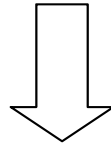


Section 5: Parsing & PCFGs

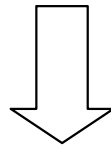
May 12, 2006
Pi-Chuan Chang

N-ary Trees in Treebank



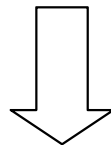
In the starter code :
binarizeTree

TreeAnnotations.annotateTree

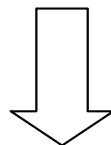


Binary Trees

TODO2:
Better vertical and horizontal
markovization



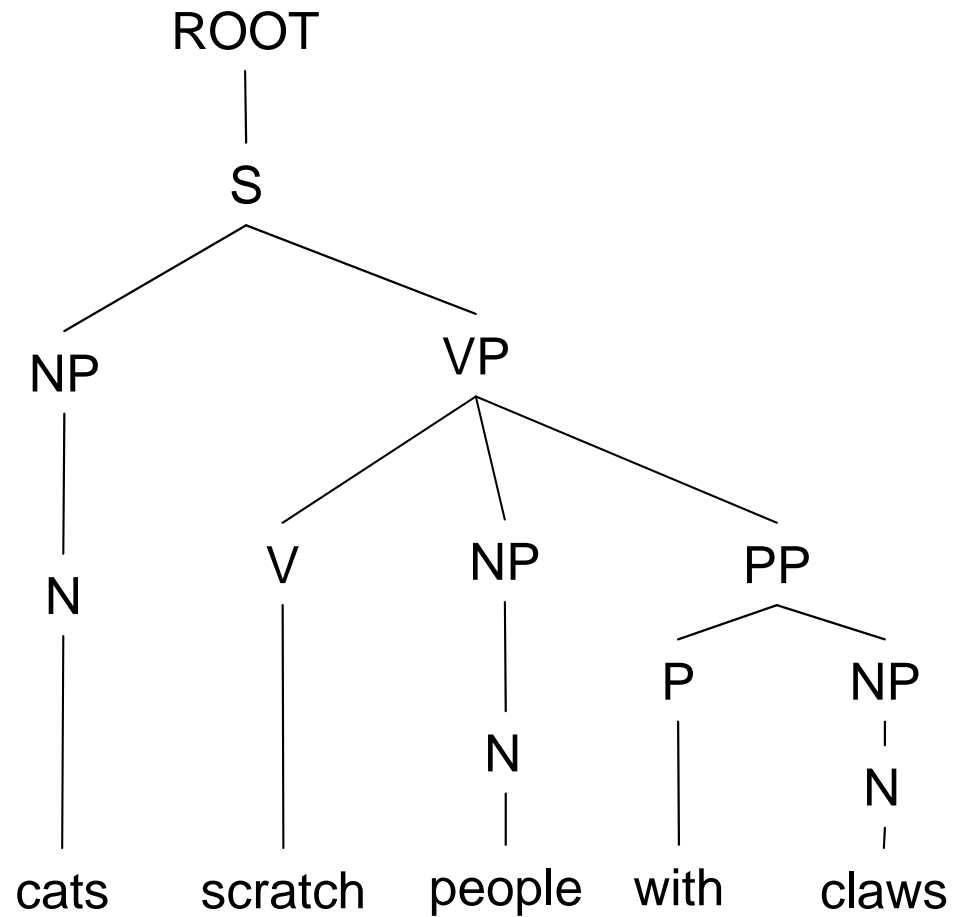
Lexicon and Grammar



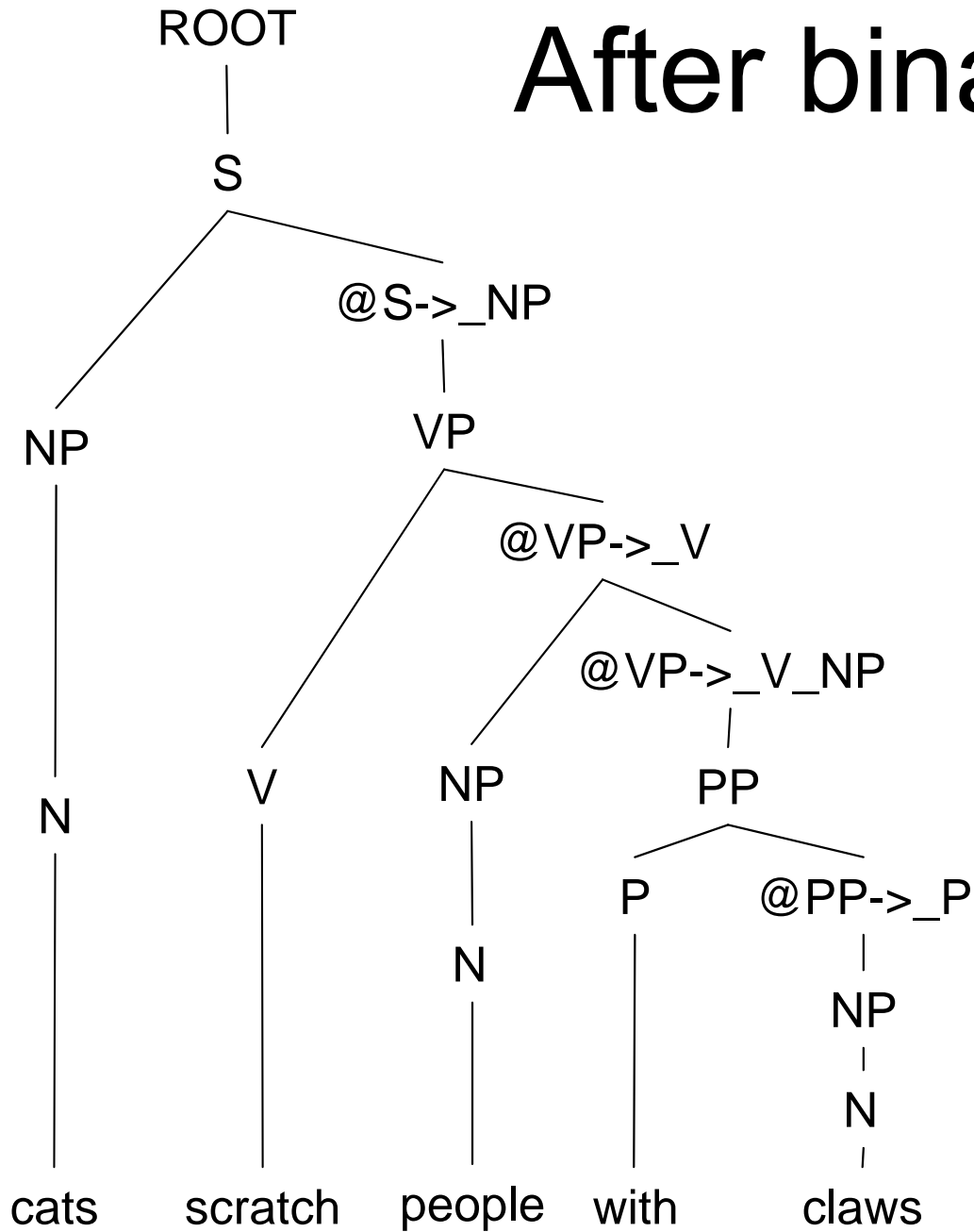
Parsing

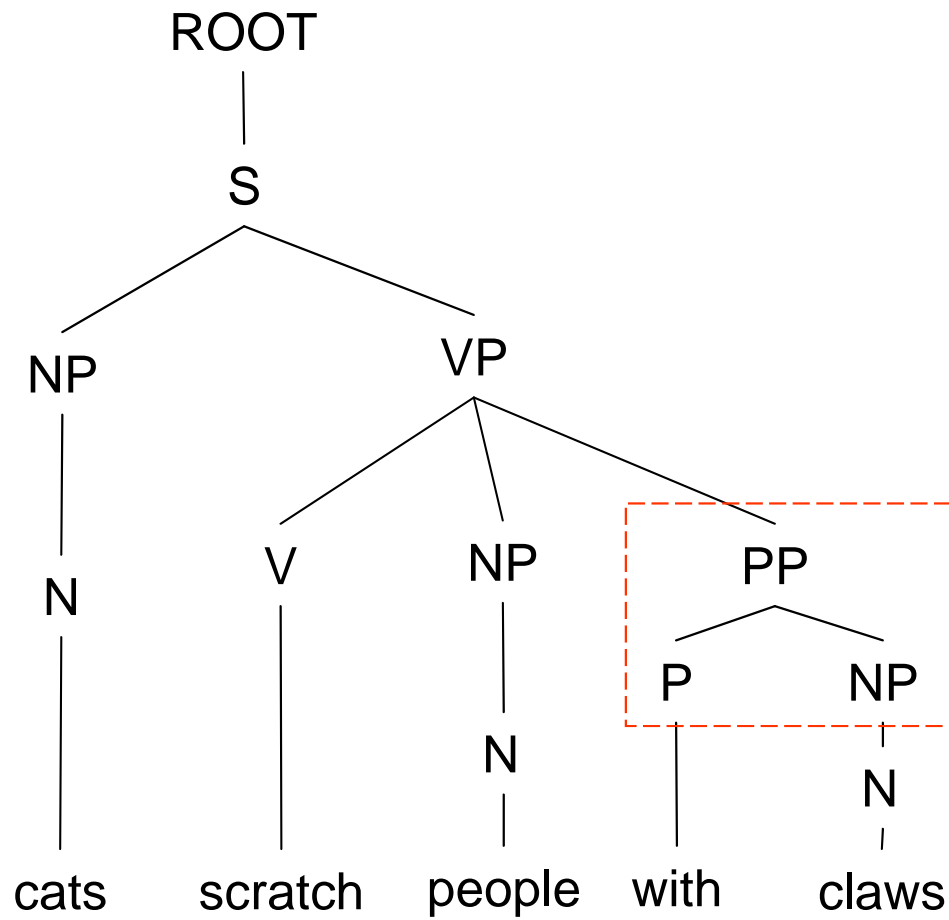
TODO1:
CKY parsing

An example: before binarization...

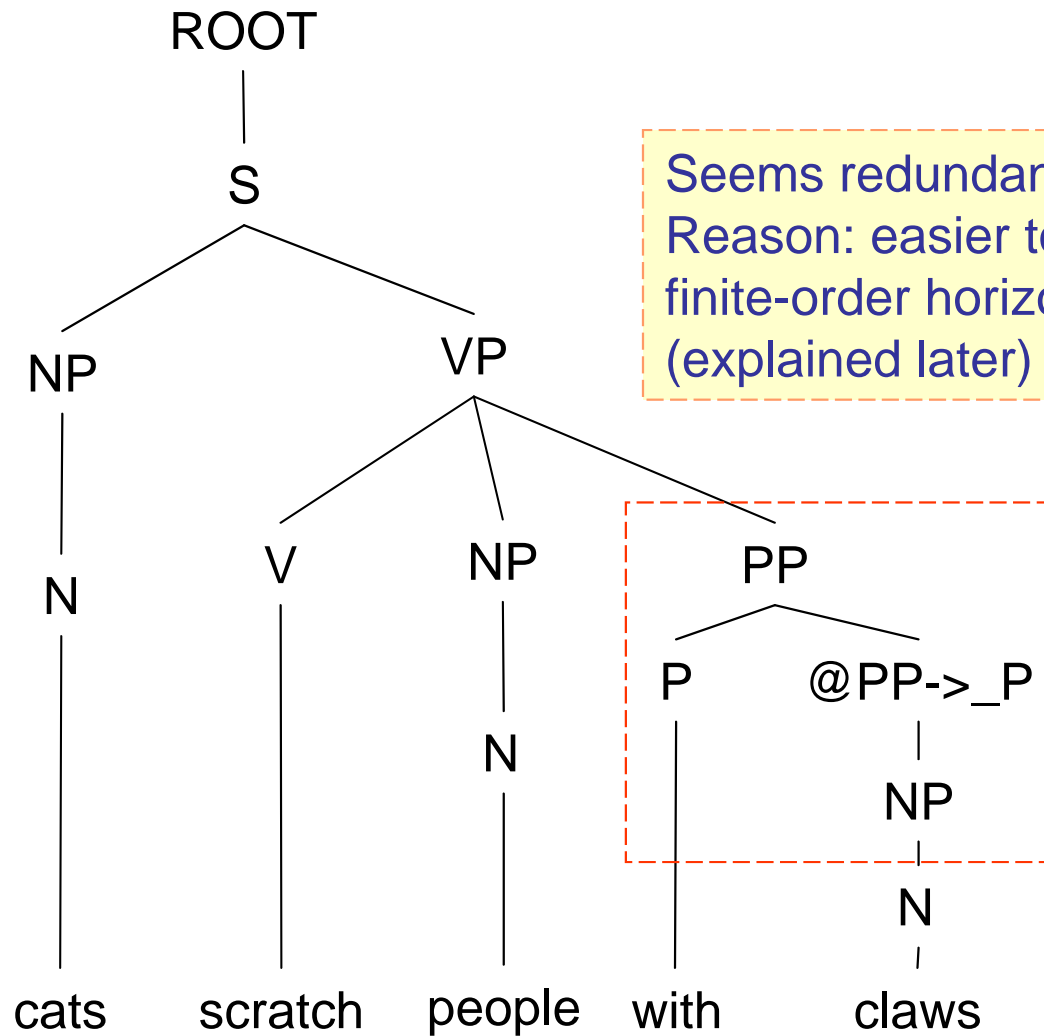


After binarization..

