullets.



$$\begin{aligned} P(t_2) &= 1.0 * 0.5 * 1.0 * 0.5 * \\ &0.4 * 1.0 * 0.2 * 0.5 * 1.0 * \\ &0.5 * 1.0 * 1.0 * 0.3 * 1.0 \\ &= \\ &0.0015 \end{aligned}$$



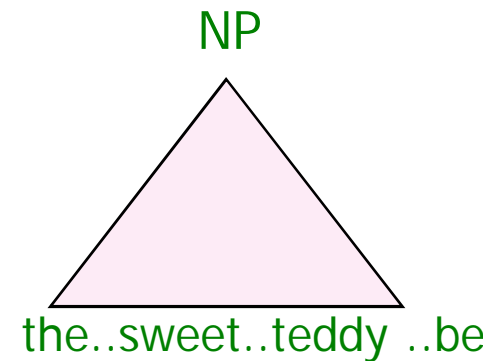
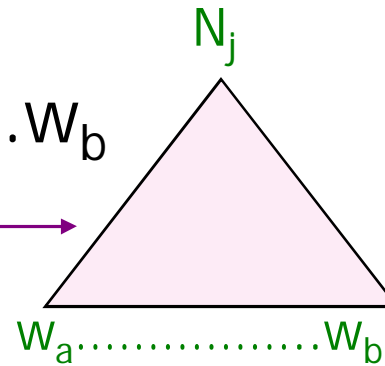
# Probability of a sentence

## ■ Notation :

■  $w_{ab}$  – subsequence  $w_a \dots w_b$

■  $N_j$  dominates  $w_a \dots w_b$  →

or  $\text{yield}(N_j) = w_a \dots w_b$



## • Probability of a sentence = $P(w_{1m})$

$$P(w_{1m}) = \sum_t P(w_{1m}, t) \quad \rightarrow \text{Where } t \text{ is a parse tree of the sentence}$$

$$= \sum_t P(t) P(w_{1m} | t)$$

$$= \sum_{t: \text{yield}(t)=w_{1m}} P(t) \quad \because P(w_{1m} | t) = 1$$

If  $t$  is a parse tree for the sentence  $w_{1m}$ , this will be 1 !!

# Assumptions of the PCFG model

- Place invariance :

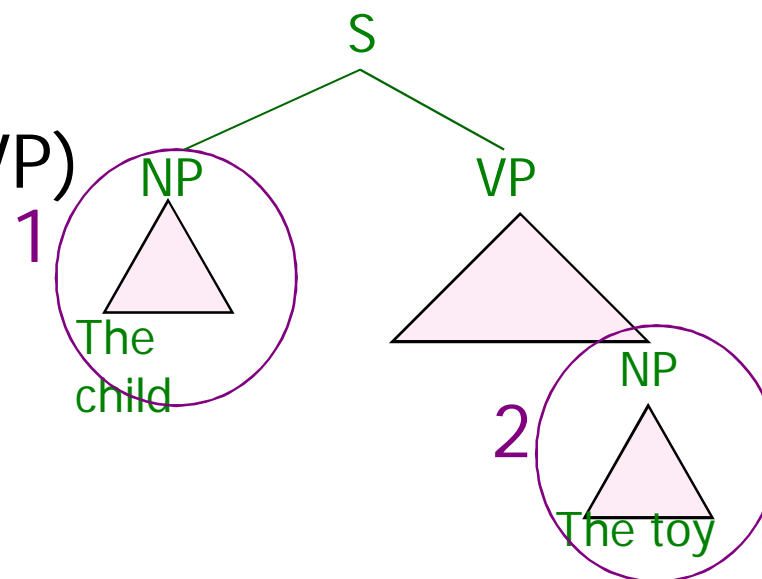
$P(\text{NP} \rightarrow \text{DT NN})$  is same in locations 1 and 2

- Context-free :

$P(\text{NP} \rightarrow \text{DT NN} \mid \text{anything outside "The child"})$   
 $= P(\text{NP} \rightarrow \text{DT NN})$

- Ancestor free : At 2,

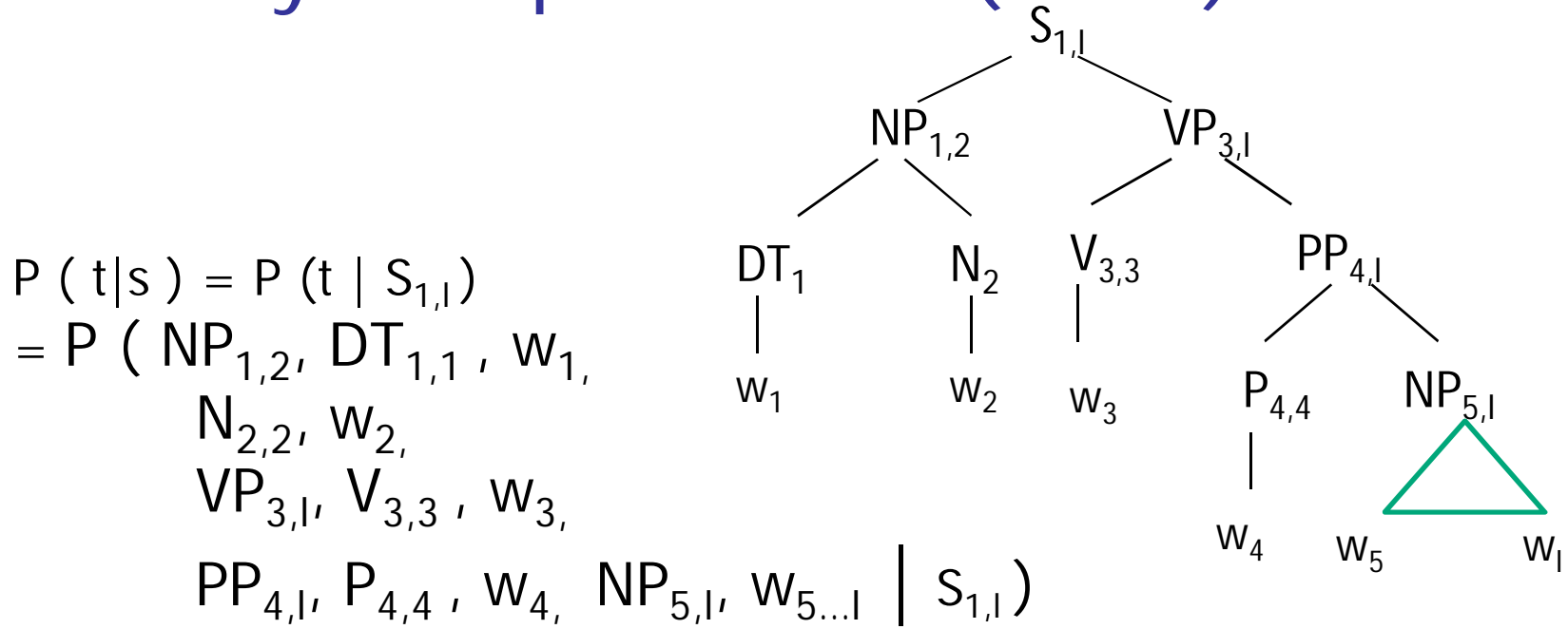
$P(\text{NP} \rightarrow \text{DT NN} \mid \text{its ancestor is VP})$   
 $= P(\text{NP} \rightarrow \text{DT NN})$



# Probability of a parse tree

- *Domination* : We say  $N_j$  dominates from  $k$  to  $l$ , symbolized as  $N_{k,l}^j$  if  $W_{k,l}$  is derived from  $N_j$
- $P(\text{tree} \mid \text{sentence}) = P(\text{tree} \mid S_{1,l})$   
where  $S_{1,l}$  means that the start symbol  $S$  dominates the word sequence  $W_{1,l}$
- $P(t \mid s)$  approximately equals joint probability of constituent non-terminals dominating the sentence fragments (next slide)

# Probability of a parse tree (cont.)



$$\begin{aligned}
 &= P(NP_{1,2}, VP_{3,1} | S_{1,1}) * P(DT_{1,1}, N_{2,2} | NP_{1,2}) * D(w_1 | DT_{1,1}) * \\
 &P(w_2 | N_{2,2}) * P(V_{3,3}, PP_{4,1} | VP_{3,1}) * P(w_3 | V_{3,3}) * P(P_{4,4}, NP_{5,1} | \\
 &PP_{4,1}) * P(w_4 | P_{4,4}) * P(w_{5..1} | NP_{5,1})
 \end{aligned}$$

*(Using Chain Rule, Context Freeness and Ancestor Freeness)*

# HMM $\leftrightarrow$ PCFG

- $O$  observed sequence  $\leftrightarrow W_{1m}$  sentence
- $X$  state sequence  $\leftrightarrow t$  parse tree
- $\mu$  model  $\leftrightarrow G$  grammar
  
- Three fundamental questions

# HMM $\leftrightarrow$ PCFG

- How likely is a certain observation given the model?  $\leftrightarrow$  How likely is a sentence given the grammar?

$$P(O | \mu) \leftrightarrow P(w_{1m} | G)$$

- How to choose a state sequence which best explains the observations?  $\leftrightarrow$  How to choose a parse which best supports the sentence?

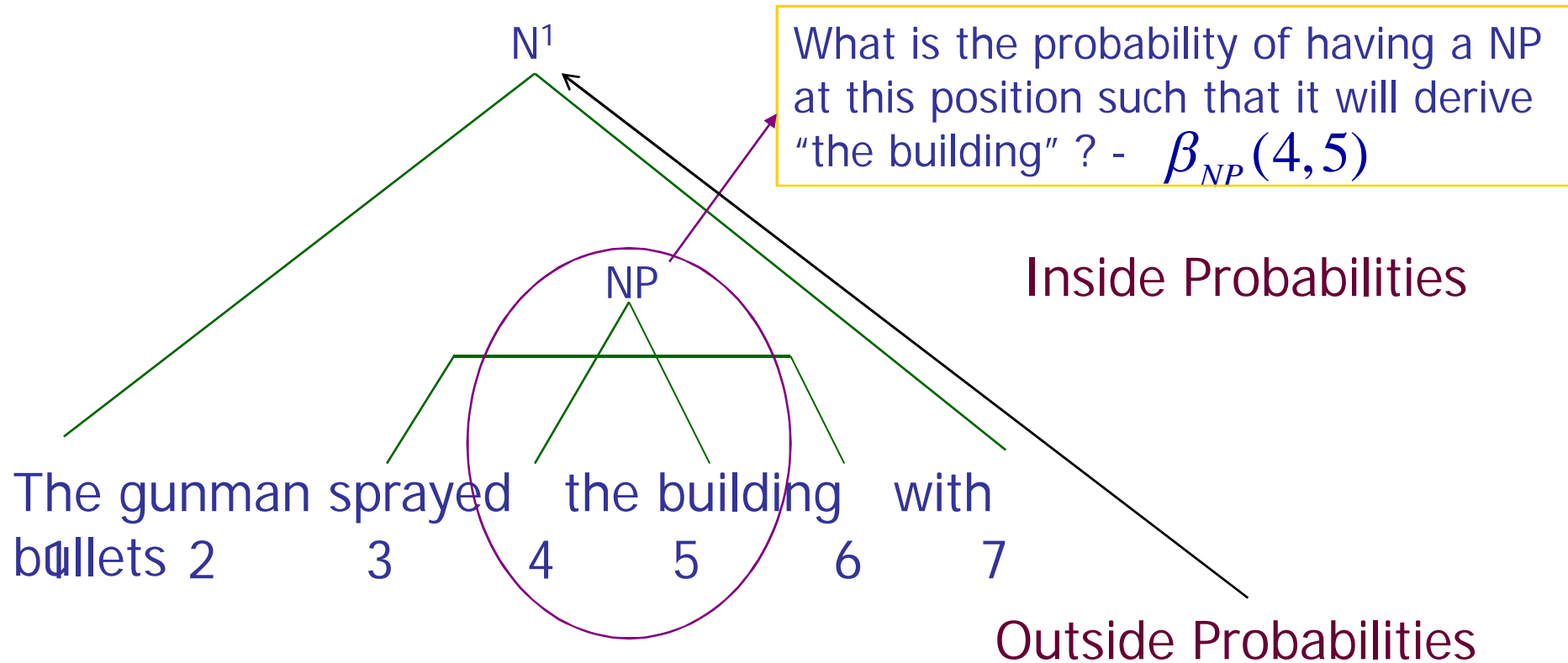
$$\arg \max_X P(X | O, \mu) \leftrightarrow \arg \max_t P(t | w_{1m}, G)$$

# HMM $\leftrightarrow$ PCFG

- How to choose the model parameters that best explain the observed data?  $\leftrightarrow$  How to choose rule probabilities which maximize the probabilities of the observed sentences?

$$\arg \max_{\mu} P(O | \mu) \leftrightarrow \arg \max_G P(w_{1m} | G)$$

# Interesting Probabilities



What is the probability of having a NP at this position such that it will derive "the building" ? -  $\beta_{NP}(4,5)$

What is the probability of starting from N<sup>1</sup> and deriving "The gunman sprayed", a NP and "with bullets" ? -  $\alpha_{NP}(4,5)$



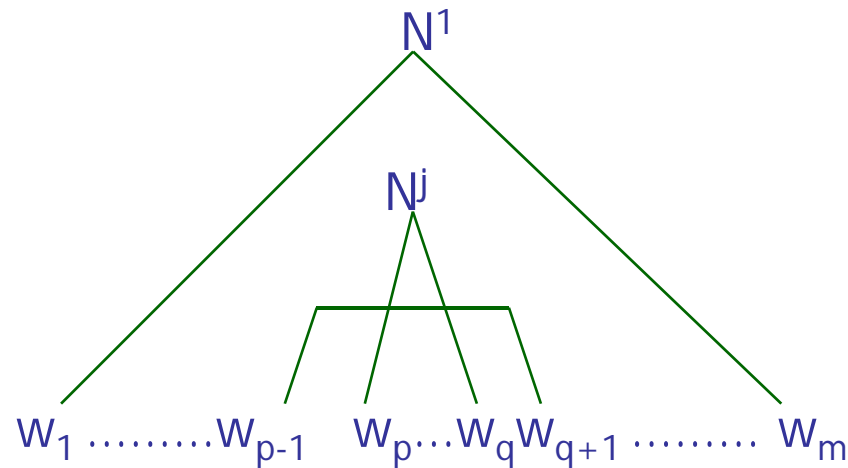
# Interesting Probabilities

- Random variables to be considered
  - The non-terminal being expanded.  
*E.g.*, NP
  - The word-span covered by the non-terminal.  
*E.g.*, (4,5) refers to words “the building”
- While calculating probabilities, consider:
  - The rule to be used for expansion :  
*E.g.*, NP → DT NN
  - The probabilities associated with the RHS non-terminals : *E.g.*, DT subtree’s inside/outside probabilities & NN subtree’s inside/outside probabilities

