

## The Expectation Maximization (EM) Algorithm

... continued!

600.465 - Intro to NLP - J. Eisner

1

## General Idea

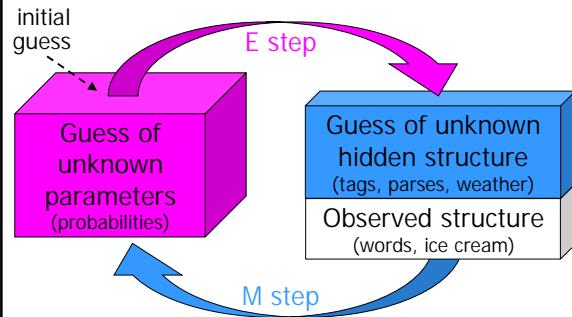
- Start by devising a noisy channel
  - Any model that predicts the corpus observations via some hidden structure (tags, parses, ...)
- Initially **guess** the parameters of the model!
  - Educated guess is best, but random can work
- **Expectation step:** Use current parameters (and observations) to reconstruct hidden structure
- **Maximization step:** Use that hidden structure (and observations) to reestimate parameters

Repeat until convergence!

600.465 - Intro to NLP - J. Eisner

2

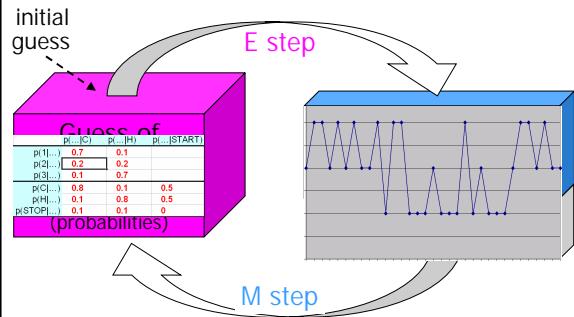
## General Idea



600.465 - Intro to NLP - J. Eisner

3

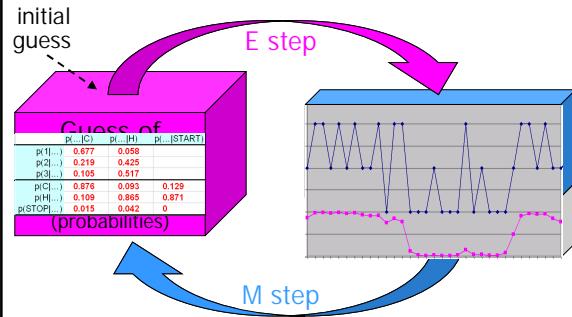
## For Hidden Markov Models



600.465 - Intro to NLP - J. Eisner

4

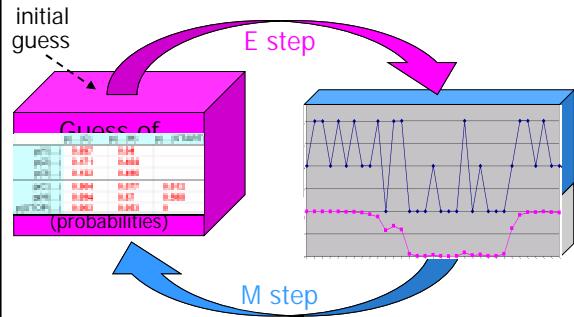
## For Hidden Markov Models



600.465 - Intro to NLP - J. Eisner

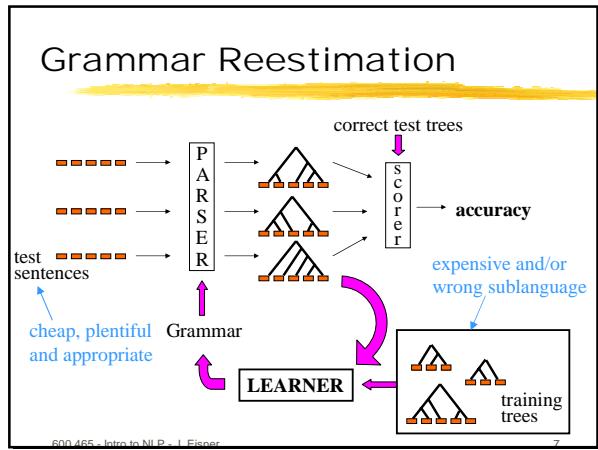
5

## For Hidden Markov Models



600.465 - Intro to NLP - J. Eisner

6



600.465 - Intro to NLP - J. Eisner

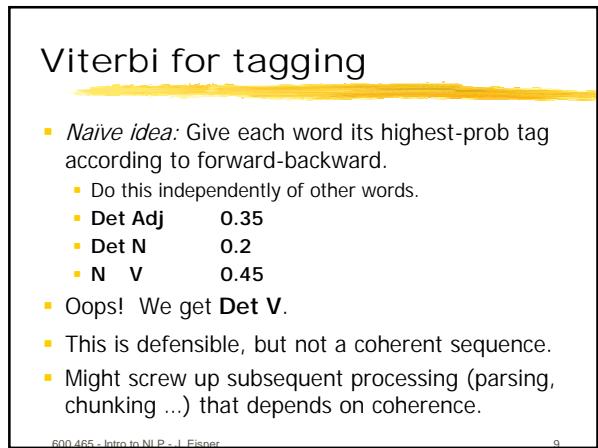
7

## EM by Dynamic Programming: Two Versions

- The Viterbi approximation
  - **Expectation:** pick the *best* parse of each sentence
  - **Maximization:** retrain on this best-parsed corpus
  - **Advantage:** Speed!
- Real EM *why slower?*