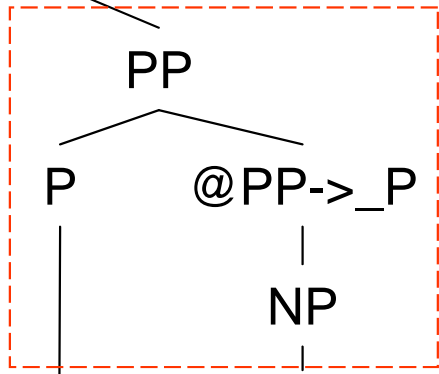


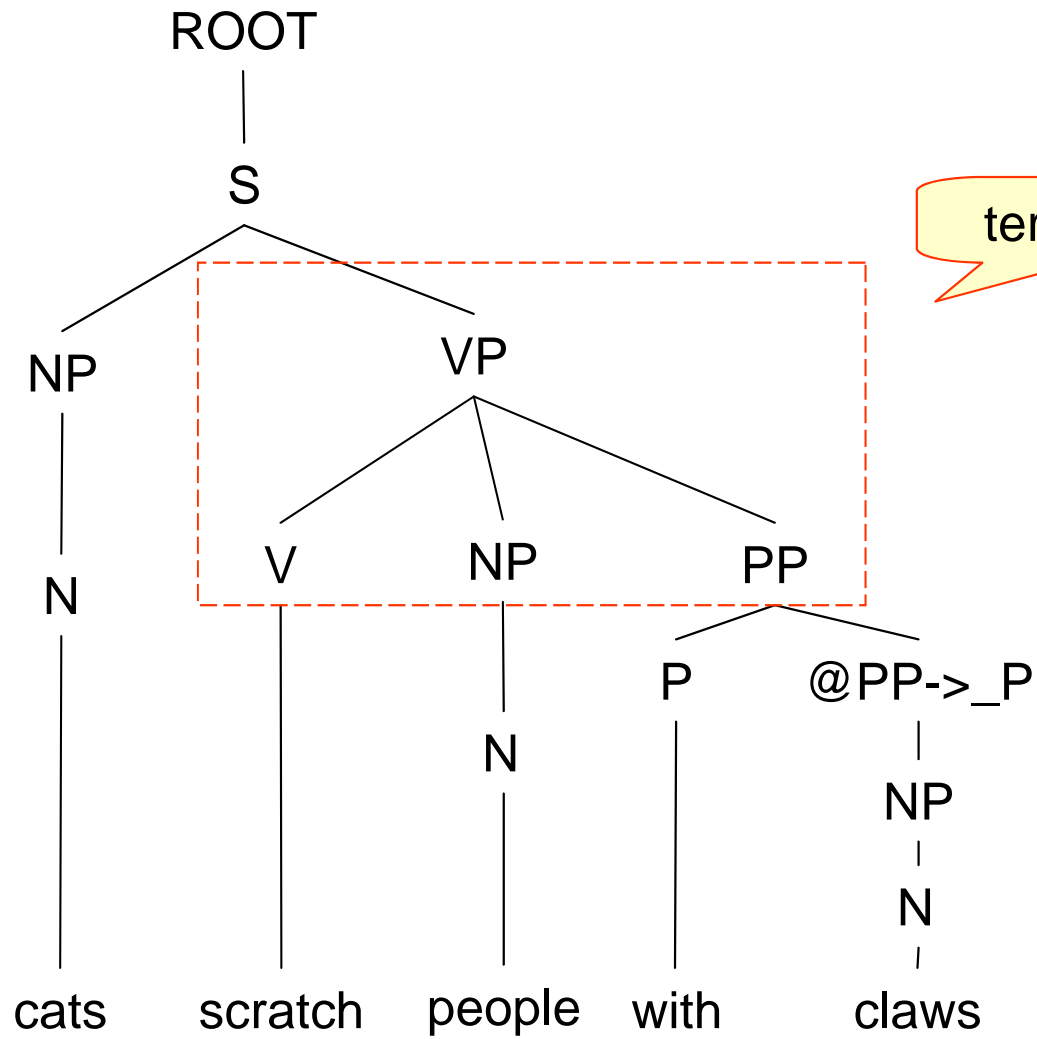


Binary rule

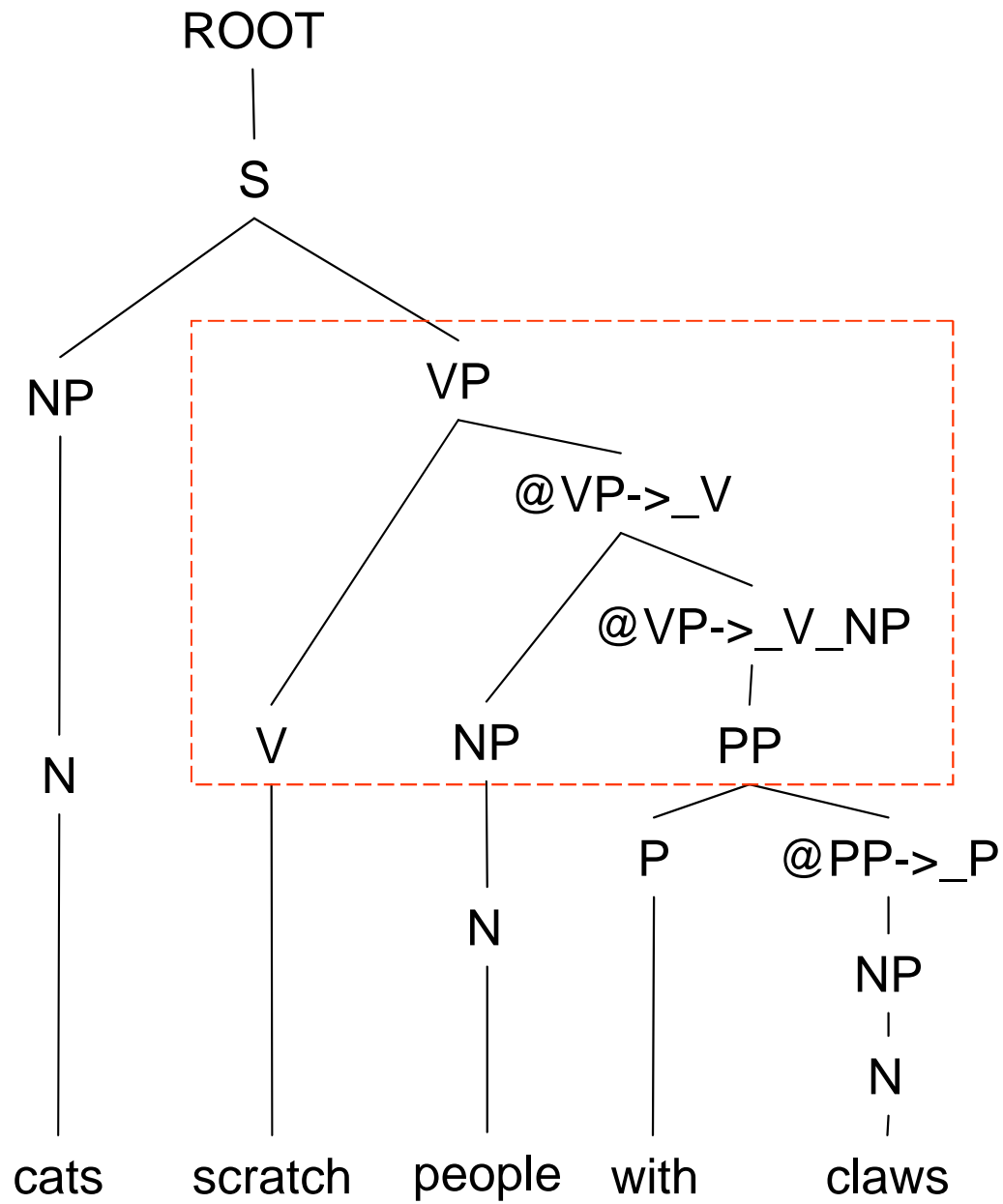


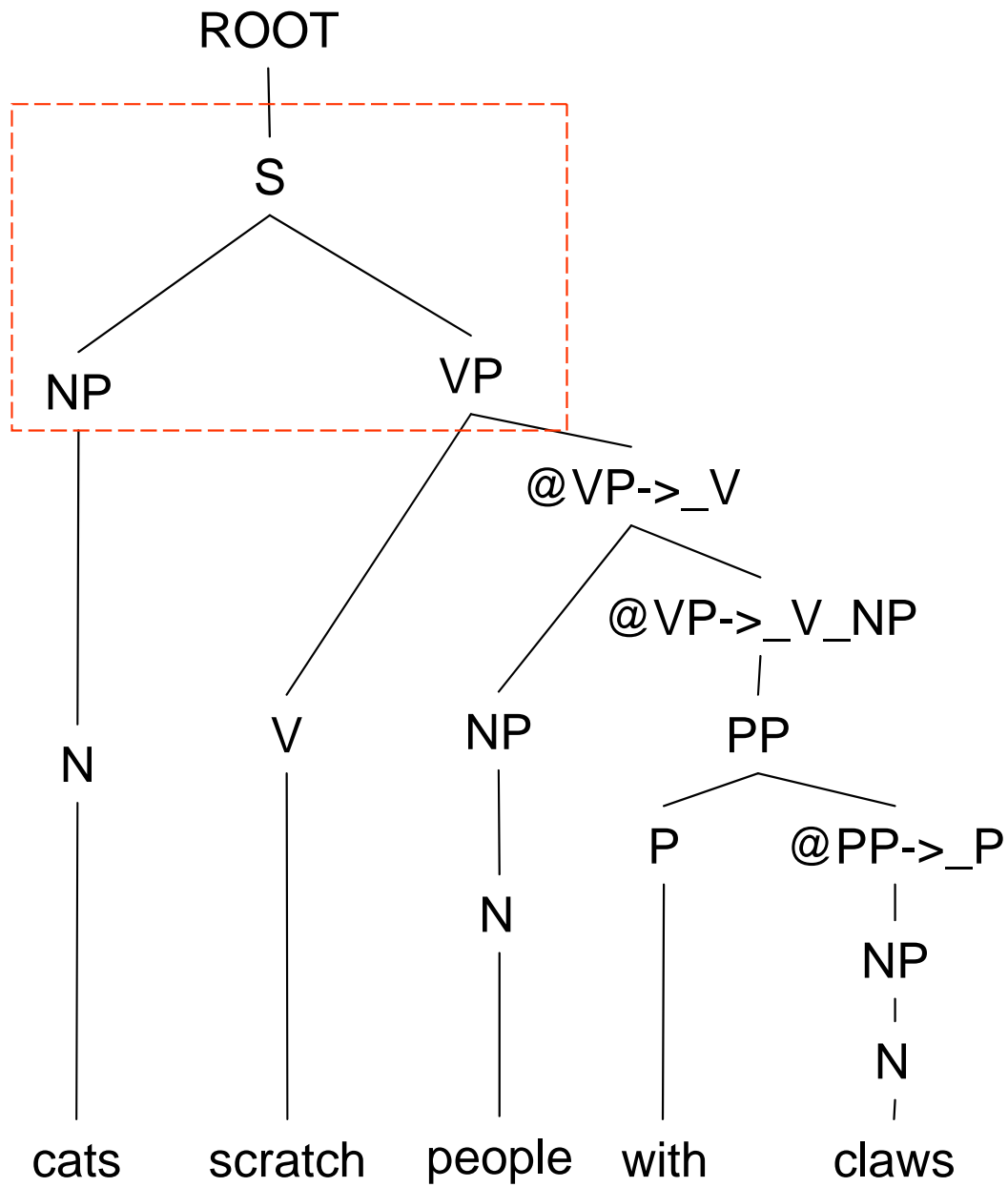
Seems redundant? (the rule was already binary)
Reason: easier to see how to make
finite-order horizontal markovizations
(explained later)

