

# Natural Language Processing

## CSE 325/425



Sihong Xie

### Lecture 17:

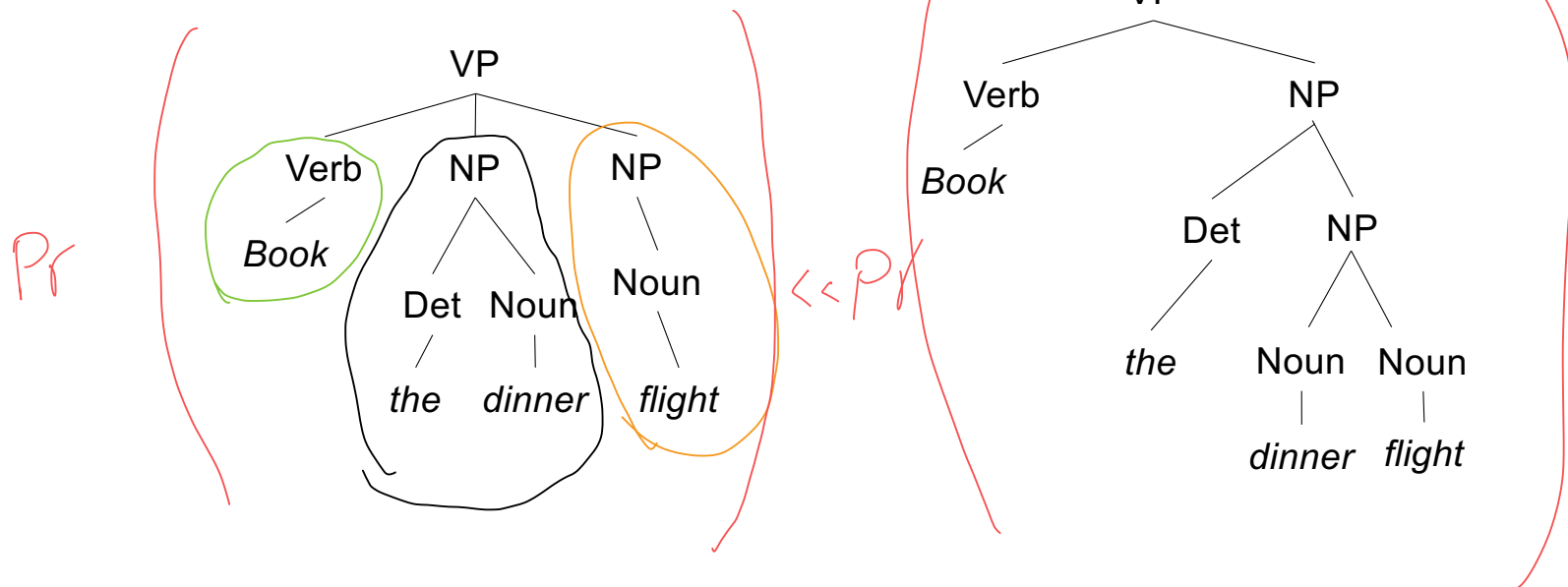
- Probabilistic CFG
- Inside probability and algorithm

# Motivations

CYK algorithm only finds all parsing trees of a given sentence, and can't

- find probability of a tree;
- tell which tree is more likely than another trees;
- identify the optimal tree.

Example: parse "*book the dinner flight*"



# Motivations

HMM can only model linear relationship and has difficulty in long-range dependencies.

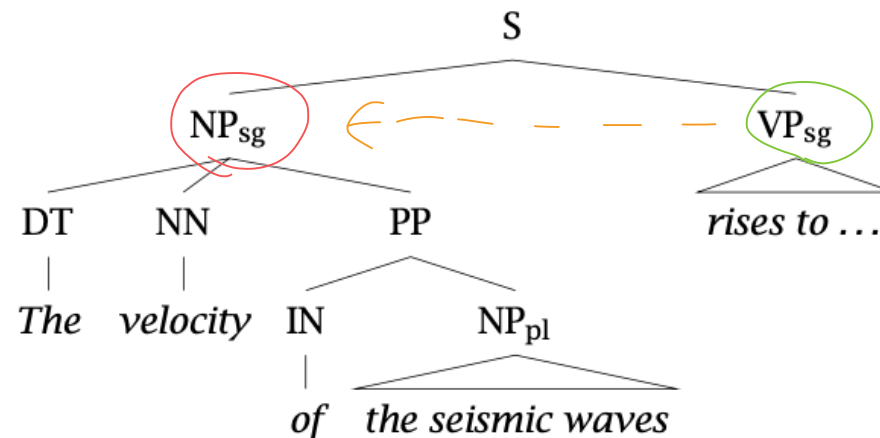
massive noun

plural

$Pr(q_t | q_{t-1})$

Example: parse “The velocity of the seismic waves rises to ...”