  - **Expectation:** find *all* parses of each sentence
  - **Maximization:** retrain on *all* parses in proportion to their probability (as if we observed fractional count)
  - **Advantage:**  $p(\text{training corpus})$  guaranteed to increase
  - Exponentially many parses, so don't extract them from chart – need some kind of clever counting

600.465 - Intro to NLP - J. Eisner

8



600.465 - Intro to NLP - J. Eisner

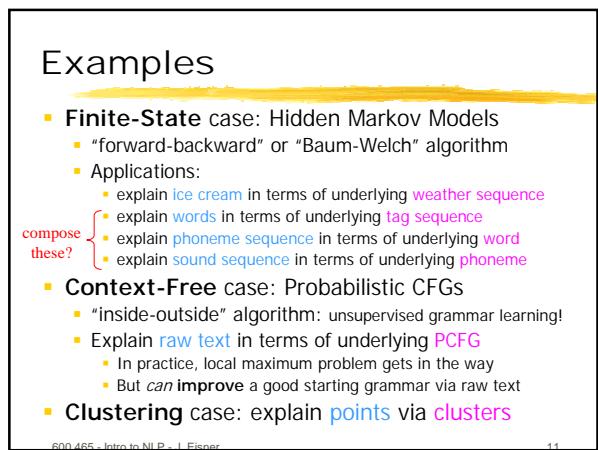
9

## Viterbi for tagging

- **Naive idea:** Give each word its highest-prob tag according to forward-backward.
  - Do this independently of other words.
  - Det Adj      0.35
  - Det N        0.2
  - N    V       0.45
- **Better idea:** Use Viterbi to pick single best tag sequence (best path).
- Can still use full EM to train the parameters, if time permits, but use Viterbi to tag once trained.

600.465 - Intro to NLP - J. Eisner

10



600.465 - Intro to NLP - J. Eisner

11

## Inside-Outside Algorithm

- How would we do the Viterbi approximation?
- How would we do full EM?

600.465 - Intro to NLP - J. Eisner

12

## PCFG

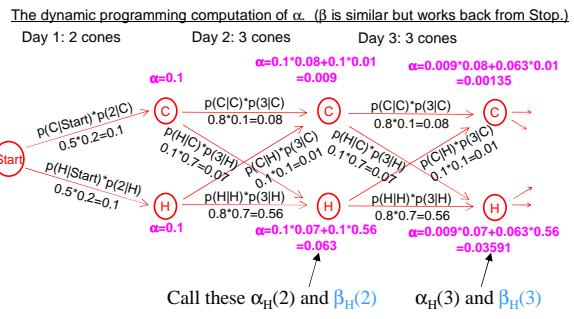
$$p(\text{S} \mid \text{time flies like an arrow}) = p(\text{S} \rightarrow \text{NP VP} \mid \text{S}) * p(\text{NP} \rightarrow \text{time} \mid \text{NP}) * p(\text{VP} \rightarrow \text{V PP} \mid \text{VP}) * p(\text{V} \rightarrow \text{flies} \mid \text{V}) * \dots$$