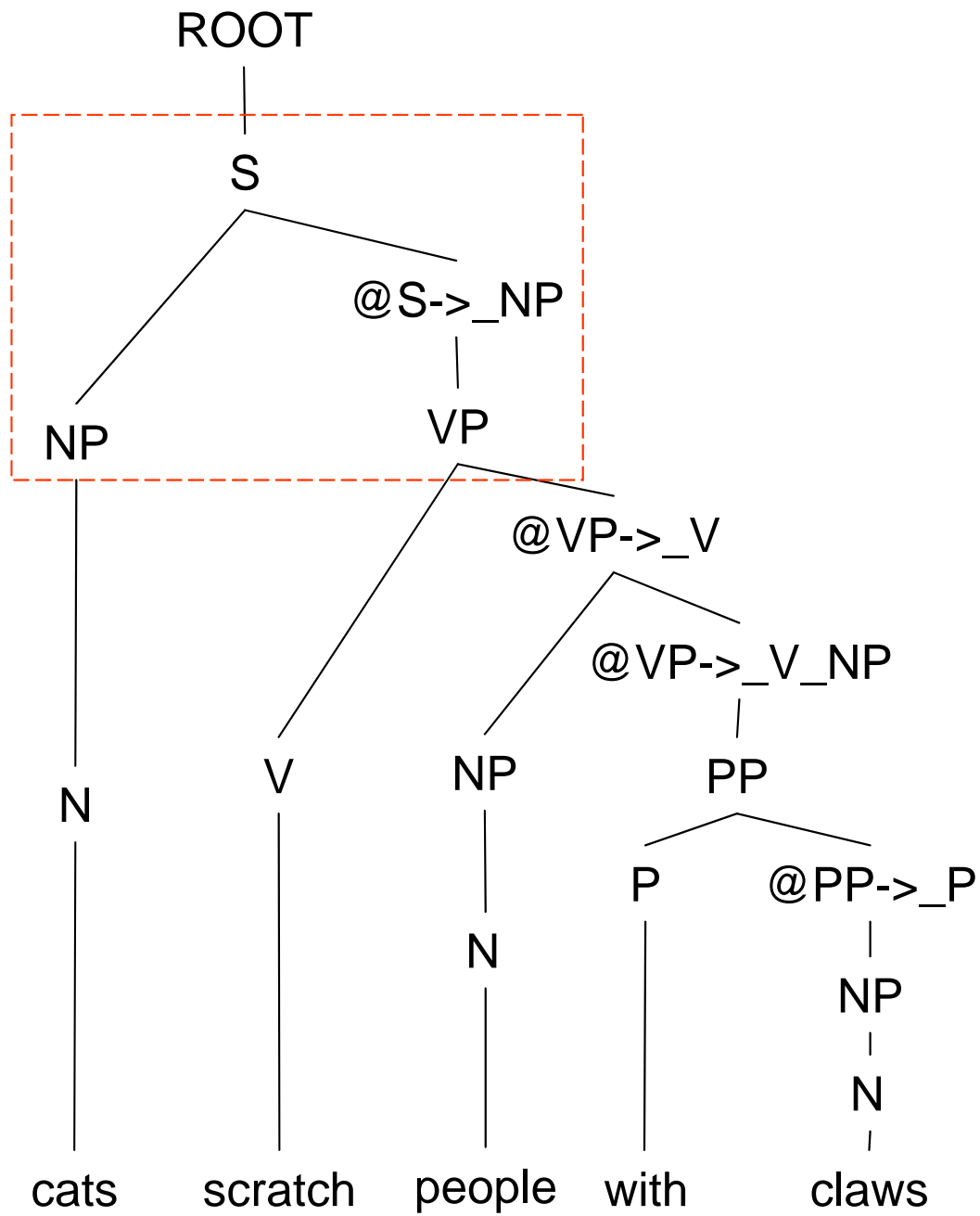


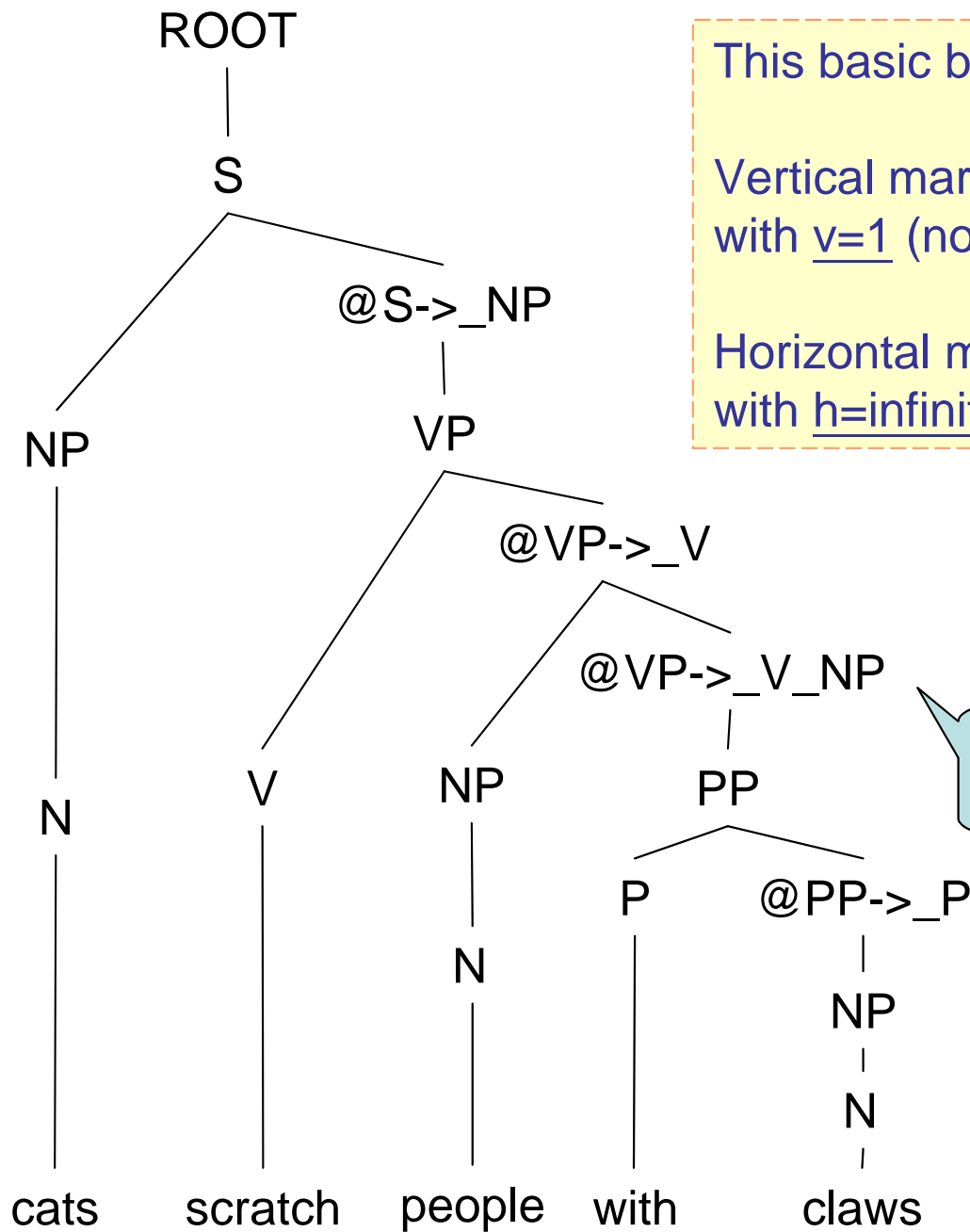


ternary rule









This basic binarization is equivalent to:
 Vertical markovization
 with v=1 (no parent annotation)
 Horizontal markovization
 with h=infinity (not forgetting any siblings)

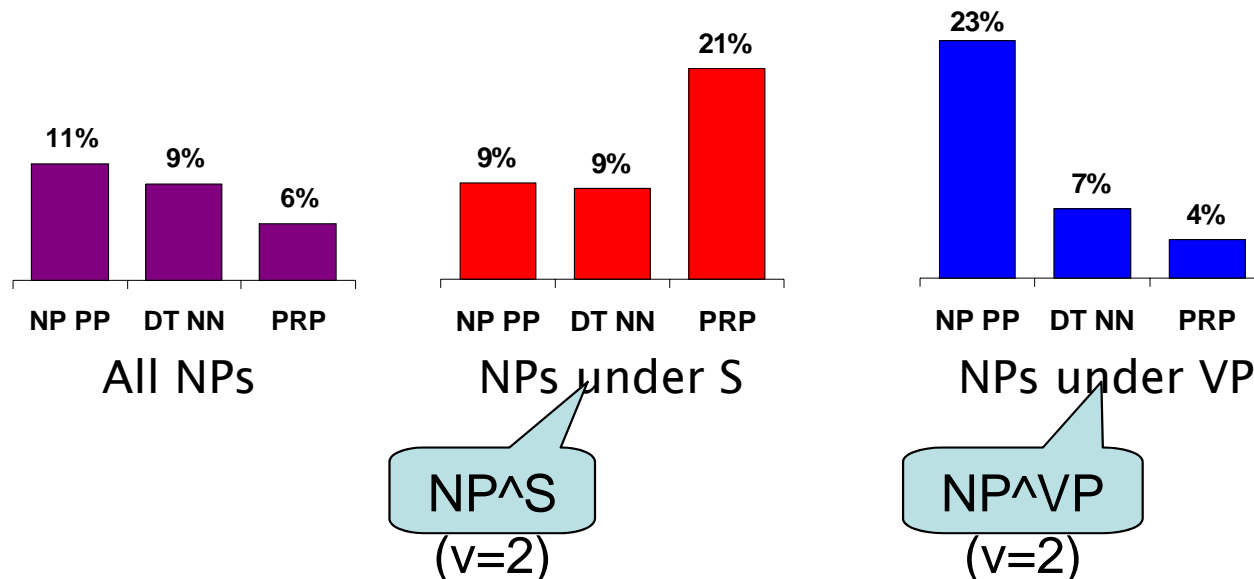
VP → V NP PP