- The singular form of the verb “rise” is attached to “velocity”, not its predecessor.



CFG  
 $\Downarrow$

# PCFG

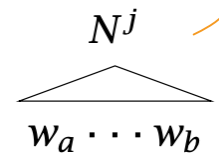
- Assign probability to each production in a CFG
  - the probabilities of the rules with the same left-hand-side sum to one.

Grammar  $G$

$S \rightarrow NP VP$	1.0	$NP \rightarrow NP PP$	0.4 $\leftarrow \Pr(NP PP   VP)$
$PP \rightarrow P NP$	1.0	$NP \rightarrow \textit{astronomers}$	0.1
$VP \rightarrow V NP$	0.7	$NP \rightarrow \textit{ears}$	0.18
$VP \rightarrow VP PP$	0.3	$NP \rightarrow \textit{saw}$	0.04
$P \rightarrow \textit{with}$	1.0	$NP \rightarrow \textit{stars}$	0.18
$V \rightarrow \textit{saw}$	1.0	$NP \rightarrow \textit{telescopes}$	0.1

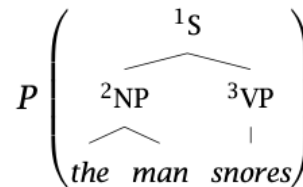
- Want to find the probability of any derivation  $N^j \xRightarrow{*} w_a \dots w_b$

$$\Pr(N^j \xRightarrow{*} w_a \dots w_b | G)$$



- In particular

$$\Pr(S \xRightarrow{*} w_a \dots w_b | G)$$



A superscript labels the non-terminal  
 A subscript indicates the range of words covered

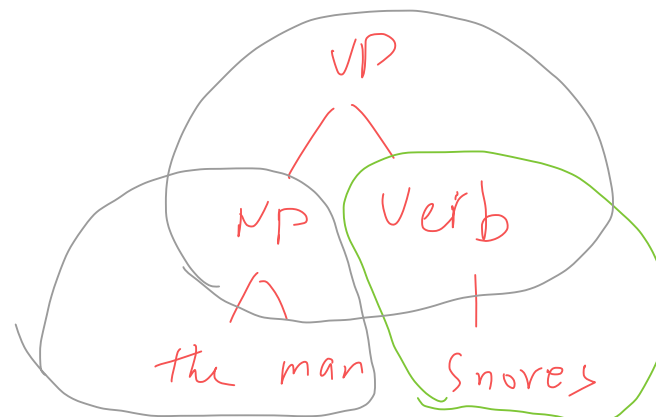
$$= P(^1S_{13} \rightarrow ^2NP_{12} \ ^3VP_{33}, ^2NP_{12} \rightarrow the_1 \ man_2, ^3VP_{33} \rightarrow snores_3)$$



# PCFG

- The probability of a parse tree  $t$  for the sentence  $w_1 \dots w_m$  given the grammar  $G$  is  $\Pr(t, w_1 \dots w_m | G) = \Pr(S \xRightarrow{*} w_1 \dots w_m | G)$
- The probability of the sentence  $w_1 \dots w_m$  is  $\Pr(w_1 \dots w_m | G) = \sum_t \Pr(t, w_1 \dots w_m | G)$   
Total Prob.
- Properties:
  - Place invariance  $\forall k \ P(N_{k(k+c)}^j \rightarrow \zeta)$  is the same
  - Context-free  $P(N_{kl}^j \rightarrow \zeta | \text{anything outside } k \text{ through } l) = P(N_{kl}^j \rightarrow \zeta)$
  - Ancestor-free  $P(N_{kl}^j \rightarrow \zeta | \text{any ancestor nodes outside } N_{kl}^j) = P(N_{kl}^j \rightarrow \zeta)$

$$P_r(x) = \sum_y P_r(x, y)$$



# PCFG

- Probability of a parsing tree that deriving a sentence:

Example continued:

$$Pr(t, w_1 \dots w_3) = P \left( \begin{array}{c} {}^1S \\ \swarrow \quad \searrow \\ {}^2NP \quad {}^3VP \\ \swarrow \quad \searrow \quad | \\ the \quad man \quad snores \end{array} \right)$$

where the tree

$$t = \begin{array}{c} S \\ / \quad \backslash \\ NP \quad verb \\ / \quad | \\ the \quad man \quad snores \end{array}$$

$$\begin{aligned}
 &= P({}^1S_{13} \rightarrow {}^2NP_{12} {}^3VP_{33}, {}^2NP_{12} \rightarrow the_1 man_2, {}^3VP_{33} \rightarrow snores_3) \\
 &= P({}^1S_{13} \rightarrow {}^2NP_{12} {}^3VP_{33}) P({}^2NP_{12} \rightarrow the_1 man_2 | {}^1S_{13} \rightarrow {}^2NP_{12} {}^3VP_{33}) \\
 &\quad P({}^3VP_{33} \rightarrow snores_3 | {}^1S_{13} \rightarrow {}^2NP_{12} {}^3VP_{33}, {}^2NP_{12} \rightarrow the_1 man_2) \\
 &= P({}^1S_{13} \rightarrow {}^2NP_{12} {}^3VP_{33}) P({}^2NP_{12} \rightarrow the_1 man_2) P({}^3VP_{33} \rightarrow snores_3) \\
 &= P(S \rightarrow NP VP) P(NP \rightarrow the man) P(VP \rightarrow snores)
 \end{aligned}$$

def of conditional prob.  
 Ancestor free  
 place invariant.