$\begin{array}{c} \text{S} \\ | \\ \text{NP} \quad \text{VP} \\ | \quad | \\ \text{time} \quad \text{flies} \\ | \quad | \\ \text{V} \quad \text{PP} \\ | \quad | \\ \text{like} \quad \text{NP} \\ | \quad | \\ \text{Det} \quad \text{N} \\ \text{an arrow} \end{array}$

600.465 - Intro to NL P - J. Eisner

13

## Analogies to $\alpha$ , $\beta$ in PCFG?



600.465 - Intro to NL P - J. Eisner

14

## "Inside Probabilities"

$$p(\text{S} \mid \text{time flies like an arrow}) = p(\text{S} \rightarrow \text{NP VP} \mid \text{S}) * p(\text{NP} \rightarrow \text{time} \mid \text{NP}) * p(\text{VP} \rightarrow \text{V PP} \mid \text{VP}) * p(\text{V} \rightarrow \text{flies} \mid \text{V}) * \dots$$

$\begin{array}{c} \text{S} \\ | \\ \text{NP} \quad \text{VP} \\ | \quad | \\ \text{time} \quad \text{flies} \\ | \quad | \\ \text{V} \quad \text{PP} \\ | \quad | \\ \text{like} \quad \text{NP} \\ | \quad | \\ \text{Det} \quad \text{N} \\ \text{an arrow} \end{array}$

- Sum over all VP parses of "flies like an arrow":  
 $\beta_{\text{VP}}(1,5) = p(\text{flies like an arrow} \mid \text{VP})$
- Sum over all S parses of "time flies like an arrow":  
 $\beta_S(0,5) = p(\text{time flies like an arrow} \mid \text{S})$

600.465 - Intro to NL P - J. Eisner

15

## Compute $\beta$ Bottom-Up by CKY

$$p(\text{S} \mid \text{time flies like an arrow}) = p(\text{S} \rightarrow \text{NP VP} \mid \text{S}) * p(\text{NP} \rightarrow \text{time} \mid \text{NP}) * p(\text{VP} \rightarrow \text{V PP} \mid \text{VP}) * p(\text{V} \rightarrow \text{flies} \mid \text{V}) * \dots$$

$\begin{array}{c} \text{S} \\ | \\ \text{NP} \quad \text{VP} \\ | \quad | \\ \text{time} \quad \text{flies} \\ | \quad | \\ \text{V} \quad \text{PP} \\ | \quad | \\ \text{like} \quad \text{NP} \\ | \quad | \\ \text{Det} \quad \text{N} \\ \text{an arrow} \end{array}$

$$\beta_{\text{VP}}(1,5) = p(\text{flies like an arrow} \mid \text{VP}) = \dots$$

$$\beta_S(0,5) = p(\text{time flies like an arrow} \mid \text{S}) = \beta_{\text{NP}}(0,1) * \beta_{\text{VP}}(1,5) * p(\text{S} \rightarrow \text{NP VP} \mid \text{S}) + \dots$$

600.465 - Intro to NL P - J. Eisner

16

## Compute $\beta$ Bottom-Up by CKY

time	1	flies	2	like	3	an	4	arrow	5
0	NP 3 Vst 3	NP 10 S 8				NP 24 NP 24 S 22 S 27 S 27			
1		NP 4 VP 4				NP 18 S 21 VP 18			
2			P 2 V 5			PP 12 VP 16			
3				Det 1	NP 10				
4					N 8		0 PP → P NP		

1 S → NP VP  
6 S → Vst NP  
2 S → S PP  
1 VP → V NP  
2 VP → VP PP  
1 NP → Det N  
2 NP → NP PP  
3 NP → NP NP  
0 PP → P NP