# Outside Probability

- $\alpha_j(p, q)$  : The probability of beginning with  $N^1$  & generating the non-terminal  $N^j_{pq}$  and all words outside  $w_p \dots w_q$

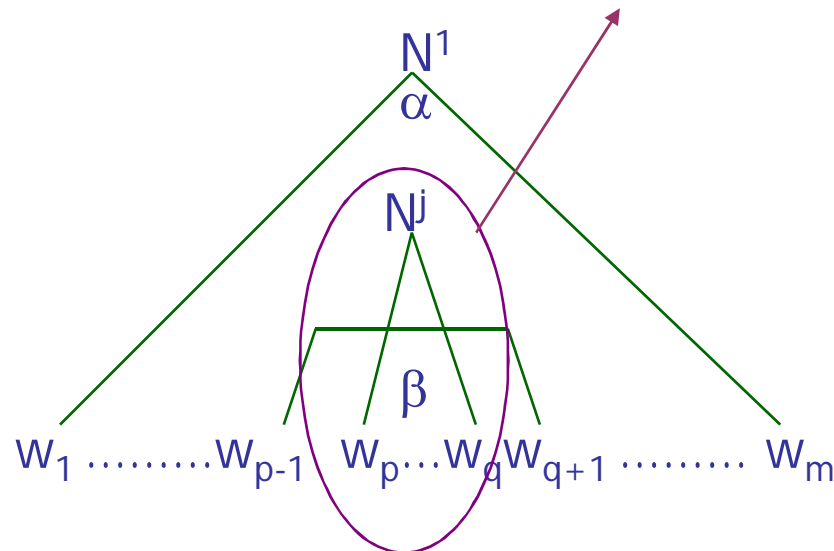
$$\alpha_j(p, q) = P(w_{1(p-1)}, N^j_{pq}, w_{(q+1)m} \mid G)$$



# Inside Probabilities

- $\beta_j(p, q)$  : The probability of generating the words  $w_p \dots w_q$  starting with the non-terminal  $N_j^j$ .

$$\beta_j(p, q) = P(w_{pq} | N_{pq}^j, G)$$

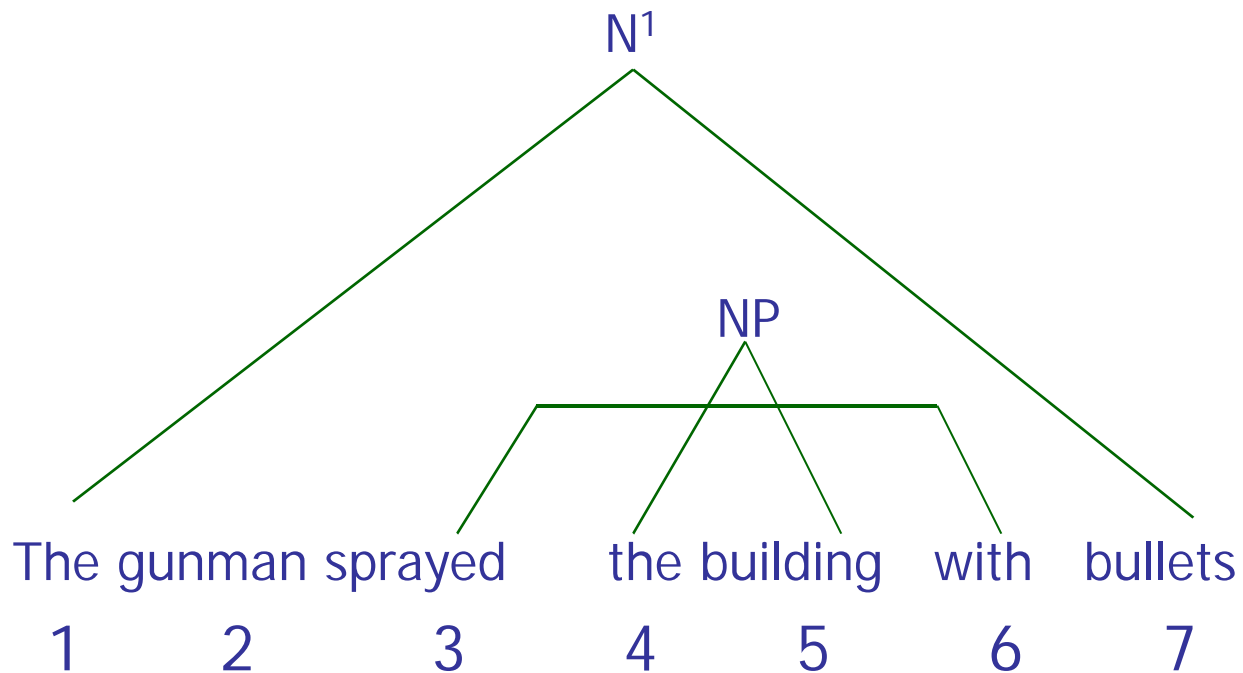


# Outside & Inside Probabilities: example

$\alpha_{NP}(4,5)$  for "the building"

$= P(\text{The gunman sprayed, } NP_{4,5}, \text{ with bullets} \mid G)$

$\beta_{NP}(4,5)$  for "the building"  $= P(\text{the building} \mid NP_{4,5}, G)$



# Parsers Comparison

*(Charniack, Collins, Stanford, RASP)*

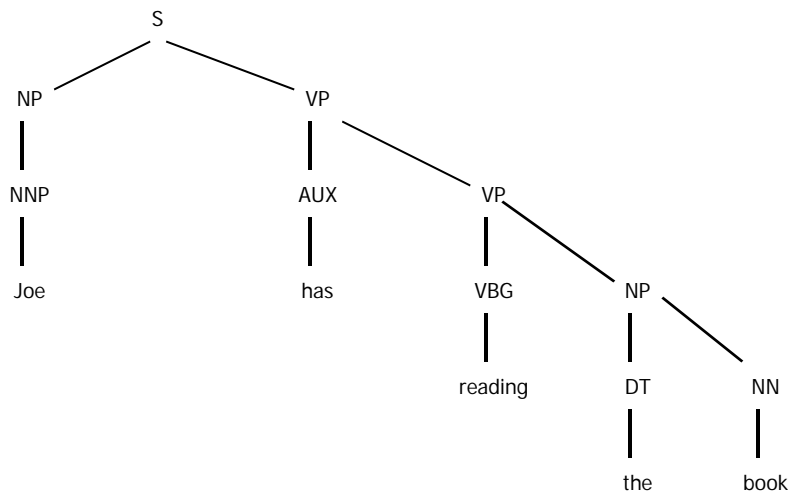
Study by many masters students,  
notably,

*Avishek, Nikhilesh, Abhishek and  
Harshada*

# Parser comparison: Handling ungrammatical sentences

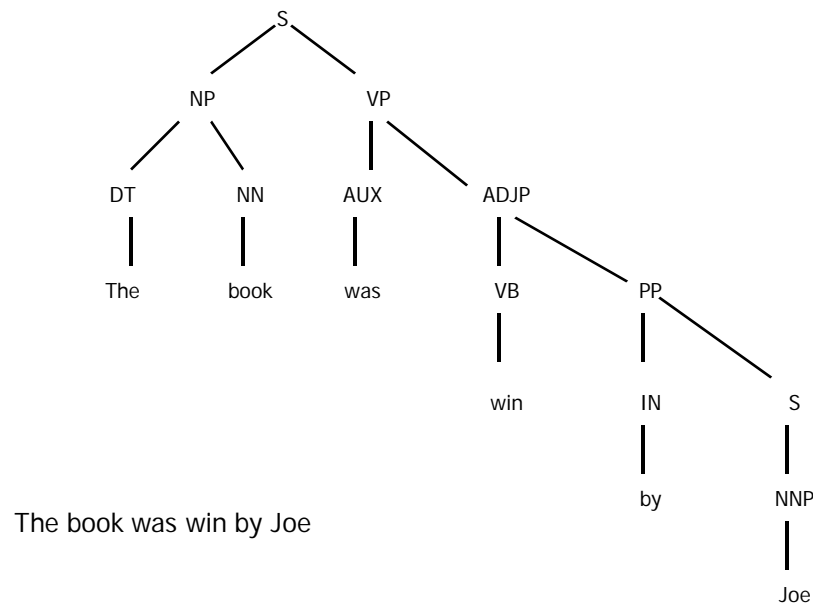
# Charniak (ungrammatical 1)

- Here *has* is tagged as AUX



Joe has reading the book

# Charniak (ungrammatical 2)

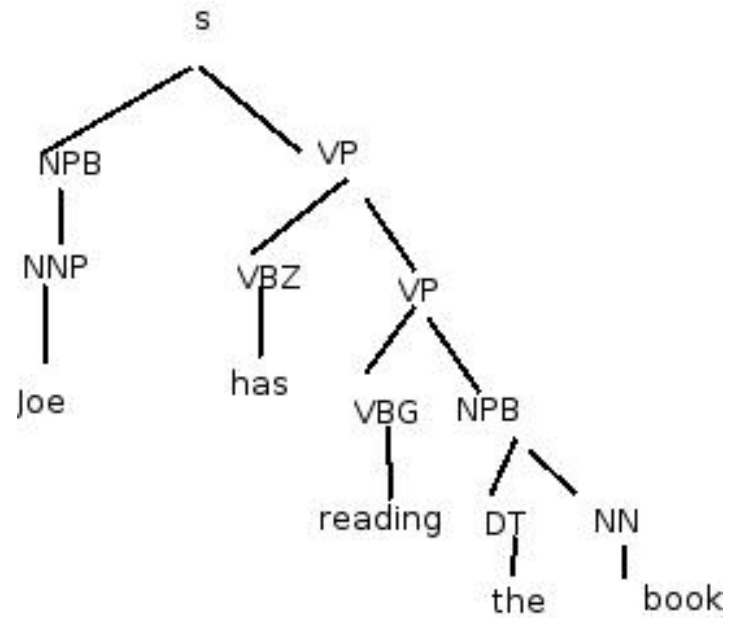


- *Win* is treated as a verb and it does not make any difference whether it is in the present or the past tense



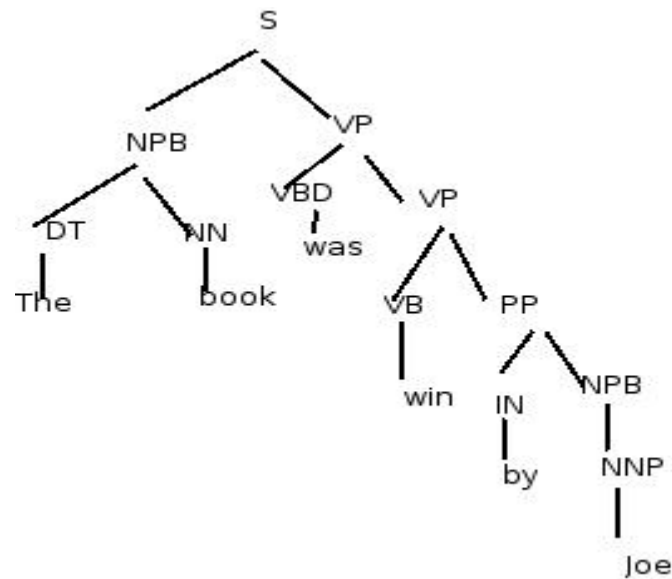
# Collins (ungrammatical 1)

- *Has* should have been AUX.



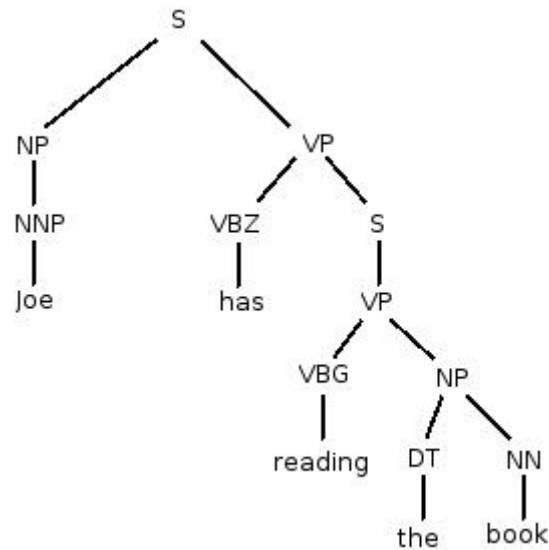
# Collins (ungrammatical 2)

- Same as charniack



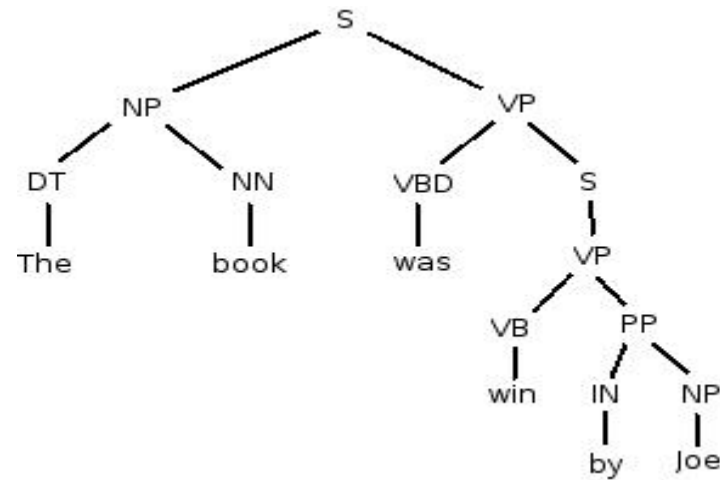
# Stanford (ungrammatical 1)

- *has* is treated as VBZ and not AUX.