# Three tasks with PCFG

- Find the probability of a sentence

$$\Pr(w_1 \dots w_m | G) = \sum_t \Pr(t, w_1 \dots w_m | G)$$

- Find the most likely parsing tree

$$t^* = \arg \max_t \Pr(t | w_1 \dots w_m, G)$$

- Learn a PCFG from a training corpus using MLE.

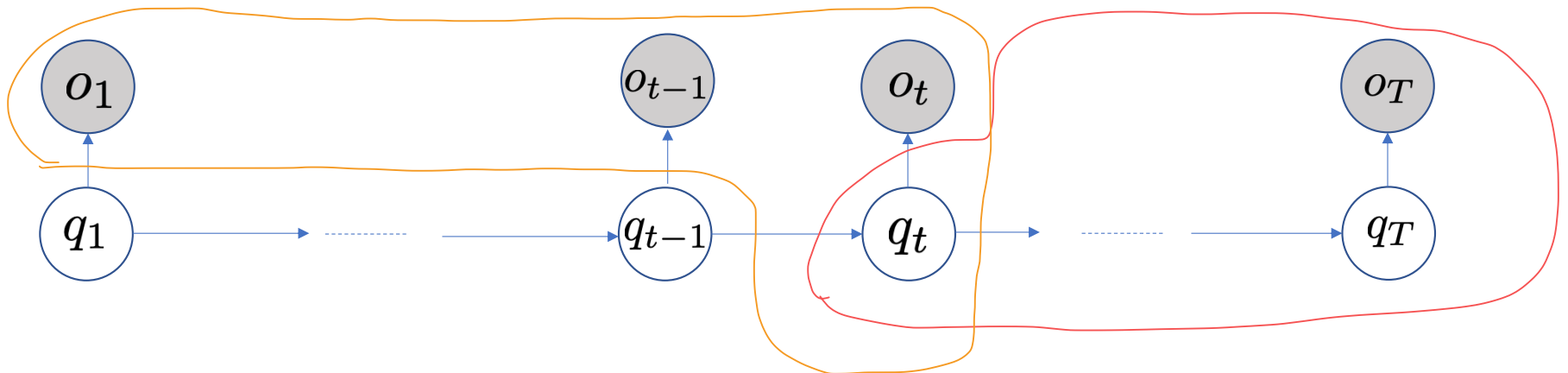
$$G^* = \arg \max_G \Pr(w_1 \dots w_m | G)$$