## Compute $\beta$ Bottom-Up by CKY

time	1	flies	2	like	3	an	4	arrow	5
0	NP $2^{-3}$ Vst $2^{-3}$	NP $2^{-10}$ S $2^{-8}$				NP $2^{-24}$ NP $2^{-24}$ S $2^{-22}$ S $2^{-27}$ S $2^{-22}$			
1			NP $2^{-4}$ VP $2^{-4}$			NP $2^{-18}$ S $2^{-21}$ VP $2^{-18}$			
2				P $2^{-2}$ V $2^{-5}$		PP $2^{-12}$ VP $2^{-16}$			
3					Det 2 $^{-1}$	NP 2 $^{-10}$			
4						N 2 $^{-8}$			

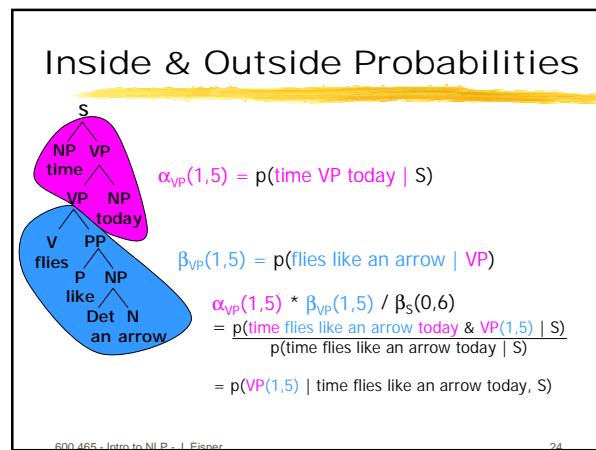
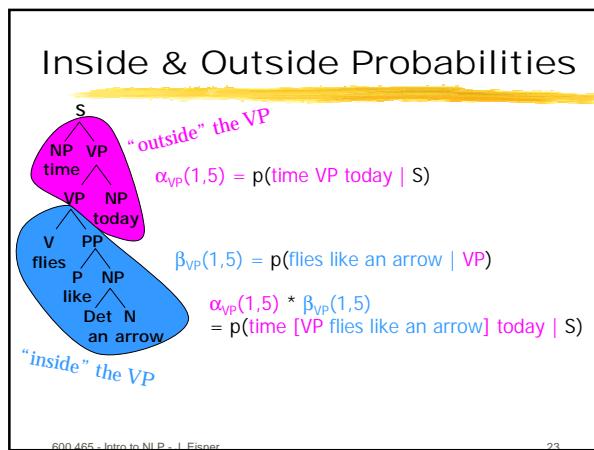
$2^{-1}$  S → NP VP  
 $2^{-6}$  S → Vst NP  
 $2^{-2}$  S → S PP  
 $2^{-1}$  VP → V NP  
 $2^{-2}$  VP → VP PP  
 $2^{-1}$  NP → Det N  
 $2^{-2}$  NP → NP PP  
 $2^{-3}$  NP → NP NP  
 $2^{-0}$  PP → P NP