 Remembers 2 siblings

If there's a rule
 VP → V NP PP PP
 ,
 @VP->_V_NP_PP
 will exist.

Two deficiencies of basic binarization (1/2)

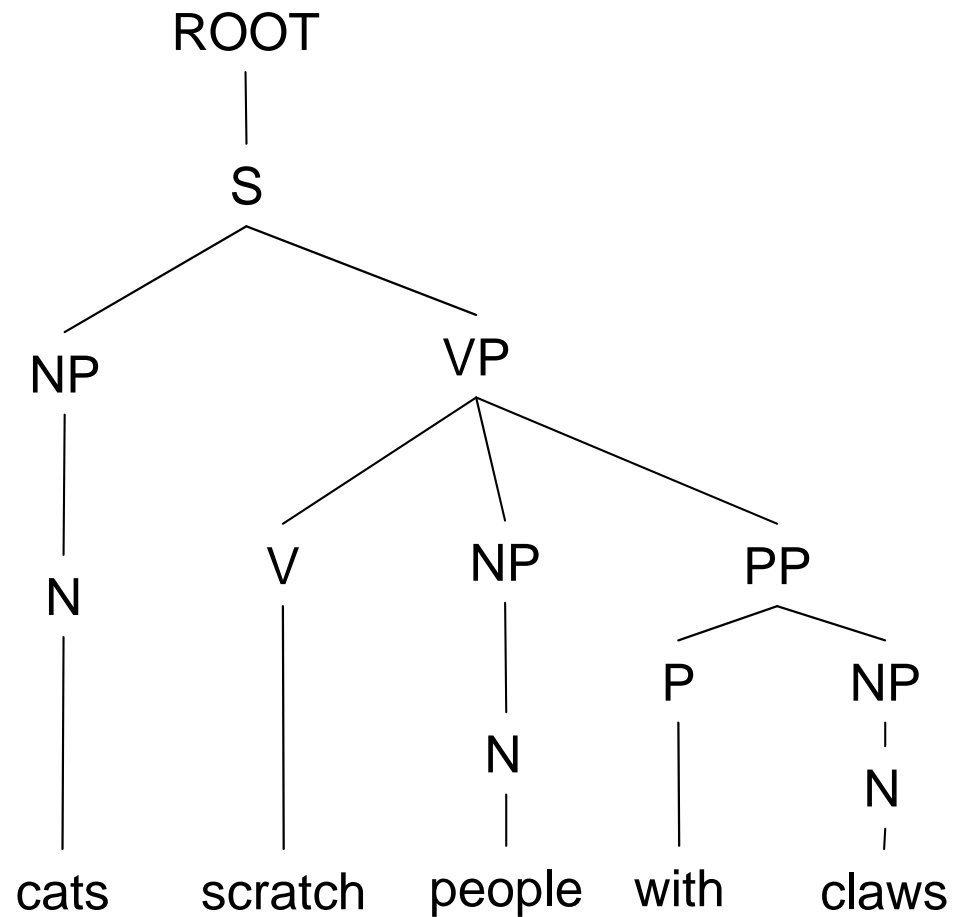
1. PCFG independence assumption

- Often too strong.
- Example: the expansion of an NP is highly dependent on the parent of the NP (i.e., subjects vs. objects).