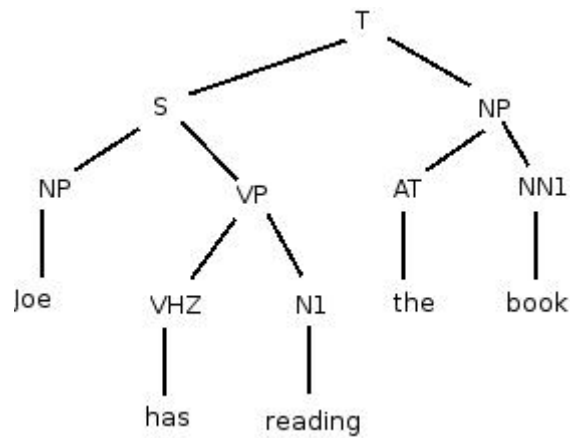
# Stanford (ungrammatical 2)

- Same as Charniak



# RASP (ungrammatical 1)

- Inaccurate tree



# Observation

- For the sentence 'Joe has reading the book' Charniak performs the best; it is able to predict that the word 'has' in the sentence should actually be an AUX
- Though both the RASP and Collins can produce a parse tree, they both cannot predict that the sentence is not grammatically correct
- Stanford performs the worst, it inserts extra 'S' nodes into the parse tree.

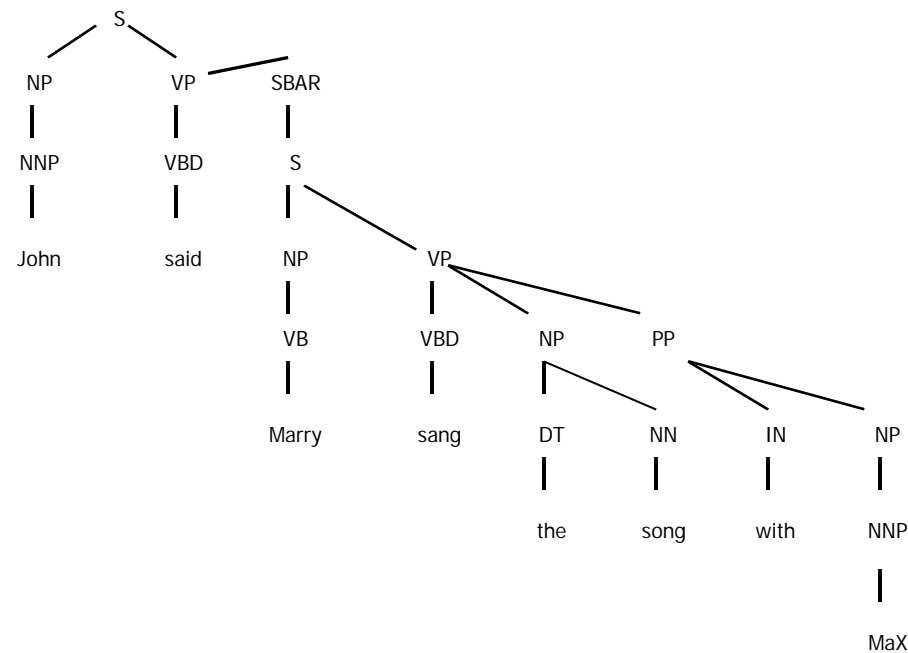
## Observation *(contd.)*

- For the sentence 'The book was win by Joe', all the parsers give the same parse structure which is correct.

Ranking in case of multiple parses



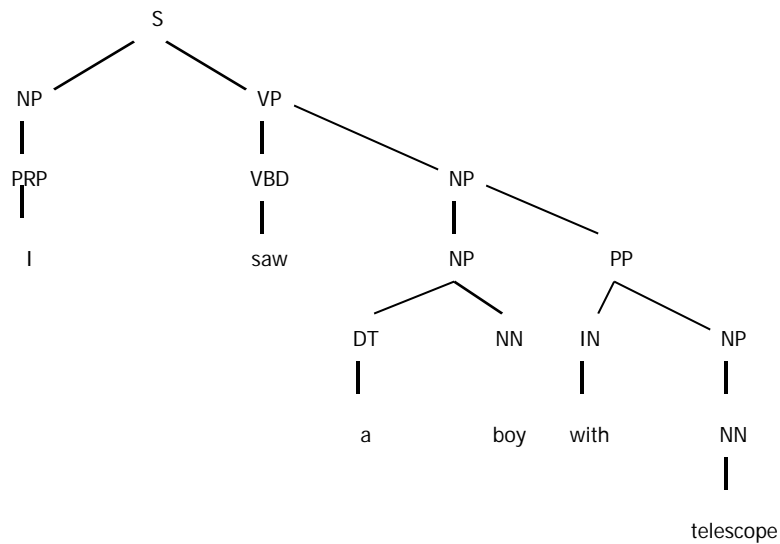
# Charniak (Multiple Parses 1)



- The parse produced is semantically correct

John said Marry sang the song with Max

# Charniak (Multiple Parses 2)

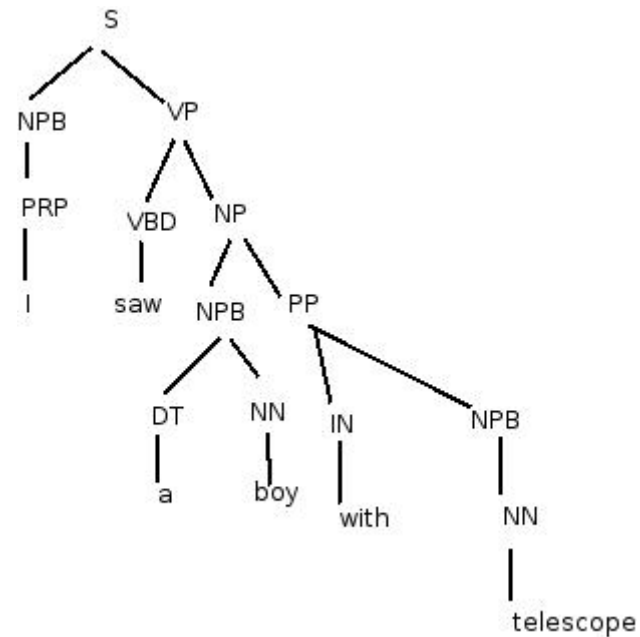


I saw a boy with telescope

- PP is attached to NP which is one of the correct meanings

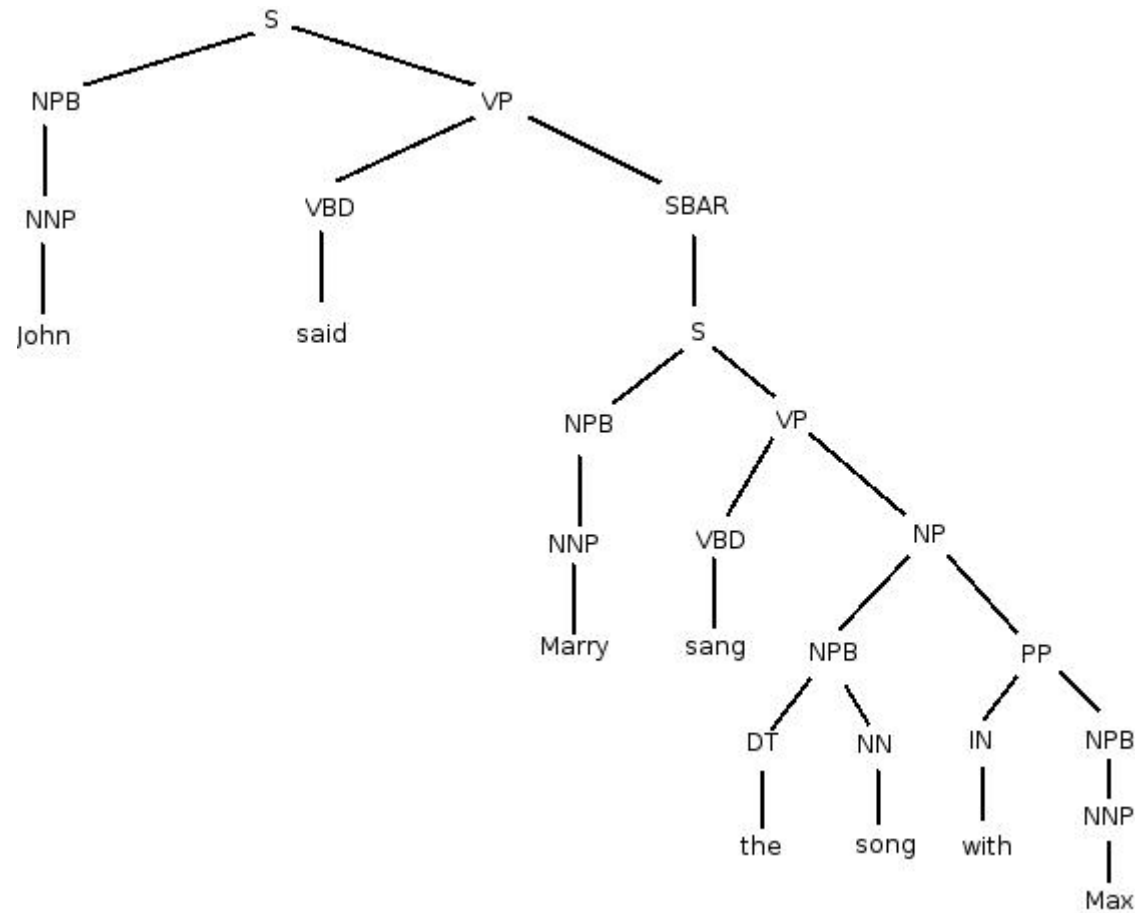
# Collins (Multiple Parses 1)

- Same as Charniak.



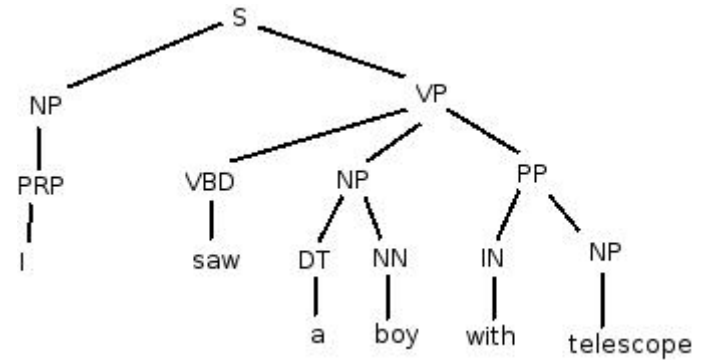
# Collins (Multiple Parses 2)

- Same as Charniak



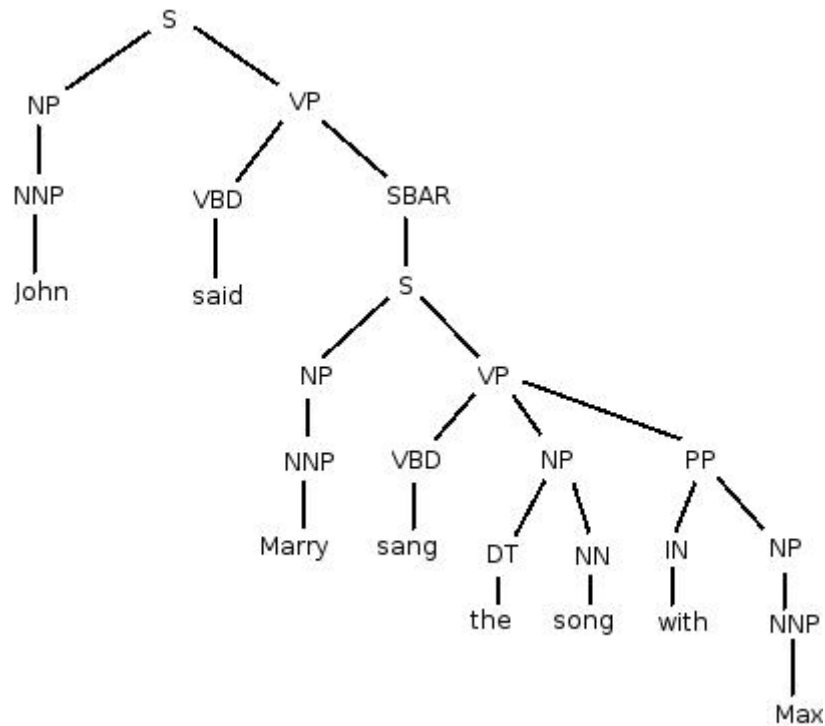
# Stanford (Multiple Parses 1)

- PP is attached to VP which is one of the correct meanings possible



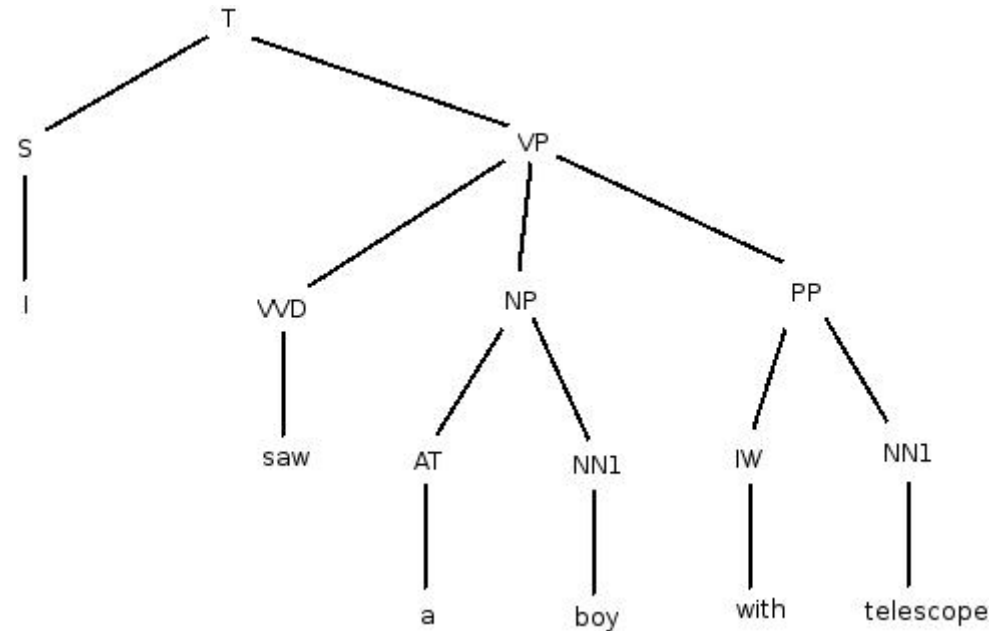
# Stanford (Multiple Parses 2)

- Same as Charniak.