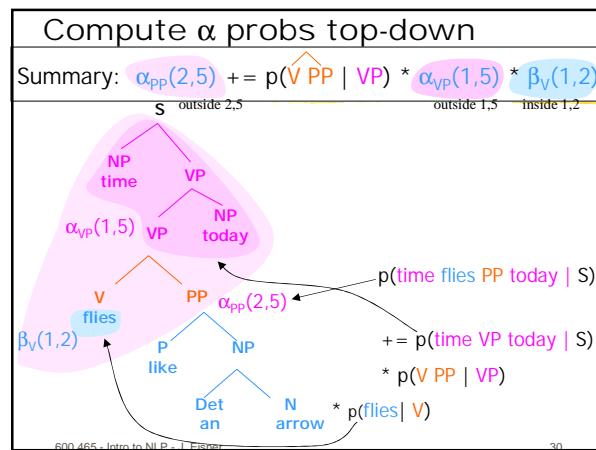
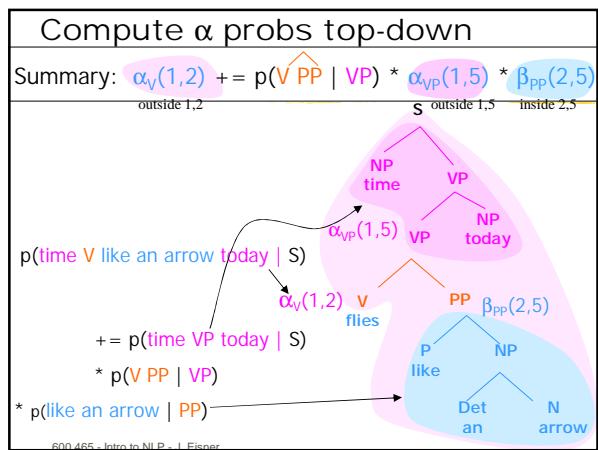
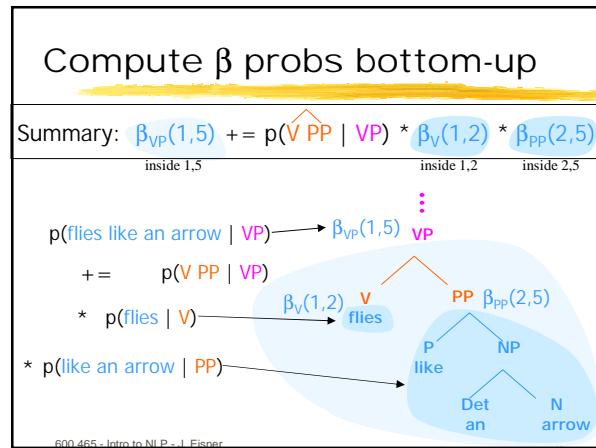
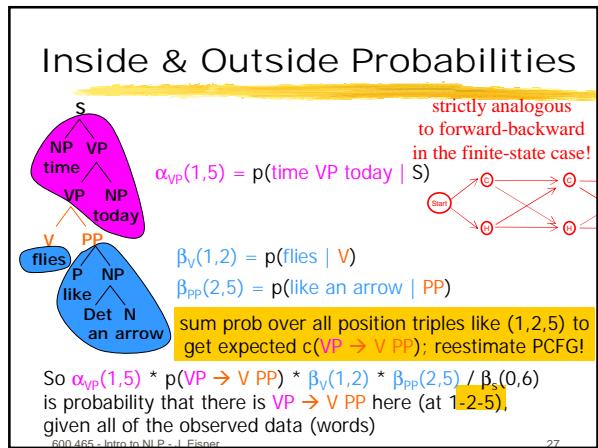
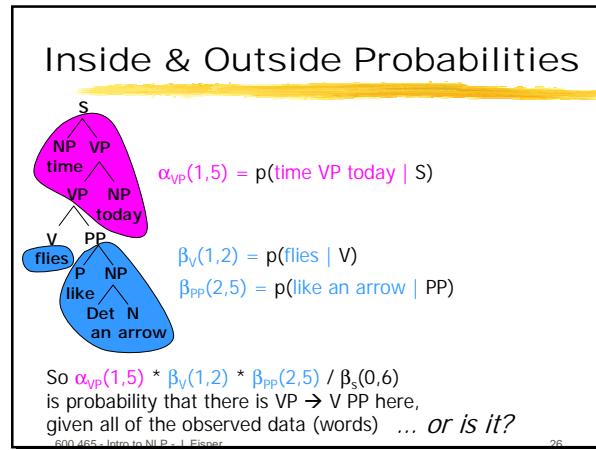
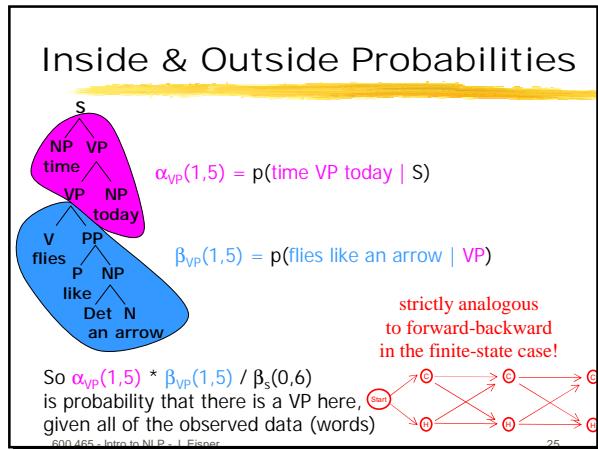
Compute $\beta$ Bottom-Up by CKY					
time	1	flies	2	like	3
0	NP $2^{-3}$ Vst $2^{-3}$	NP $2^{-10}$ S $2^{-8}$ S $2^{-13}$			NP $2^{-24}$ NP $2^{-24}$ S $2^{-22}$ S $2^{-27}$ S $2^{-27}$ S $2^{-22}$ S $2^{-27}$ S $2^{-27}$
1					$2^{-1} S \rightarrow NP VP$ $2^{-6} S \rightarrow Vst NP$ $2^{-2} S \rightarrow S PP$
2		NP $2^{-4}$ VP $2^{-4}$			NP $2^{-18}$ S $2^{-21}$ VP $2^{-18}$
3		P $2^{-2}$ V $2^{-5}$			VP $2^{-12}$ VP $2^{-16}$
4			Det $2^{-1}$	NP $2^{-10}$	NP $2^{-24}$ NP $2^{-24}$ S $2^{-22}$ S $2^{-27}$ S $2^{-27}$
				N $2^{-8}$	$2^{-1} VP \rightarrow V NP$ $2^{-2} VP \rightarrow VP PP$ $2^{-1} NP \rightarrow Det N$ $2^{-2} NP \rightarrow NP PP$ $2^{-3} NP \rightarrow NP NP$ $2^{-0} PP \rightarrow P NP$