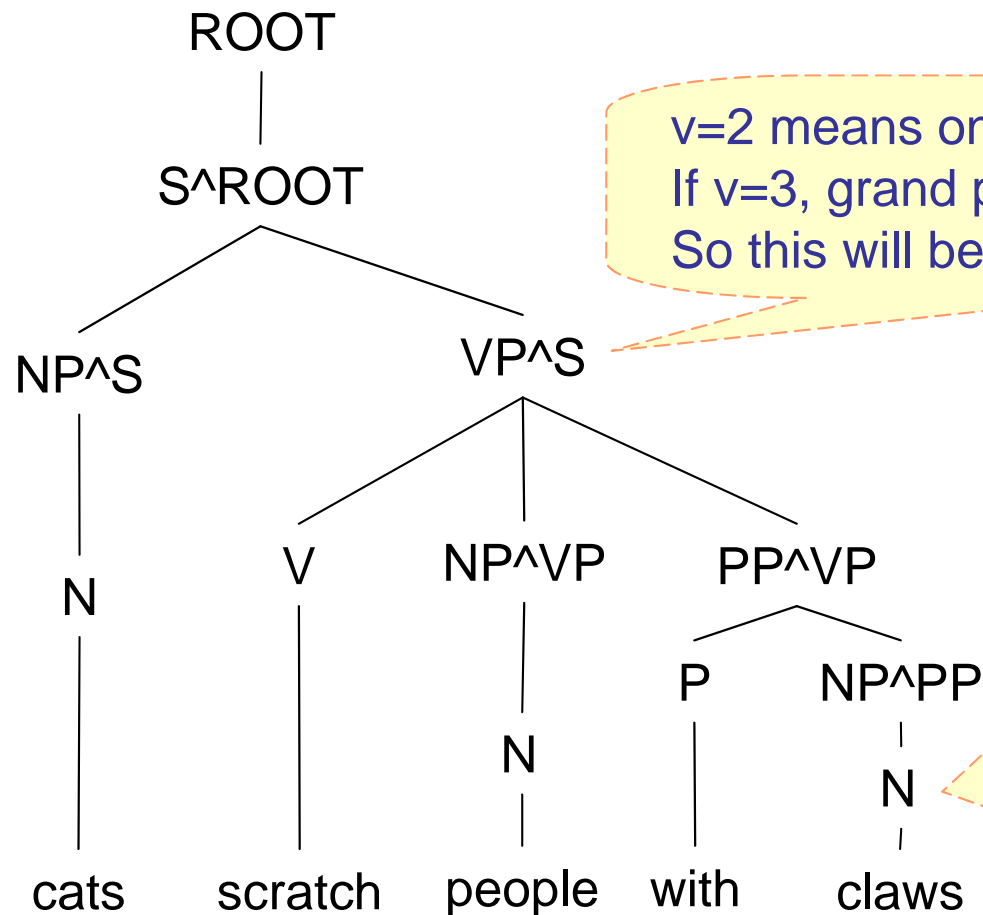


- Can be solved by “vertical markovization” (parent annotation)

Vertical Markovization (v=2): Before...



Vertical Markovization (v=2): After...

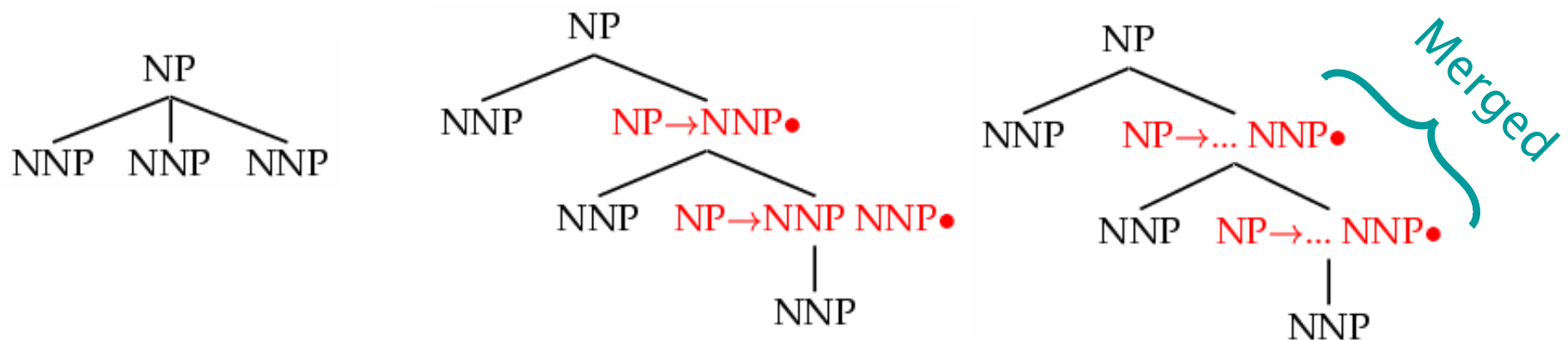


v=2 means only annotate with the parent node.
If v=3, grand parent node is also included.
So this will be **VP^S^ROOT**

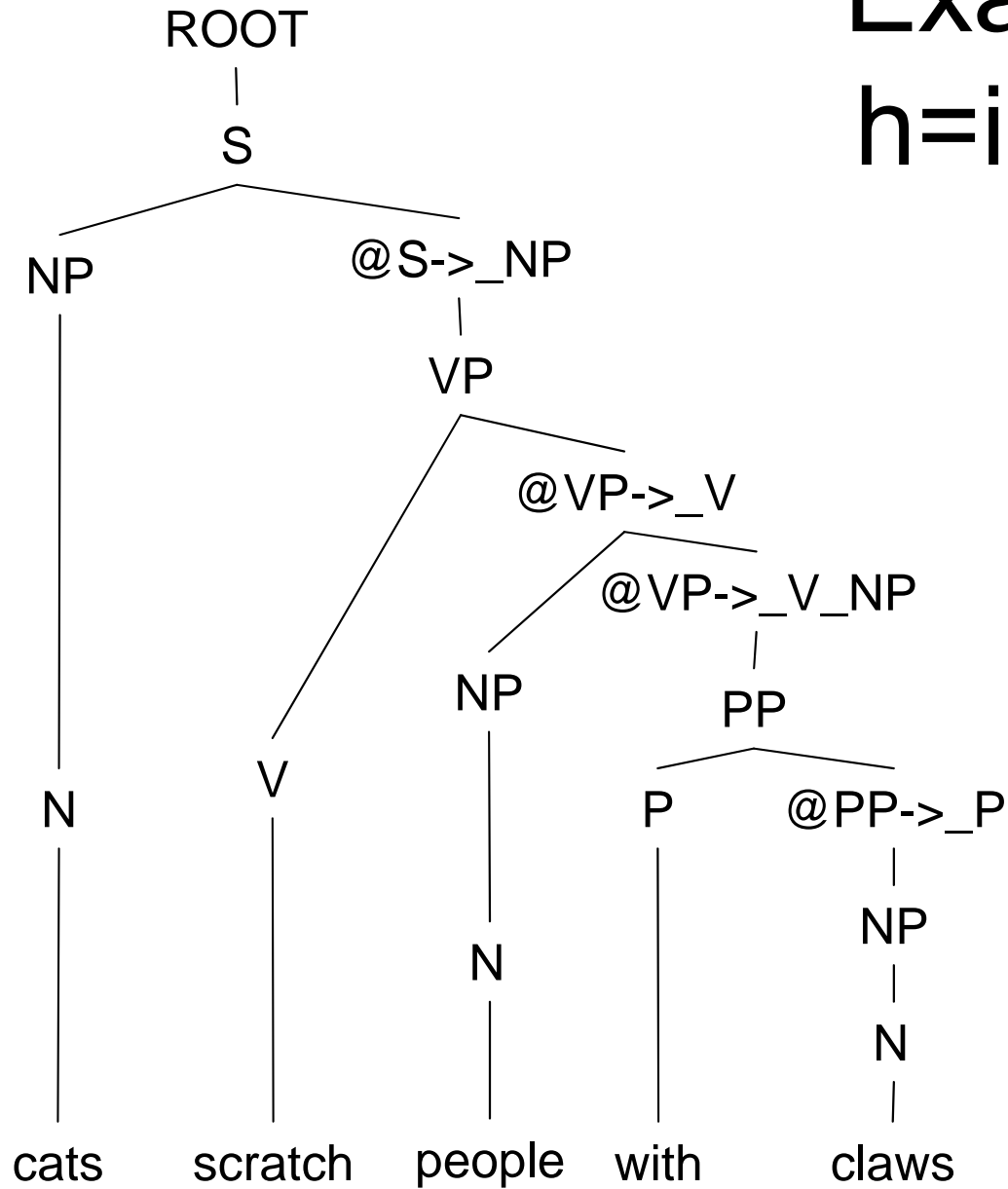
Suggestion:
Don't parent annotate
POS tags.
You'll have to be
more careful about
smoothing the lexicon
if you really want to.