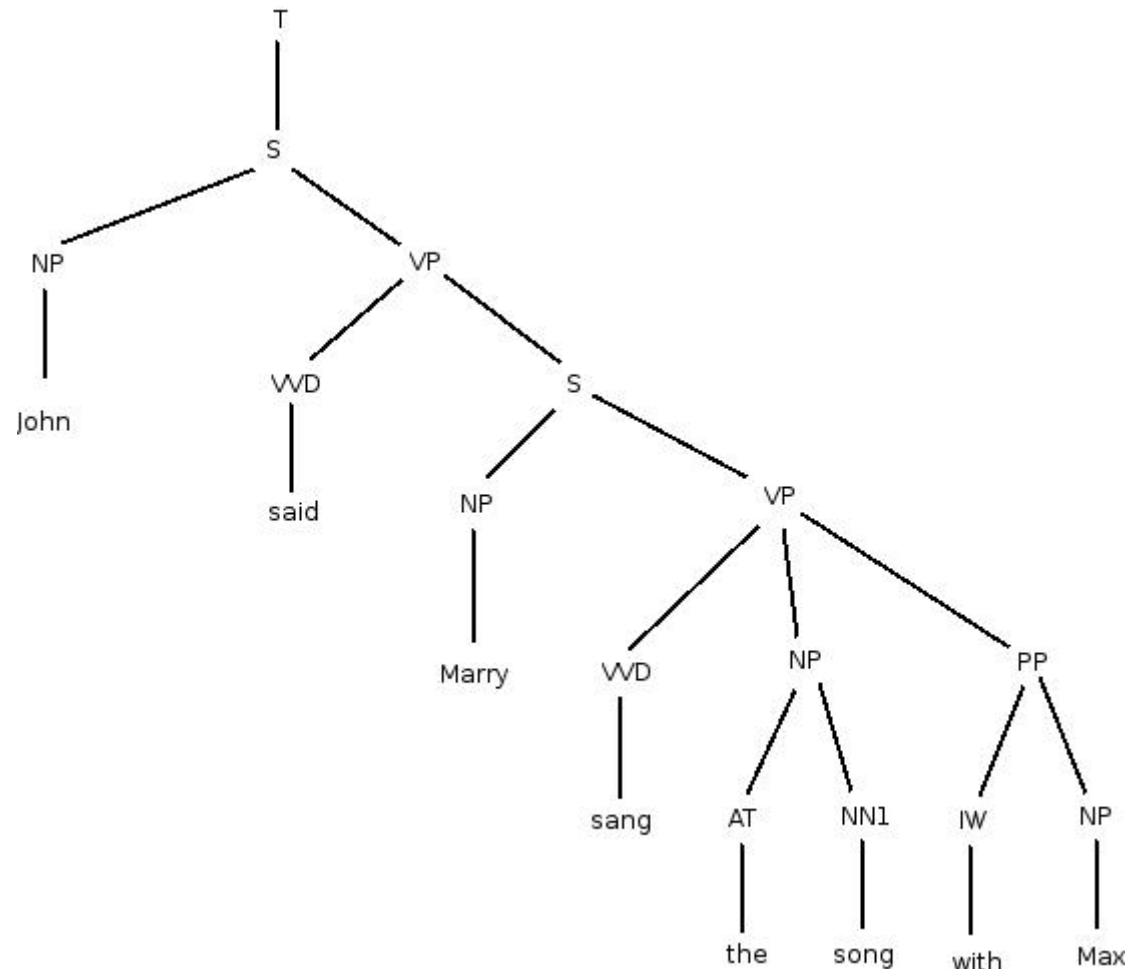
# RASP (Multiple Parses 1)

- PP is attached to VP.



# RASP (Multiple Parses 2)

- The change in the pos tags as compared to charniak is due to the different corpora but the parse trees are comparable.





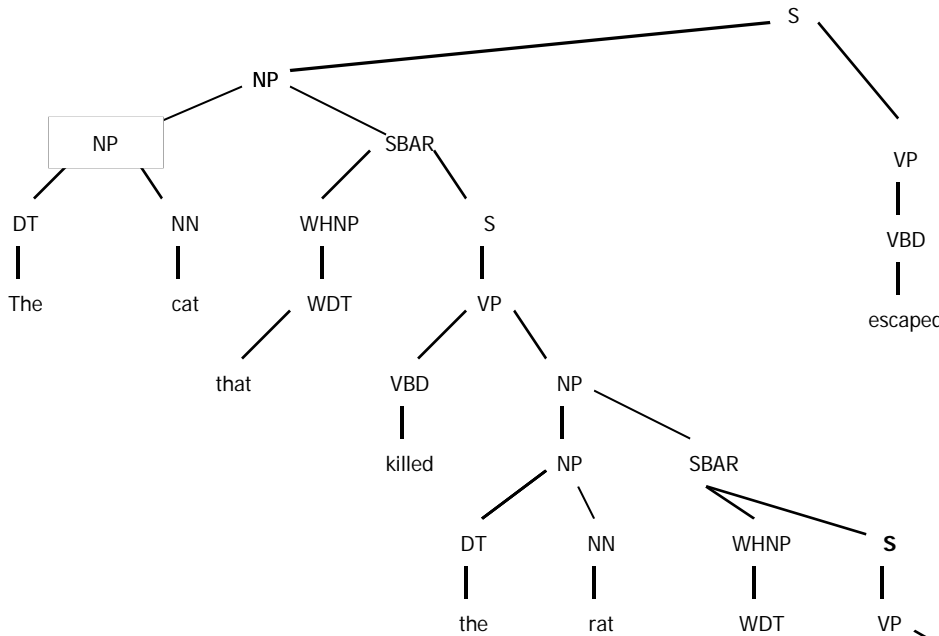


# Time taken

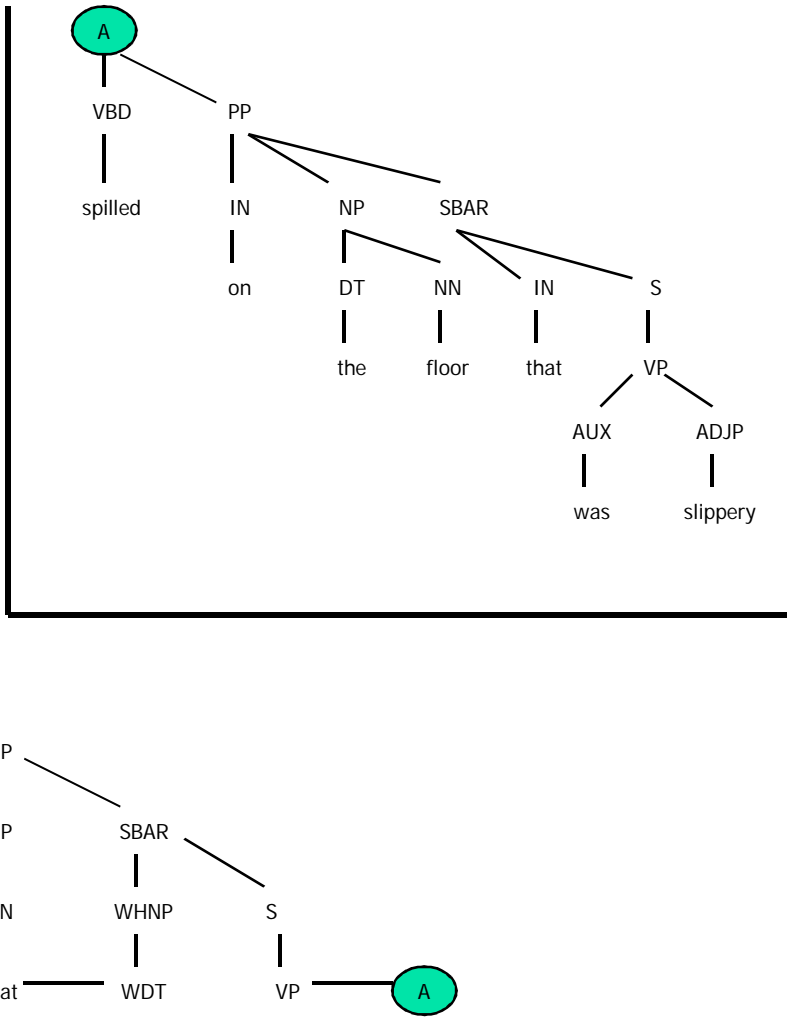
- 54 instances of the sentence 'This is just to check the time' is used to check the time
- Time taken
  - Collins : 40s
  - Stanford : 14s
  - Charniak : 8s
  - RASP : 5s

# Embedding Handling

# Charniak (Embedding 1)

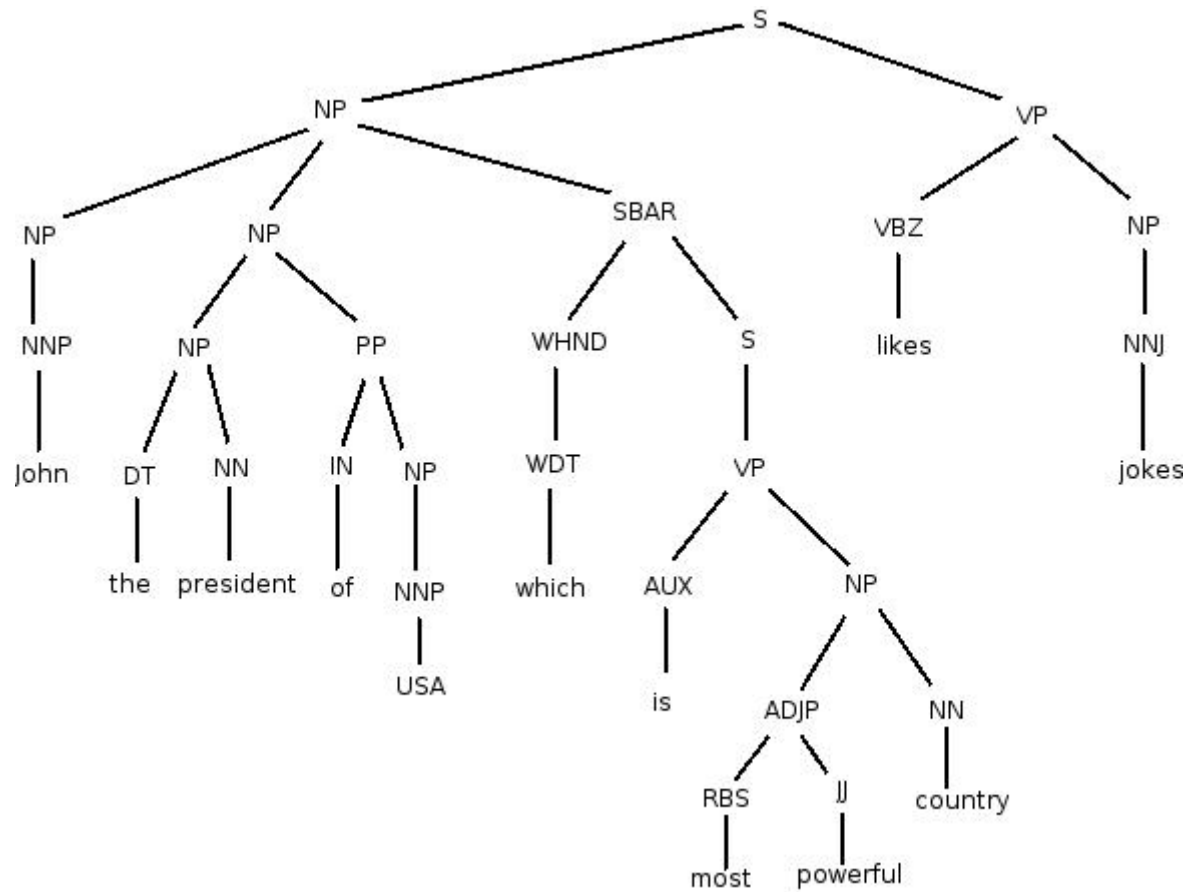


The cat that killed the rat that stole the milk that spilled on the floor that was slippery escaped.