The Efficient Version: Add as we go					
time	1	flies	2	like	3
0	NP $2^{-3}$ Vst $2^{-3}$	NP $2^{-10}$ S $2^{-8}$ S $2^{-13}$			NP $2^{-24}$ NP $2^{-24}$ S $2^{-22}$ S $2^{-27}$ S $2^{-27}$
1					$2^{-1} S \rightarrow NP VP$ $2^{-6} S \rightarrow Vst NP$ $2^{-2} S \rightarrow S PP$
2		NP $2^{-4}$ VP $2^{-4}$			NP $2^{-18}$ S $2^{-21}$ VP $2^{-18}$
3		P $2^{-2}$ V $2^{-5}$			PP $2^{-12}$ VP $2^{-16}$
4			Det $2^{-1}$	NP $2^{-10}$	NP $2^{-24}$ NP $2^{-24}$ S $2^{-22}$ S $2^{-27}$ S $2^{-27}$
				N $2^{-8}$	$2^{-1} VP \rightarrow V NP$ $2^{-2} VP \rightarrow VP PP$ $2^{-1} NP \rightarrow Det N$ $2^{-2} NP \rightarrow NP PP$ $2^{-3} NP \rightarrow NP NP$ $2^{-0} PP \rightarrow P NP$

The Efficient Version: Add as we go					
time	1	flies	2	like	3
0	NP $2^{-3}$ Vst $2^{-3}$	NP $2^{-10}$ S $2^{-8}$ S $+2^{-13}$			NP $2^{-24}$ + $2^{-24}$ S $2^{-22}$ + $2^{-27}$ + $2^{-27}$ S $2^{-27}$
1		$\beta_S(0,2)$ $\beta_S(0,2) * \beta_{PP}(2,5) * p(S \rightarrow S PP   S)$			$2^{-1} S \rightarrow NP VP$ $2^{-6} S \rightarrow Vst NP$ $2^{-2} S \rightarrow S PP$
2		NP $2^{-4}$ VP $2^{-4}$			NP $2^{-18}$ S $2^{-21}$ VP $2^{-18}$
3		P $2^{-2}$ V $2^{-5}$			PP $2^{-12}$ VP $2^{-16}$
4			Det $2^{-1}$	NP $2^{-10}$	NP $2^{-24}$ NP $2^{-24}$ S $2^{-22}$ S $2^{-27}$ S $2^{-27}$
				N $2^{-8}$	$2^{-1} VP \rightarrow V NP$ $2^{-2} VP \rightarrow VP PP$ $2^{-1} NP \rightarrow Det N$ $2^{-2} NP \rightarrow NP PP$ $2^{-3} NP \rightarrow NP NP$ $2^{-0} PP \rightarrow P NP$

Compute $\beta$ probs bottom-up (CKY)					
need some initialization up here for the width-1 case					
for width := 2 to n (* build smallest first *)					
for i := 0 to n-width (* start *)					
let k := i + width (* end *)					
for j := i+1 to k-1 (* middle *)					
for all grammar rules X → Y Z					
$\beta_X(i,k) = p(X \rightarrow Y Z   X) * \beta_Y(i,j) * \beta_Z(j,k)$					
what if you changed + to max?					
what if you replaced all rule probabilities by 1?					