Two deficiencies of basic binarization (2/2)

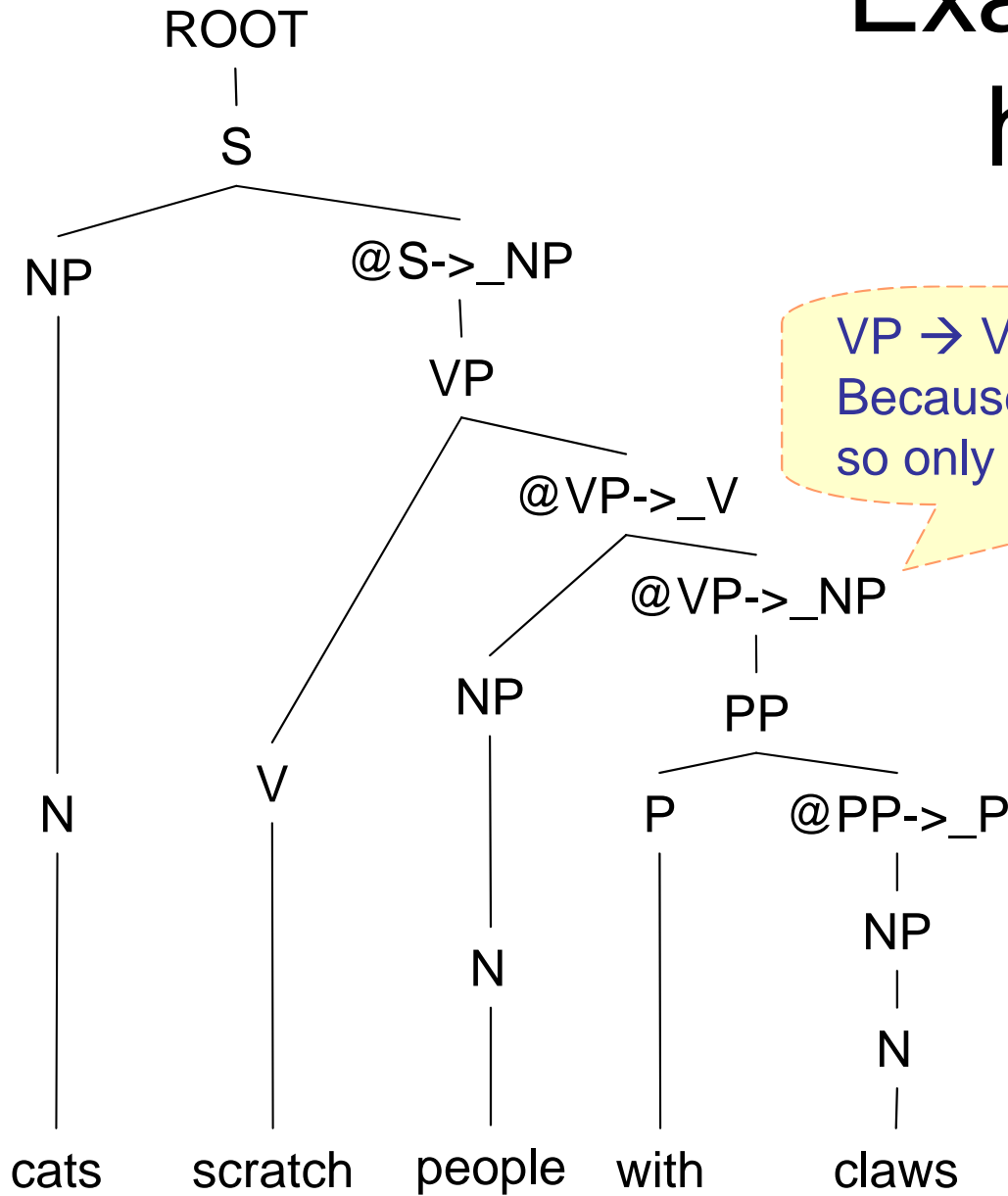
2. Many rules have been seen only once
- Sparseness
 - We can make the horizontal markovization more forgetful.



Example:
h=infinity



Example: h=1



VP → V NP PP
Because h=1,
so only remember "NP" but forget "V"

Some tips on Markovization

1. Vertical & horizontal

- In the Stanford Parser, the order is: first do vertical markovization, and then horizontal markovization.

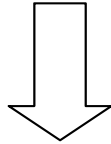
2. “unAnnotateTree” method

- Although the comment said the unannotation cuts at the leftmost -, ^, or : character, but it actually cuts at ‘-’ or ‘=’.
- One solution: instead of “NP^S”, use “NP-^S” or “NP=S”

3. Don’t parent annotate POS tags.

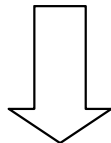
- Can be useful as well. But you need to do some fancier smoothing to get it to work well, and leaving it out will keep your grammar more compact.

N-ary Trees in Treebank



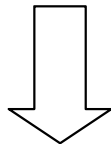
In the starter code :
binarizeTree

TreeAnnotations.annotateTree

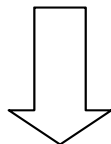


Binary Trees

DONE:
Better vertical and horizontal
markovization



Lexicon and Grammar



Parsing

TODO1:
CKY parsing

CKY algorithm

```

function CKY(words, grammar) returns most probable parse/probability
  score = new double[#(words)+1][#(words)+1][#(nonterms)]
  back = new Pair[#(words)+1][#(words)+1][#nonterms]]
  for i=0; i<#(words); i++
    for A in nonterms
      if A -> words[i] in grammar
        score[i][i+1][A] = P(A -> words[i])
  //handle unaries
  boolean added = true
  while added
    added = false
    for A, B in nonterms
      if score[i][i+1][B] > 0 && A->B in grammar
        prob = P(A->B)*score[i][i+1][B]
        if(prob > score[i][i+1][A])
          score[i][i+1][A] = prob
          back[i][i+1] [A] = B
          added = true

```

```

for span = 2 to #(words)
  for begin = 0 to #(words)- span
    end = begin + span
    for split = begin+1 to end-1
      for A,B,C in nonterms
        prob=score[begin][split][B]*score[split][end][C]*P(A->BC)
        if(prob > score[begin][end][A])
          score[begin][end][A] = prob
          back[begin][end][A] = new Triple(split,B,C)
//handle unaries
boolean added = true
while added
  added = false
  for A, B in nonterms
    prob = P(A->B)*score[begin][end][B];
    if(prob > score[begin][end] [A])
      score[begin][end] [A] = prob
      back[begin][end] [A] = B
      added = true
return buildTree(score, back)

```


	cats	1	scratch	2	walls	3	with	4	claws	5
0	N→cats P→cats V→cats									
1		N→scratch P→scratch V→scratch								
2			N→walls P→walls V→walls							
3				N→with P→with V→with						
4								N→claws P→claws V→claws		
5										