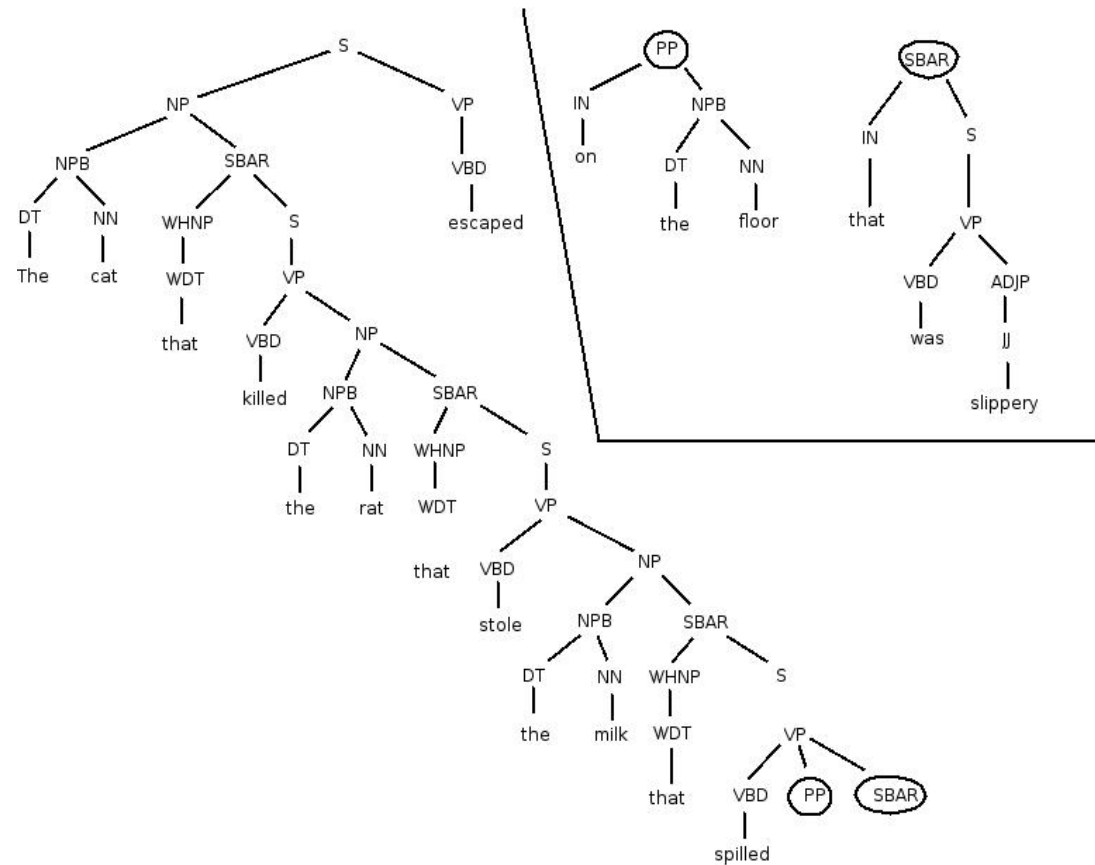


A

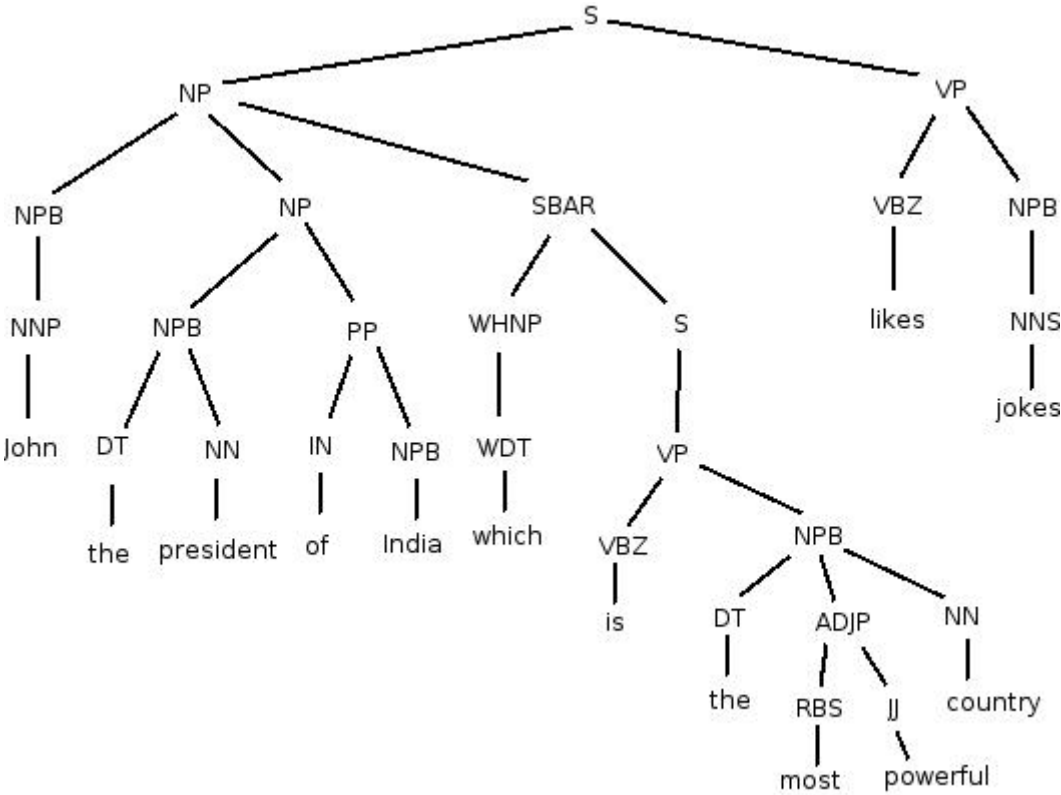
# Charniak (Embedding 2)



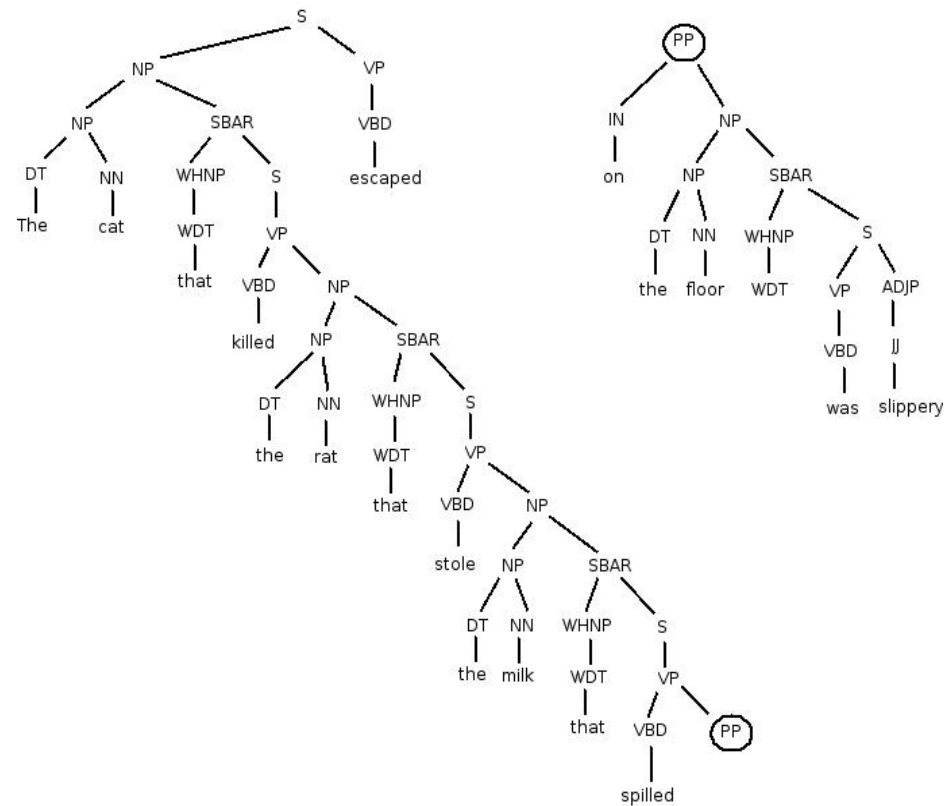
# Collins (Embedding 1)



# Collins (Embedding 2)



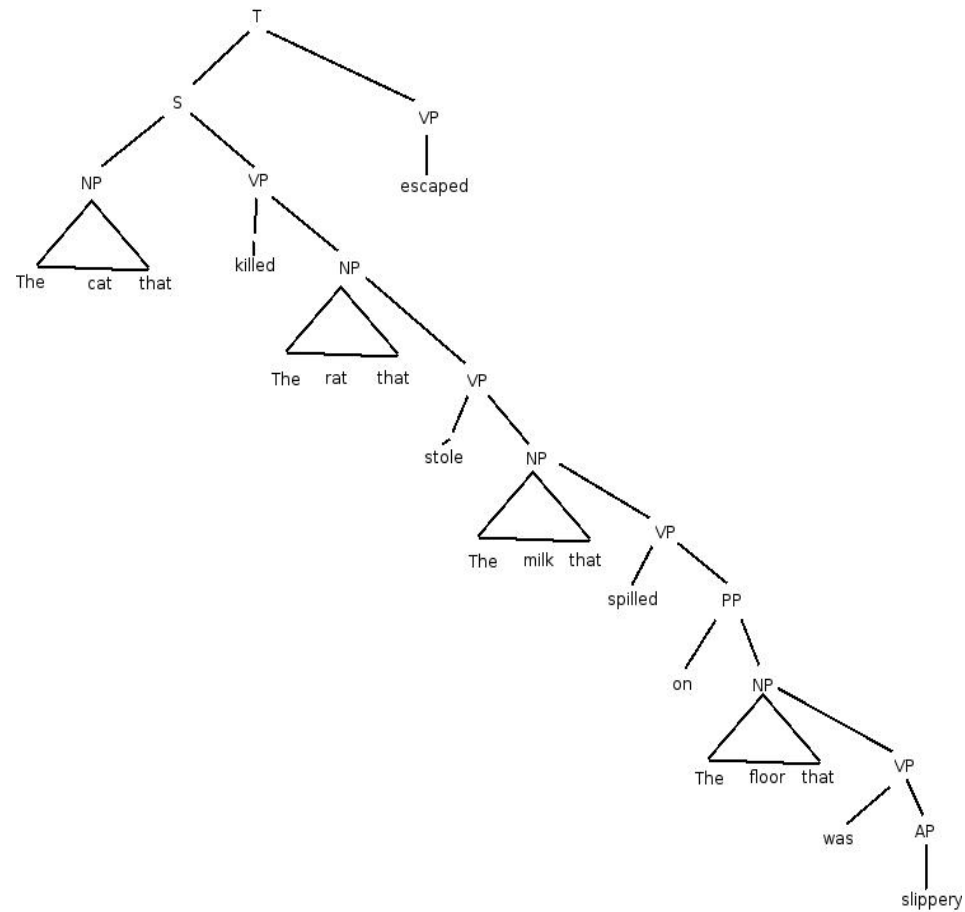
# Stanford (Embedding 1)



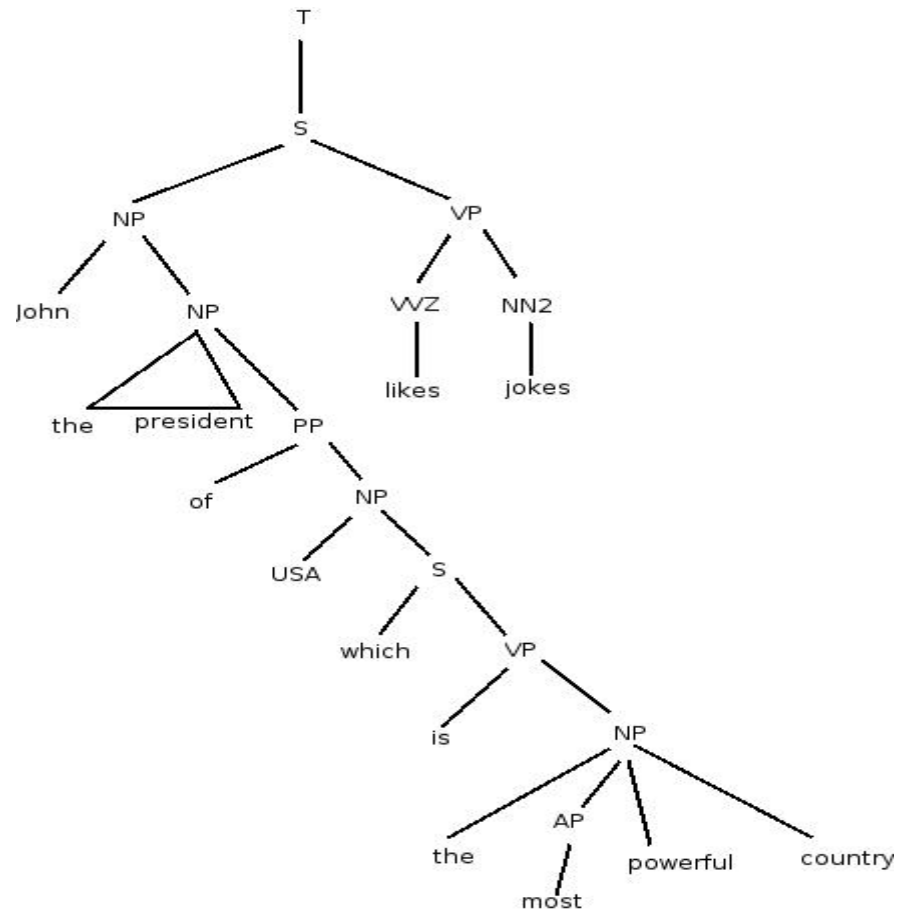




# RASP (Embedding 1)



# RASP (Embedding 2)

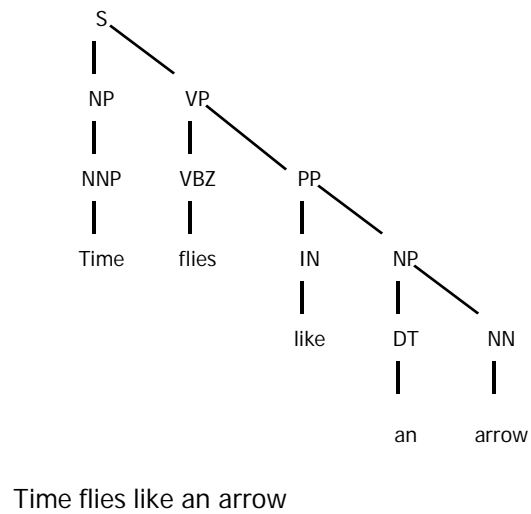


# Observation

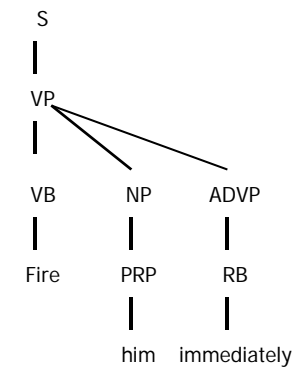
- For the sentence 'The cat that killed the rat that stole the milk that spilled on the floor that was slippery escaped.' all the parsers give the correct results.
- For the sentence 'John the president of USA which is the most powerful country likes jokes': RASP, Charniak and Collins give correct parse, *i.e.*, it attaches the verb phrase 'likes jokes' to the top NP 'John'.
- Stanford produces incorrect parse tree; attaches the VP 'likes' to wrong NP 'the president of ...'