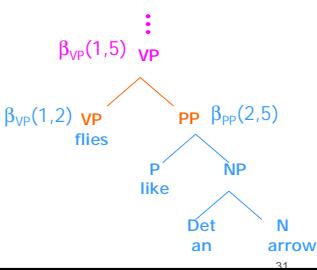




## Compute $\beta$ probs bottom-up

When you build  $VP(1,5)$ ,  
from  $VP(1,2)$  and  $VP(2,5)$   
during CKY,  
increment  $\beta_{VP}(1,5)$  by  
 $p(VP \rightarrow VP PP) * \beta_{VP}(1,2) * \beta_{PP}(2,5)$

Why?  $\beta_{VP}(1,5)$  is total probability of all derivations  
 $p(files \text{ like an arrow} | VP)$  and we just found another.  
(See earlier slide of CKY chart.)

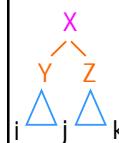


600.465 - Intro to NLP - J. Eisner

31

## Compute $\beta$ probs bottom-up (CKY)

```
for width := 2 to n          (* build smallest first *)
  for i := 0 to n-width      (* start *)
    let k := i + width       (* end *)
    for j := i+1 to k-1       (* middle *)
      for all grammar rules X → Y Z
         $\beta_X(i,k) += p(X \rightarrow Y Z) * \beta_Y(i,j) * \beta_Z(j,k)$ 
```

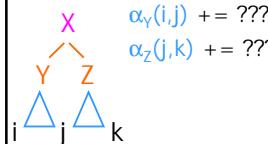


600.465 - Intro to NLP - J. Eisner

32

## Compute $\alpha$ probs top-down (reverse CKY)

```
for width := 2 to n          "unbuild" biggest first
  for i := 0 to n-width      (* build smallest first *)
    let k := i + width       (* start *)
    for j := i+1 to k-1       (* end *)
      for all grammar rules X → Y Z
         $\alpha_Y(i,j) += ???$ 
         $\alpha_Z(j,k) += ???$ 
```



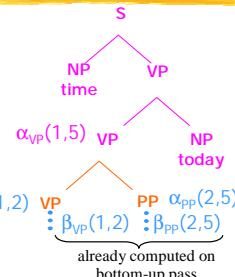
600.465 - Intro to NLP - J. Eisner

33

## Compute $\alpha$ probs top-down

After computing  $\beta$  during CKY,  
revisit constituents in reverse order (i.e.,  
bigger constituents first).

When you "unbuild"  $VP(1,5)$   
from  $VP(1,2)$  and  $VP(2,5)$ ,  
increment  $\alpha_{VP}(1,2)$  by  
 $\alpha_{VP}(1,5) * p(VP \rightarrow VP PP) * \beta_{PP}(2,5)$   
and increment  $\alpha_{PP}(2,5)$  by  
 $\alpha_{VP}(1,5) * p(VP \rightarrow VP PP) * \beta_{VP}(1,2)$



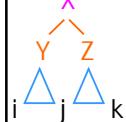
$\alpha_{VP}(1,2)$  is total prob of all ways to gen VP(1,2) and all outside words.

600.465 - Intro to NLP - J. Eisner

34

## Compute $\alpha$ probs top-down (reverse CKY)

```
for width := 2 to n          "unbuild" biggest first
  for i := 0 to n-width      (* build smallest first *)
    let k := i + width       (* start *)
    for j := i+1 to k-1       (* end *)
      for all grammar rules X → Y Z
         $\alpha_Y(i,j) += \alpha_X(i,k) * p(X \rightarrow Y Z) * \beta_Z(j,k)$ 
         $\alpha_Z(j,k) += \alpha_X(i,k) * p(X \rightarrow Y Z) * \beta_Y(i,j)$ 
```



600.465 - Intro to NLP - J. Eisner

35