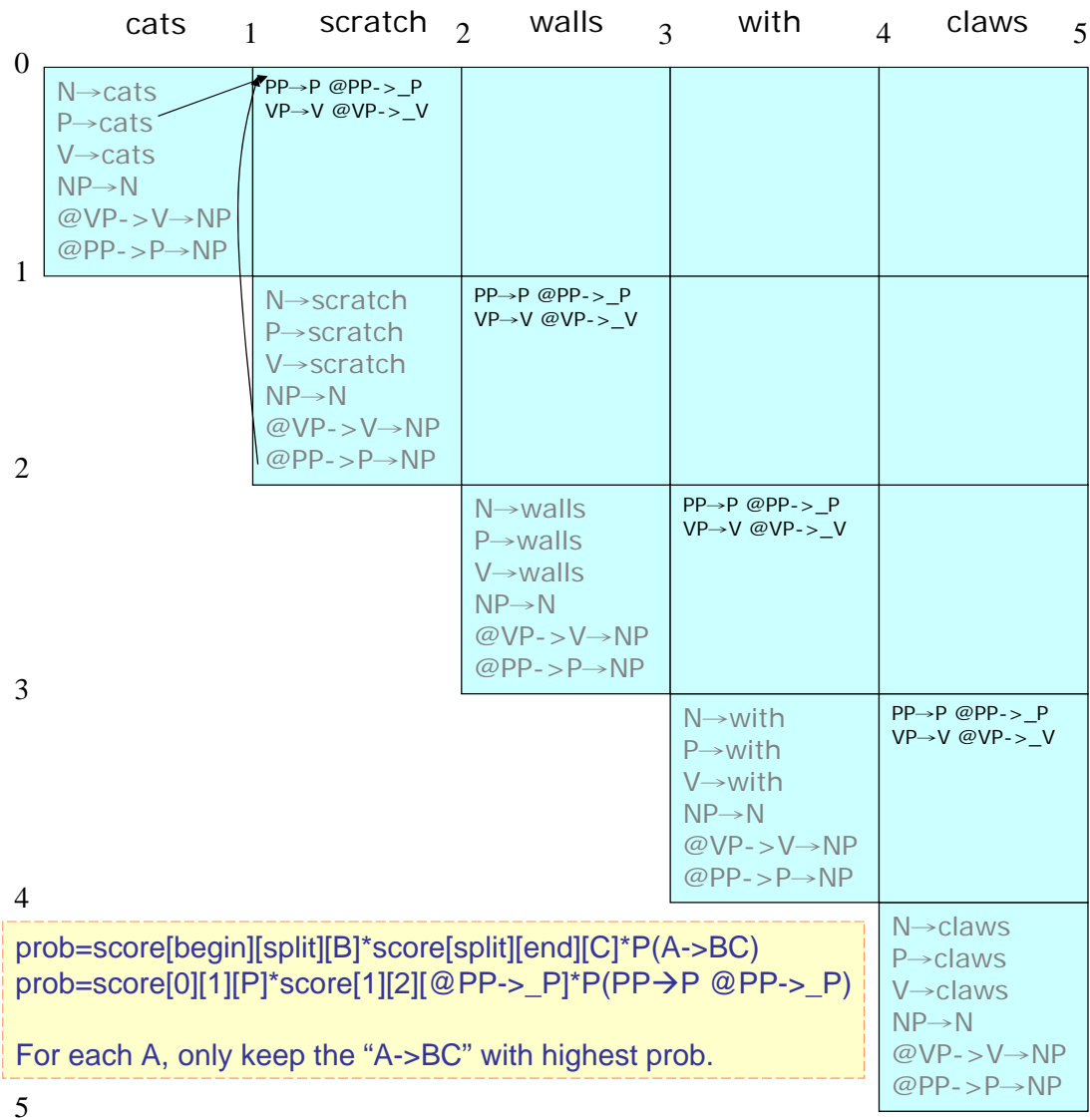
```

for i=0; i<#(words); i++
  for A in nonterms
    if A -> words[i] in grammar
      score[i][i+1][A] = P(A -> words[i]);

```


	cats	1	scratch	2	walls	3	with	4	claws	5
0	N→cats P→cats V→cats NP→N @VP->V→NP @PP->P→NP									
1		N→scratch P→scratch V→scratch NP→N @VP->V→NP @PP->P→NP								
2			N→walls P→walls V→walls NP→N @VP->V→NP @PP->P→NP							
3				N→with P→with V→with NP→N @VP->V→NP @PP->P→NP						
4							N→claws P→claws V→claws NP→N @VP->V→NP @PP->P→NP			
5										

// handle unaries



	cats	1	scratch	2	walls	3	with	4	claws	5
0	N→cats P→cats V→cats NP→N @VP->V→NP @PP->P→NP	PP→P @PP->_P VP→V @VP->_V @S->_NP→VP @NP->_NP→PP @VP->_V_NP→PP								
1		N→scratch P→scratch V→scratch NP→N @VP->V→NP @PP->P→NP	PP→P @PP->_P VP→V @VP->_V @S->_NP→VP @NP->_NP→PP @VP->_V_NP→PP							
2			N→walls P→walls V→walls NP→N @VP->V→NP @PP->P→NP	PP→P @PP->_P VP→V @VP->_V @S->_NP→VP @NP->_NP→PP @VP->_V_NP→PP						
3				N→with P→with V→with NP→N @VP->V→NP @PP->P→NP	PP→P @PP->_P VP→V @VP->_V @S->_NP→VP @NP->_NP→PP @VP->_V_NP→PP					
4								N→claws P→claws V→claws NP→N @VP->V→NP @PP->P→NP		
5										

// handle unaries

.....

	cats	1	scratch	2	walls	3	with	4	claws	5
0	N→cats 0.5259 P→cats 0.0725 V→cats 0.0967 NP→N 0.4675 @VP->V→NP 0.3116 @PP->P→NP 0.4675	PP→P @PP->_P 0.0062 VP→V @VP->_V 0.0055 @S->_NP→VP 0.0055 @NP->_NP→PP 0.0062 @VP->_V_NP→PP 0.0062	@VP->_V→NP @VP->_V_NP 0.0030 NP→NP @NP->_NP 0.0010 S→NP @S->_NP 0.0727 ROOT→S 0.0727 @PP->_P→NP 0.0010	PP→P @PP->_P 5.187E-6 VP→V @VP->_V 2.074E-5 @S->_NP→VP 2.074E-5 @NP->_NP→PP 5.187E-6 @VP->_V_NP→PP 5.187E-6	@VP->_V→NP @VP->_V_NP 1.600E-4 NP→NP @NP->_NP 5.335E-5 S→NP @S->_NP 0.0172 ROOT→S 0.0172 @PP->_P→NP 5.335E-5					
1		N→scratch 0.0967 P→scratch 0.0773 V→scratch 0.9285 NP→N 0.0859 @VP->V→NP 0.0573 @PP->P→NP 0.0859	PP→P @PP->_P 0.0194 VP→V @VP->_V 0.1556 @S->_NP→VP 0.1556 @NP->_NP→PP 0.0194 @VP->_V_NP→PP 0.0194	@VP->_V→NP @VP->_V_NP 2.145E-4 NP→NP @NP->_NP 7.150E-5 S→NP @S->_NP 5.720E-4 ROOT→S 5.720E-4 @PP->_P→NP 7.150E-5	PP→P @PP->_P 0.0010 VP→V @VP->_V 0.0369 @S->_NP→VP 0.0369 @NP->_NP→PP 0.0010 @VP->_V_NP→PP 0.0010					
2			N→walls 0.2829 P→walls 0.0870 V→walls 0.1160 NP→N 0.2514 @VP->V→NP 0.1676 @PP->P→NP 0.2514	PP→P @PP->_P 0.0074 VP→V @VP->_V 0.0066 @S->_NP→VP 0.0066 @NP->_NP→PP 0.0074 @VP->_V_NP→PP 0.0074	@VP->_V→NP @VP->_V_NP 0.0398 NP→NP @NP->_NP 0.0132 S→NP @S->_NP 0.0062 ROOT→S 0.0062 @PP->_P→NP 0.0132					
3				N→with 0.0967 P→with 1.3154 V→with 0.1031 NP→N 0.0859 @VP->V→NP 0.0573 @PP->P→NP 0.0859	PP→P @PP->_P 0.4750 VP→V @VP->_V 0.0248 @S->_NP→VP 0.0248 @NP->_NP→PP 0.4750 @VP->_V_NP→PP 0.4750					
4								N→claws 0.4062 P→claws 0.0773 V→claws 0.1031 NP→N 0.3611 @VP->V→NP 0.2407 @PP->P→NP 0.3611		
5										

Call buildTree(score, back) to get the best parse