Handling multiple POS tags

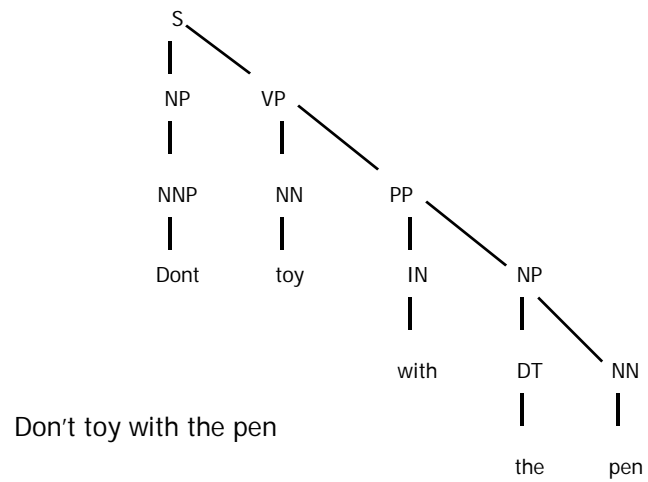
# Charniak (multiple pos 1)



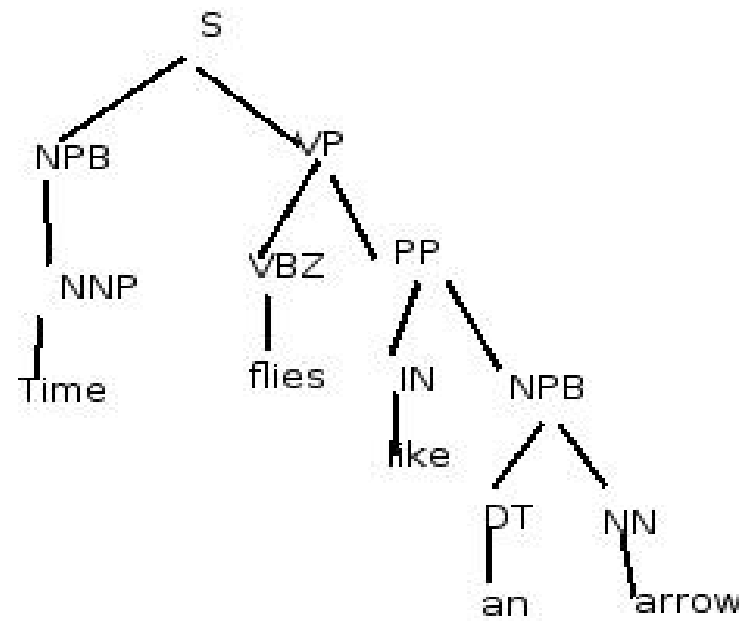
Fire him immediately



# Charniak (multiple pos 2)

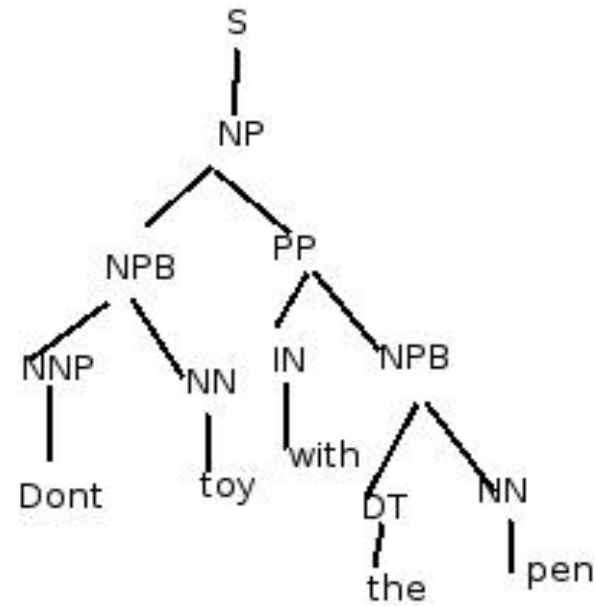


# Collins (multiple pos 1)

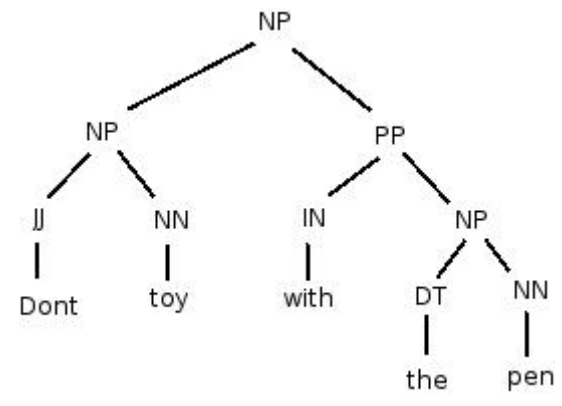
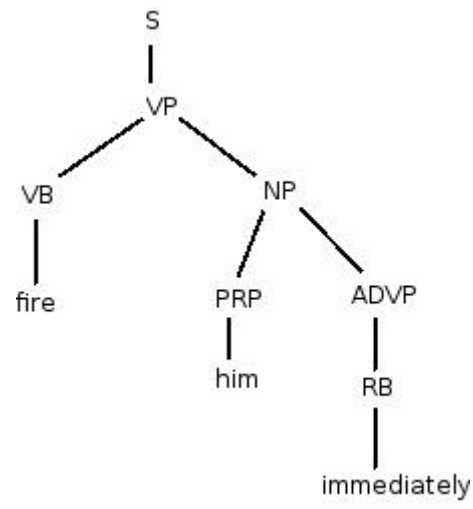




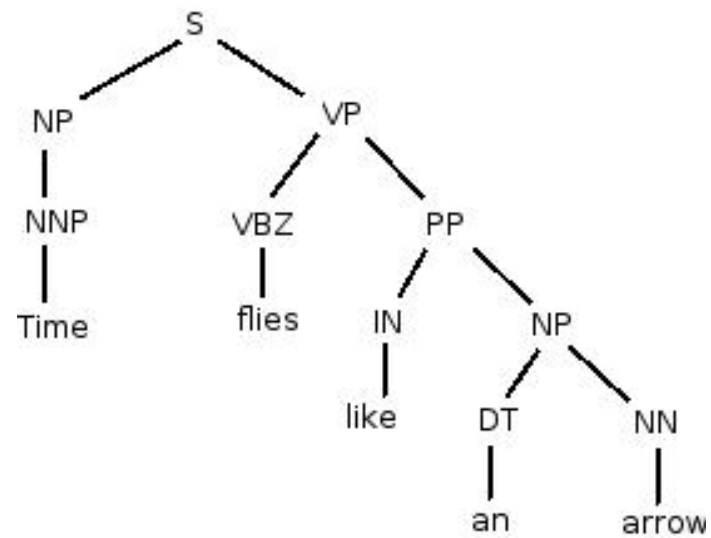
# Collins (multiple pos 2)



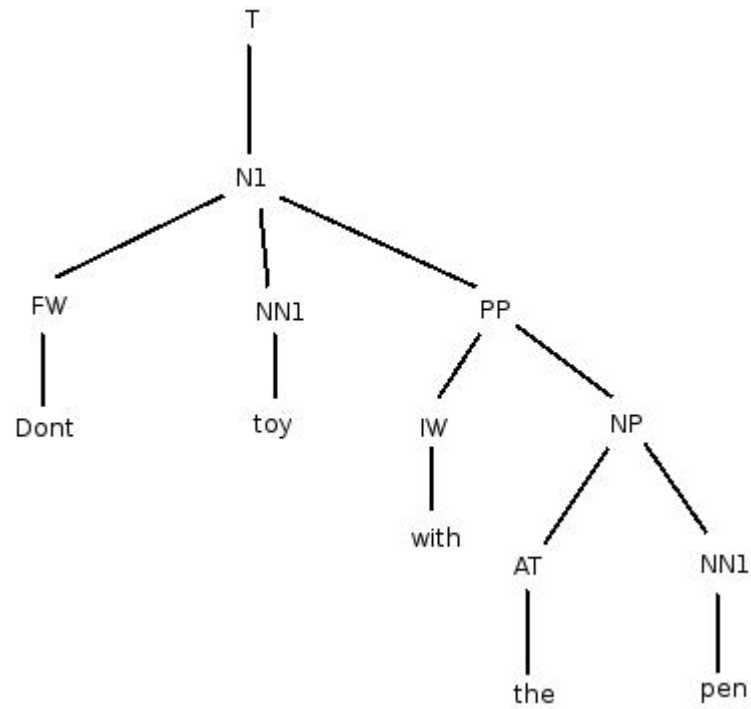
# Stanford (multiple pos 1)



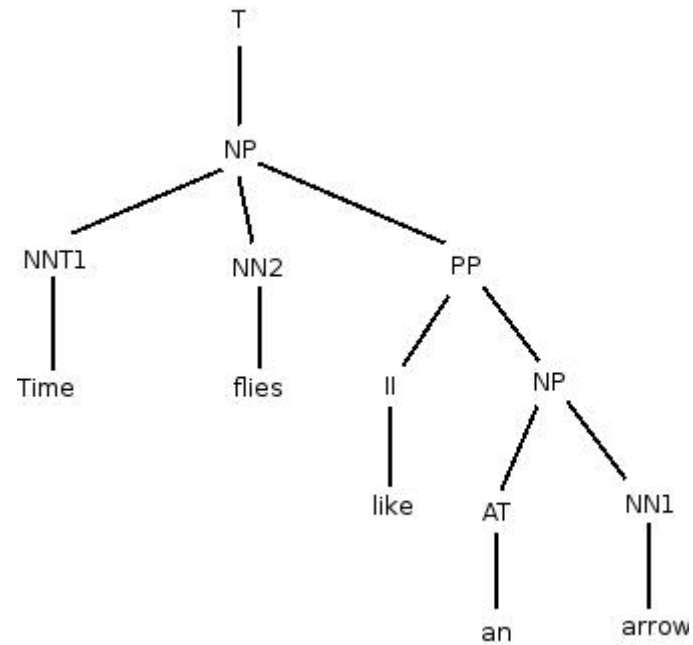
# Stanford (multiple pos 2)



# RASP (multiple pos 1)



# RASP (multiple pos 2)

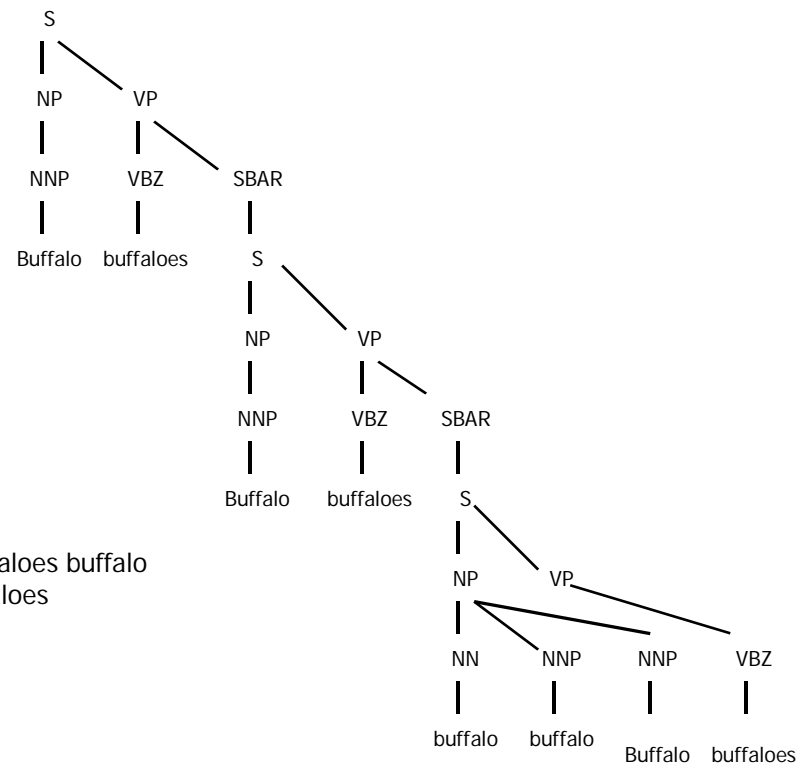


# Observation

- All but RASP give comparable pos tags. In the sentence 'Time flies like an arrow' RASP give flies as noun.
- In sentence 'Don't toy with the pen', all parsers are tagging 'toy' as noun.

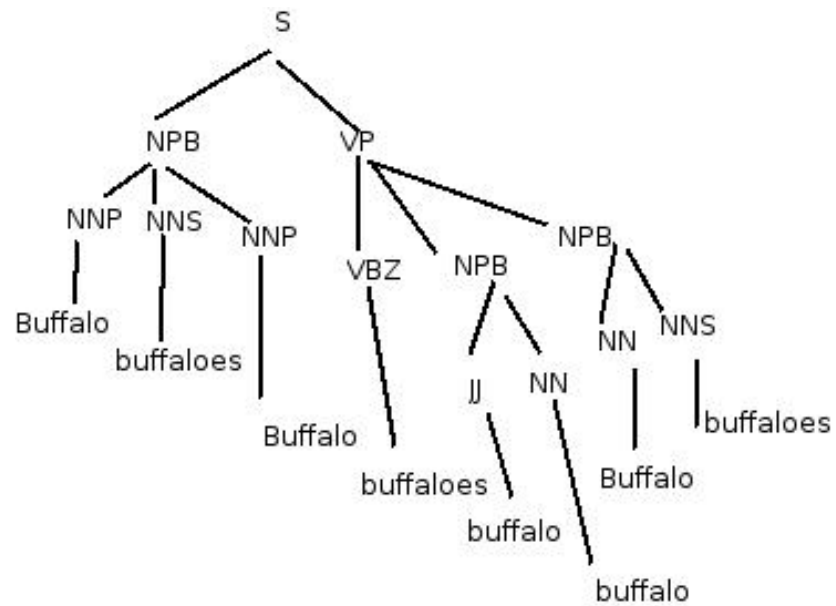
Repeated Word handling

# Charniak

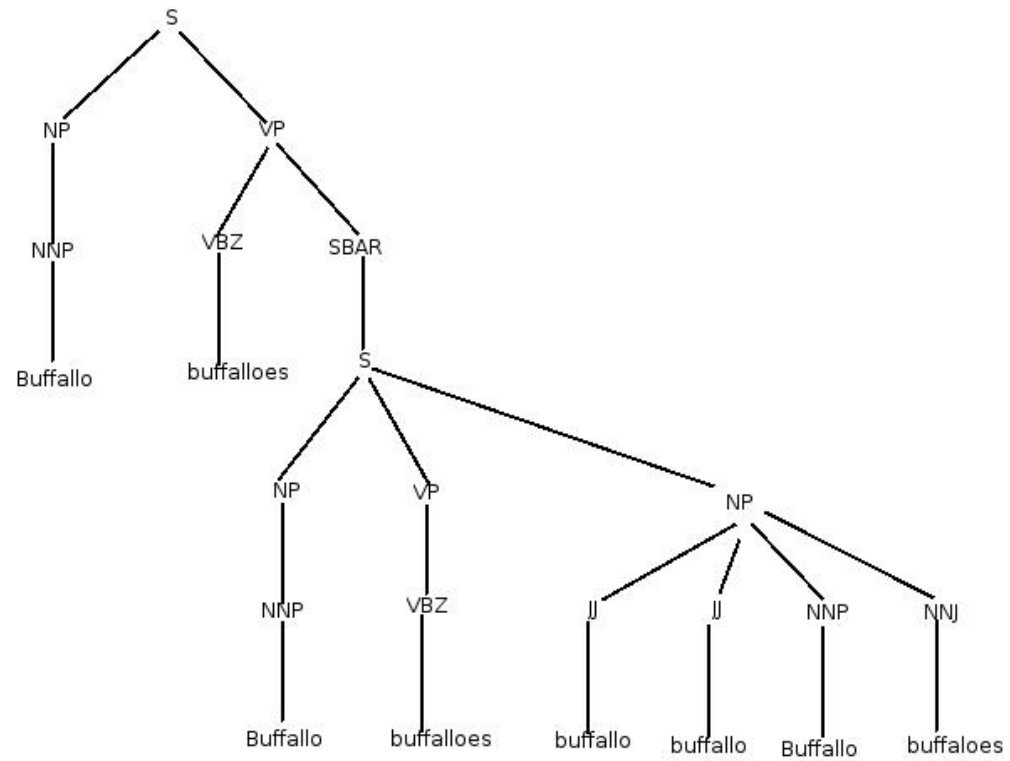




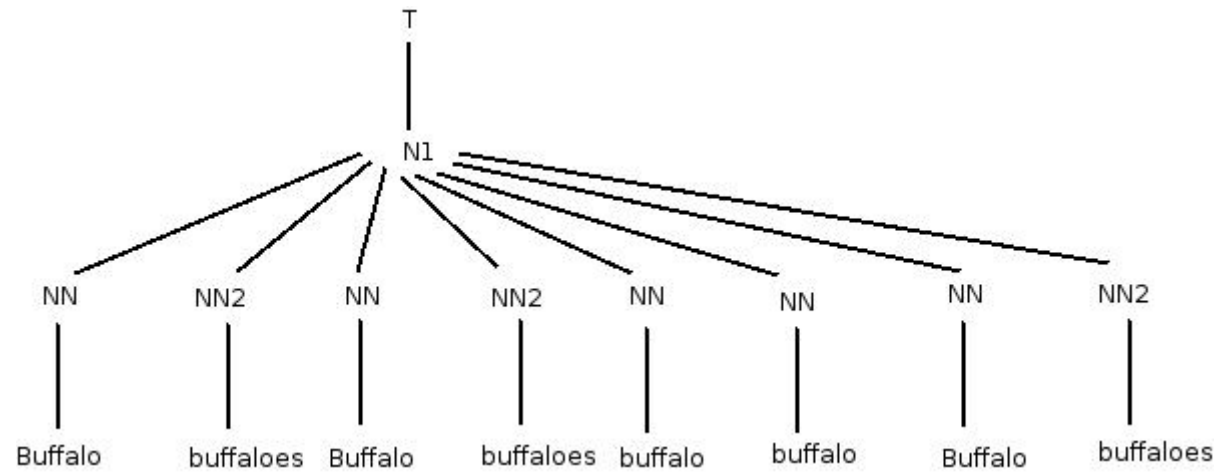
# Collins



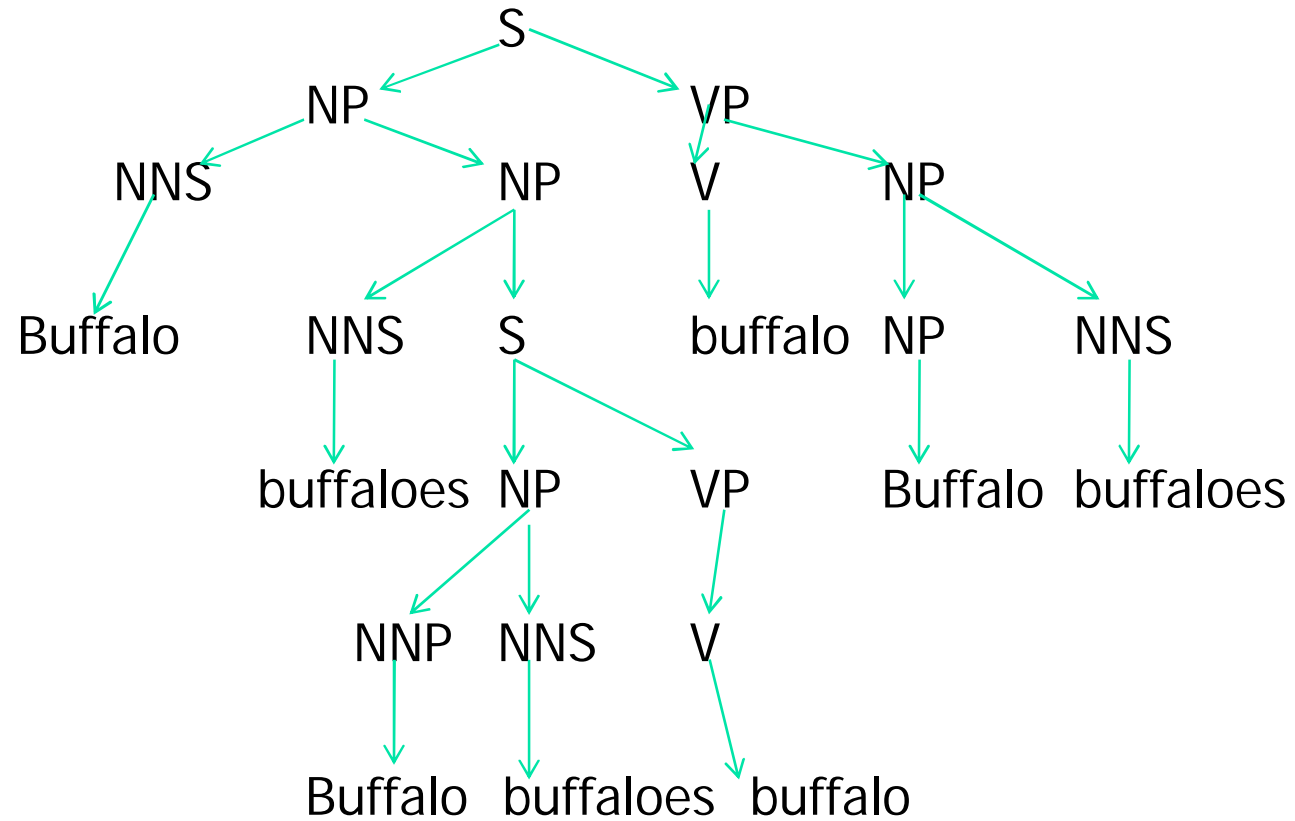
# Stanford



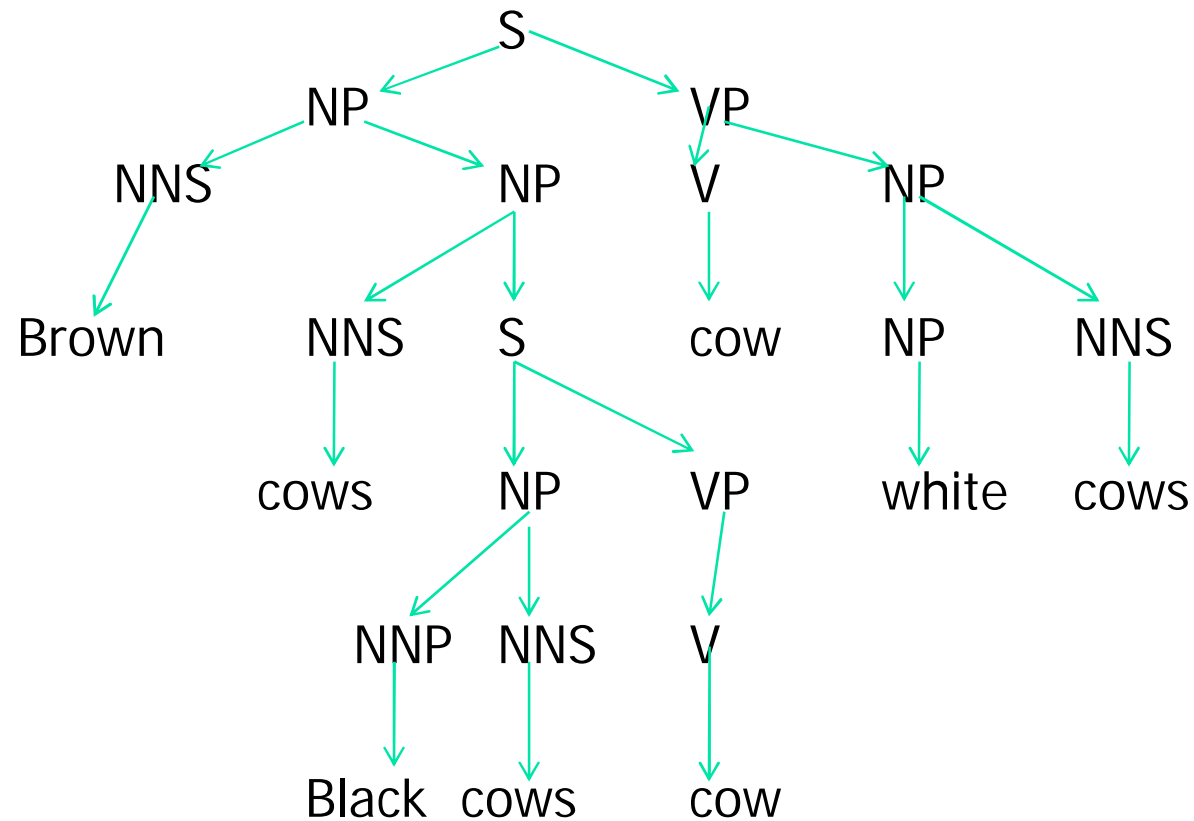
# RASP



# Correct parse



# Another sentence of same structure



# Observation

- Collins and Charniak come close to producing the correct parse.
- RASP tags all the words as nouns.

# Long sentences

- Given a sentence of 394 words, only RASP was able to parse.

# Lengthy sentence

- One day, Sam left his small, yellow home to head towards the meat-packing plant where he worked, a task which was never completed, as on his way, he tripped, fell, and went careening off of a cliff, landing on and destroying Max, who, incidentally, was also heading to his job at the meat-packing plant, though not the same plant at which Sam worked, which he would be heading to, if he had been aware that that the plant he was currently heading towards had been destroyed just this morning by a mysterious figure clad in black, who hailed from the small, remote country of France, and who took every opportunity he could to destroy small meat-packing plants, due to the fact that as a child, he was tormented, and frightened, and beaten savagely by a family of meat-packing plants who lived next door, and scarred his little mind to the point where he became a twisted and sadistic creature, capable of anything, but specifically capable of destroying meat-packing plants, which he did, and did quite often, much to the chagrin of the people who worked there, such as Max, who was not feeling quite so much chagrin as most others would feel at this point, because he was dead as a result of an individual named Sam, who worked at a competing meat-packing plant, which was no longer a competing plant, because the plant that it would be competing against was, as has already been mentioned, destroyed in, as has not quite yet been mentioned, a massive, mushroom cloud of an explosion, resulting from a heretofore unmentioned horse manure bomb manufactured from manure harvested from the farm of one farmer J. P. Harvenkirk, and more specifically harvested from a large, ungainly, incontinent horse named Seabiscuit, who really wasn't named Seabiscuit, but was actually named Harold, and it completely baffled him why anyone, particularly the author of a very long sentence, would call him Seabiscuit; actually, it didn't baffle him, as he was just a stupid, manure-making horse, who was incapable of cognitive thought for a variety of reasons, one of which was that he was a horse, and the other of which was that he was just knocked unconscious by a flying chunk of a meat-packing plant, which had been blown to pieces just a few moments ago by a shifty character from France.



# Partial RASP Parse of the sentence

- ([One\_MC1 | day\_NNT1 | | | Sam\_NP1 | leave+ed\_VVD | his\_APP\$ | small\_JJ | | yellow\_JJ | home\_NN1 | to\_TO | head\_VV0 | towards\_II | the\_AT | meat-packing\_JJ | plant\_NN1 | where\_RRQ | he\_PPHS1 | work+ed\_VVD | | | a\_AT1 | task\_NN1 | which\_DDO | be+ed\_VBDZ | never\_RR | complete+ed\_VVN | | | as\_CSA | on\_II | his\_APP\$ | way\_NN1 | | | he\_PPHS1 | trip+ed\_VVD | | | fall+ed\_VVD | | | and\_CC | go+ed\_VVD | careen+ing\_VVG | off\_RP | of\_IO | a\_AT1 | cliff\_NN1 | | | land+ing\_VVG | on\_RP | and\_CC | destroy+ing\_VVG | Max\_NP1 | | | who\_PNQS | | | incidentally\_RR | | | be+ed\_VBDZ | also\_RR | head+ing\_VVG | to\_II | his\_APP\$ | job\_NN1 | at\_II | the\_AT | meat-packing\_JB | plant\_NN1 | | | though\_CS | not+\_XX | the\_AT | same\_DA | plant\_NN1 | at\_II | which\_DDO | Sam\_NP1 | work+ed\_VVD | | | which\_DDO | he\_PPHS1 | would\_VM | be\_VB0 | head+ing\_VVG | to\_II | | | if\_CS | he\_PPHS1 | have+ed\_VHD | be+en\_VBN | aware\_JJ | that\_CST | that\_CST | the\_AT | plant\_NN1 | he\_PPHS1 | be+ed\_VBDZ | currently\_RR | head+ing\_VVG | towards\_II | have+ed\_VHD | be+en\_VBN | destroy+ed\_VVN | just\_RR | this\_DD1 | morning\_NNT1 | by\_II | a\_AT1 | mysterious\_JJ | figure\_NN1 | clothe+ed\_VVN | in\_II | black\_JJ | | | who\_PNQS | hail+ed\_VVD | from\_II | the\_AT | small\_JJ | | | remote\_JJ | country\_NN1 | of\_IO | France\_NP1 | | | and\_CC | who\_PNQS | take+ed\_VVD | every\_AT1 | opportunity\_NN1 | he\_PPHS1 | could\_VM | to\_TO | destroy\_VV0 | small\_JJ | meat-packing\_NN1 | plant+s\_NN2 | | | due\_JJ | to\_II | the\_AT | fact\_NN1 | that\_CST | as\_CSA | a\_AT1 | child\_NN1 | | | he\_PPHS1 | be+ed\_VBDZ | torment+ed\_VVN | | | and\_CC | frighten+ed\_VVD | | | and\_CC | beat+en\_VVN | savagely\_RR | by\_II | a\_AT1 | family\_NN1 | of\_IO | meat-packing\_JJ | plant+s\_NN2 | who\_PNQS | live+ed\_VVD | next\_MD | door\_NN1 | | | and\_CC | scar+ed\_VVD | his\_APP\$ | little\_DD1 | mind\_NN1 | to\_II | the\_AT | point\_NNL1 | where\_RRQ | he\_PPHS1 | become+ed\_VVD | a\_AT1 | twist+ed\_VVN | and\_CC | sadistic\_JJ | creature\_NN1 | | | capable\_JJ | of\_IO | anything\_PN1 | | | but\_CCB | specifically\_RR | capable\_JJ | of\_IO | destroy+ing\_VVG | meat-packing\_JJ | plant+s\_NN2 | | | which\_DDO | he\_PPHS1 | do+ed\_VDD | | | and\_CC | do+ed\_VDD | quite\_RG | often\_RR | | | much\_DA1 | to\_II | the\_AT | chagrin\_NN1 | of\_IO | the\_AT | people\_NN | who\_PNQS | work+ed\_VVD | there\_RL | | | such\_DA | as\_CSA | Max\_NP1 | | | who\_PNQS | be+ed\_VBDZ | not+\_XX | feel+ing\_VVG | quite\_RG | so\_RG | much\_DA1 | chagrin\_NN1 | as\_CSA | most\_DAT | other+s\_NN2 | would\_VM | feel\_VV0 | at\_II | this\_DD1 | point\_NNL1 | | | because\_CS | he\_PPHS1 | be+ed\_VBDZ | dead\_JJ | as\_CSA | a\_AT1 | result\_NN1 | of\_IO | an\_AT1 | individual\_NN1 | name+ed\_VVN | Sam\_NP1 | | | who\_PNQS | work+ed\_VVD | at\_II | a\_AT1 | compete+ing\_VVG | meat-packing\_JJ | plant\_NN1 | | | which\_DDO | be+ed\_VBDZ | no\_AT | longer\_RRR | a\_AT1 | compete+ing\_VVG | plant\_NN1 | | | because\_CS | the\_AT | plant\_NN1 | that\_CST | it\_PPH1 | would\_VM | be\_VB0 | compete+ing\_VVG | against\_II | be+ed\_VBDZ | | | as\_CSA | have+s\_VHZ | already\_RR | be+en\_VBN | mention+ed\_VVN | | | destroy+ed\_VVN | in\_RP | | | as\_CSA | have+s\_VHZ | not+\_XX | quite\_RG | yet\_RR | be+en\_VBN | mention+ed\_VVN | | | a\_AT1 | massive\_JJ | | | mushroom\_NN1 | cloud\_NN1 | of\_IO | an\_AT1 | explosion\_NN1 | | | result+ing\_VVG | from\_II | a\_AT1 | heretofore\_RR | unmentioned\_JJ | horse\_NN1 | manure\_NN1 | bomb\_NN1 | manufacture+ed\_VVN | from\_II | manure\_NN1 | harvest+ed\_VVN | from\_II | the\_AT | farm\_NN1 | of\_IO | one\_MC1 | farmer\_NN1 | J\_NP1 P\_NP1 Harvenkirk\_NP1 | | | and\_CC | more\_DAR | specifically\_RR | harvest+ed\_VVN | from\_II | a\_AT1 | large\_JJ | | | ungainly\_JJ | | | incontinent\_NN1 | horse\_NN1 | name+ed\_VVN | Seabiscuit\_NP1 | | | who\_PNQS | really\_RR | be+ed\_VBDZ | not+\_XX | name+ed\_VVN | Seabiscuit\_NP1 | | | but\_CCB | be+ed\_VBDZ | actually\_RR | name+ed\_VVN | Harold\_NP1 | | | and\_CC | it\_PPH1 | completely\_RR | baffle+ed\_VVD | he+\_PPHO1 | why\_RRQ | anyone\_PN1 | | | particularly\_RR | the\_AT | author\_NN1 | of\_IO | a\_AT1 | very\_RG | long\_JJ | sentence\_NN1 | | | would\_VM | call\_VV0 | he+\_PPHO1 | Seabiscuit\_NP1 | | | actually\_RR | | | it\_PPH1 | do+ed\_VDD | not+\_XX | baffle\_VV0 | he+\_PPHO1 | | | as\_CSA | he\_PPHS1 | be+ed\_VBDZ | just\_RR | a\_AT1 | stupid\_JJ | | | manure-making\_NN1 | horse\_NN1 | | | who\_PNQS | be+ed\_VBDZ | incapable\_JJ | of\_IO | cognitive\_JJ | thought\_NN1 | for\_IF | a\_AT1 | variety\_NN1 | of\_IO | reason+s\_NN2 | | | one\_MC1 | of\_IO | which\_DDO | be+ed\_VBDZ | that\_CST | he\_PPHS1 | be+ed\_VBDZ | a\_AT1 | horse\_NN1 | | | and\_CC | the\_AT | other\_JB | of\_IO | which\_DDO | be+ed\_VBDZ | that\_CST | he\_PPHS1 | be+ed\_VBDZ | just\_RR | knock+ed\_VVN | unconscious\_JJ | by\_II | a\_AT1 | flying\_NN1 | chunk\_NN1 | of\_IO | a\_AT1 | meat-packing\_JJ | plant\_NN1 | | | which\_DDO | have+ed\_VHD | be+en\_VBN | blow+en\_VVN | to\_II | piece+s\_NN2 | just\_RR | a\_AT1 | few\_DA2 | moment+s\_NNT2 | ago\_RA | by\_II | a\_AT1 | shifty\_JJ | character\_NN1 | from\_II | France\_NP1 | | | -1 ; 0

An important parsing algo

# Illustrating CYK [Cocke, Younger, Kashmi] Algo

- $S \rightarrow NP VP$       1.0
- $NP \rightarrow DT NN$       0.5
- $NP \rightarrow NNS$       0.3
- $NP \rightarrow NP PP$       0.2
- $PP \rightarrow P NP$       1.0
- $VP \rightarrow VP PP$       0.6
- $VP \rightarrow VBD NP$       0.4
- $DT \rightarrow the$       1.0
- $NN \rightarrow gunman$       0.5
- $NN \rightarrow building$       0.5
- $VBD \rightarrow sprayed$       1.0
- $NNS \rightarrow bullets$       1.0

# CYK: Start with (0,1)

o *The* 1 *gunman* 2 *sprayed* 3 *the* 4 *building* 5 *with* 6 *bullets* 7.

To From	1	2	3	4	5	6	7
0	DT						
1	→ -----						
2	↓ -----	----- --					
3	-----	----- --	----- -				
4	----- -	----- --	----- -	----- -			
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK: Keep filling diagonals

o The 1 *gunman* 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0	DT						
1 →	-----	NN					
2 ↓	-----	----- --					
3	-----	----- --	----- -				
4	----- -	----- --	----- -	----- -			
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK: Try getting higher level structures

o *The* 1 *gunman* 2 *sprayed* 3 *the* 4 *building* 5 *with* 6 *bullets* 7.

To From	1	2	3	4	5	6	7
0 →	DT	NP					
1	-----	NN					
2 ↓	-----	----- --					
3	-----	----- --	----- -				
4	----- -	----- --	----- -	----- -			
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK: Diagonal continues

o The 1 gunman 2 *sprayed* 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0	DT	NP					
1 →	-----	NN					
2 ↓	-----	----- --	VBD				
3	-----	----- --	----- -				
4	----- -	----- --	----- -	----- -			
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK (cont...)

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0 →	DT	NP	----- -				
1 ↓	-----	NN	----- -				
2	-----	----- --	VBD				
3	-----	----- --	----- -				
4	----- -	----- --	----- -	----- -			
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	



# CYK (cont...)

o The 1 gunman 2 sprayed 3 *the* 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0	DT	NP	----- -				
1 <span style="margin-left: 20px;">→</span> <span style="margin-left: 20px;">↓</span>	-----	NN	----- -				
2	-----	----- --	VBD				
3	-----	----- --	----- -	DT			
4	----- -	----- --	----- -	----- -			
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK (cont...)

o The 1 gunman 2 sprayed 3 the 4 *building* 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0 →	DT	NP	----- -	----- -			
1 ↓	-----	NN	----- -	----- -			
2	-----	----- --	VBD	----- -			
3	-----	----- --	----- -	DT			
4	----- -	----- --	----- -	----- -	NN		
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK: starts filling the 5<sup>th</sup> column

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0 →	DT	NP	----- -	----- -			
1 ↓	-----	NN	----- -	----- -			
2	-----	----- --	VBD	----- -			
3	-----	----- --	----- -	DT	NP		
4	----- -	----- --	----- -	----- -	NN		
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK (cont...)

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0	DT	NP	----- -	----- -			
1	----- ↓	NN	----- -	----- -			
2	-----	----- --	VBD	----- -	VP		
3	-----	----- --	----- -	DT	NP		
4	----- -	----- --	----- -	----- -	NN		
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK (cont...)

o The 1 *gunman* 2 *sprayed* 3 *the* 4 *building* 5 *with* 6 *bullets* 7.

To From	1	2	3	4	5	6	7
0	DT	NP	----- -	----- -			
1	----- ↓	NN	----- -	----- -	----- -		
2	-----	----- --	VBD	----- -	VP		
3	-----	----- --	----- -	DT	NP		
4	----- -	----- --	----- -	----- -	NN		
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK: S found, but NO termination!

o *The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.*

To From	1	2	3	4	5	6	7
0	DT	NP	----- -	----- -	S		
1	----- ↓	NN	----- -	----- -	----- -		
2	-----	----- --	VBD	----- -	VP		
3	-----	----- --	----- -	DT	NP		
4	----- -	----- --	----- -	----- -	NN		
5	----- -	----- --	----- -	----- -	----- -		
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK (cont...)

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0	DT	NP	----- -	----- -	S		
1	----- ↓	NN	----- -	----- -	----- -		
2	-----	----- --	VBD	----- -	VP		
3	-----	----- --	----- -	DT	NP		
4	----- -	----- --	----- -	----- -	NN		
5	----- -	----- --	----- -	----- -	----- -	P	
6	----- -	----- --	----- -	----- -	----- -	----- -	

# CYK (cont...)

- The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0 →	DT	NP	----- -	----- -	S	----- -	
1 ↓	-----	NN	----- -	----- -	----- -	----- -	
2	-----	----- --	VBD	----- -	VP	----- -	
3	-----	----- --	----- -	DT	NP	----- -	
4	----- -	----- --	----- -	----- -	NN	----- -	
5	----- -	----- --	----- -	----- -	----- -	P	
6	----- -	----- --	----- -	----- -	----- -	----- -	



# CYK: Control moves to last column

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0	DT	NP	----- -	----- -	S	----- -	
1 →	-----	NN	----- -	----- -	----- -	----- -	
↓ 2	-----	----- --	VBD	----- -	VP	----- -	
3	-----	----- --	----- -	DT	NP	----- -	
4	----- -	----- --	----- -	----- -	NN	----- -	
5	----- -	----- --	----- -	----- -	----- -	P	
6	----- -	----- --	----- -	----- -	----- -	----- -	NP NNS

# CYK (cont...)

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0	DT	NP	----- -	----- -	S	----- -	
1	----- ↓	NN	----- -	----- -	----- -	----- -	
2	-----	----- --	VBD	----- -	VP	----- -	
3	-----	----- --	----- -	DT	NP	----- -	
4	----- -	----- --	----- -	----- -	NN	----- -	
5	----- -	----- --	----- -	----- -	----- -	P	PP
6	----- -	----- --	----- -	----- -	----- -	----- -	NP NNS

# CYK (cont...)

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0	DT	NP	----- -	----- -	S	----- -	
1 →	-----	NN	----- -	----- -	----- -	----- -	
2 ↓	-----	----- --	VBD	----- -	VP	----- -	
3	-----	----- --	----- -	DT	NP	----- -	NP
4	----- -	----- --	----- -	----- -	NN	----- -	----- -
5	----- -	----- --	----- -	----- -	----- -	P	PP
6	----- -	----- --	----- -	----- -	----- -	----- -	NP NNS

# CYK (cont...)

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0 →	DT	NP	----- -	----- -	S	----- -	
1 ↓	-----	NN	----- -	----- -	----- -	----- -	
2	-----	----- --	VBD	----- -	VP	----- -	VP
3	-----	----- --	----- -	DT	NP	----- -	NP
4	----- -	----- --	----- -	----- -	NN	----- -	----- -
5	----- -	----- --	----- -	----- -	----- -	P	PP
6	----- -	----- --	----- -	----- -	----- -	----- -	NP NNS

# CYK: filling the last column

o The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.

To From	1	2	3	4	5	6	7
0 →	DT	NP	----- -	----- -	S	----- -	
1 ↓	-----	NN	----- -	----- -	----- -	----- -	----- -
2	-----	----- --	VBD	----- -	VP	----- -	VP
3	-----	----- --	----- -	DT	NP	----- -	NP
4	----- -	----- --	----- -	----- -	NN	----- -	----- -
5	----- -	----- --	----- -	----- -	----- -	P	PP
6	----- -	----- --	----- -	----- -	----- -	----- -	NP NNS

# CYK: terminates with S in (0,7)

o *The 1 gunman 2 sprayed 3 the 4 building 5 with 6 bullets 7.*

To From	1	2	3	4	5	6	7
0	DT	NP	----- -	----- -	S	----- -	S
1	----- ↓	NN	----- -	----- -	----- -	----- -	----- -
2	-----	----- --	VBD	----- -	VP	----- -	VP
3	-----	----- --	----- -	DT	NP	----- -	NP
4	----- -	----- --	----- -	----- -	NN	----- -	----- -
5	----- -	----- --	----- -	----- -	----- -	P	PP
6	----- -	----- --	----- -	----- -	----- -	----- -	NP NNS

# CYK: Extracting the Parse Tree

- The parse tree is obtained by keeping back pointers.