# From HMM to PCFG

- Assuming Chomsky Normal Form of the PCFG
- From HMM to PCFG

$$\alpha_t(j) = P(o_1, \dots, o_t, q_t = j)$$

$$\beta_t(j) = P(o_{t+1}, \dots, o_T | q_t = j)$$

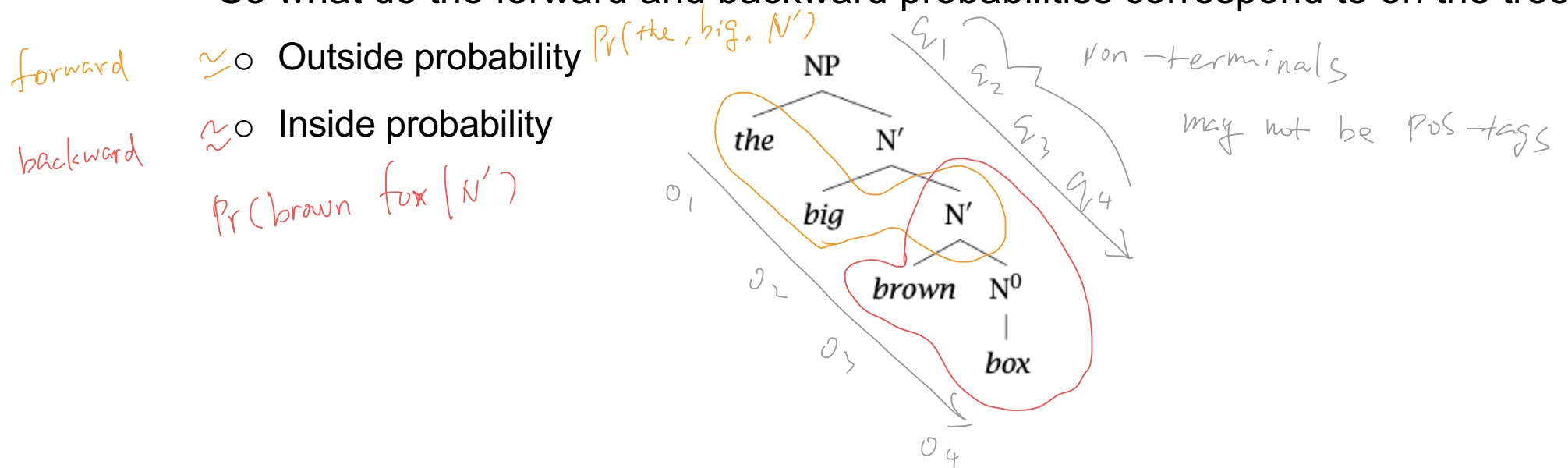


- Can we generalize a HMM sequence to a PCFG tree?

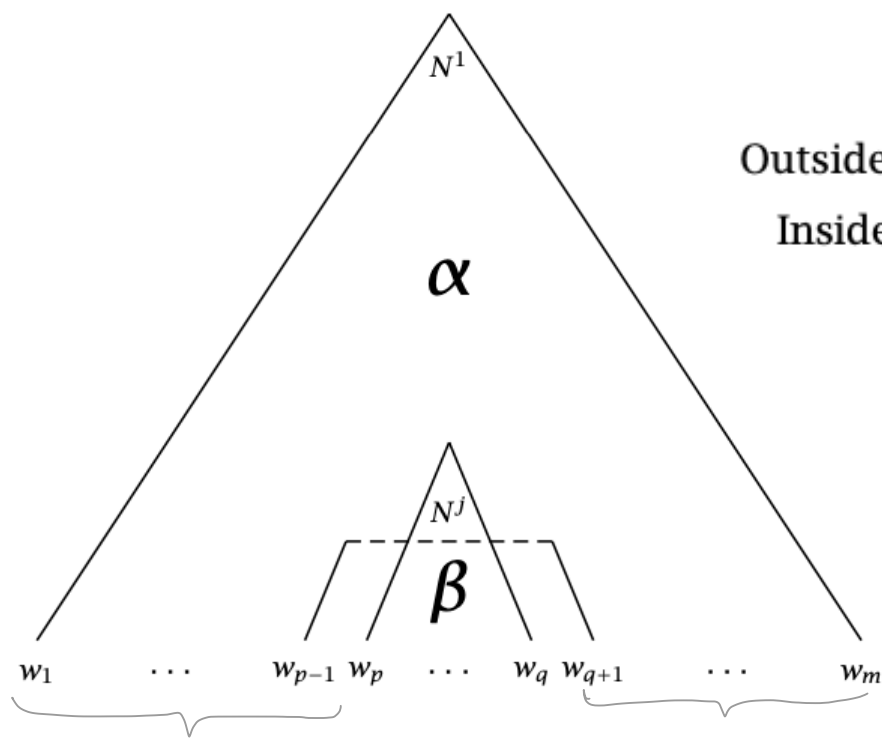


# From HMM to PCFG

- Assuming Chomsky Normal Form of the PCFG
- From HMM to PCFG
- Between HMM and PCFG, there is a grammar called Probabilistic Regular Grammar (PRG):  $N^i \rightarrow w^j N^k$  or  $N^i \rightarrow w^j$
- So what do the forward and backward probabilities correspond to on the tree?



# Inside and outside probabilities



Outside probability  $\alpha_j(p, q) = P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G)$

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

*start* *end*

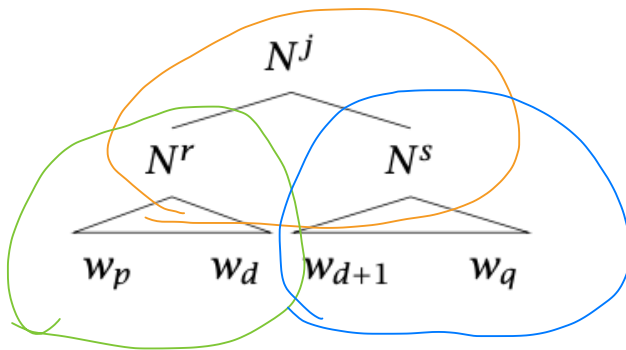
$$\begin{aligned} & \sum_j \alpha_j(p, q) \beta_j(p, q) \\ &= \sum_j \Pr(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G) \\ & \quad \times \Pr(w_{pq} | N_{pq}^j, G) \\ &= \sum_j \Pr(w_{1m}, N_{pq}^j | G) = \Pr(w_{1..m} | G) \end{aligned}$$

# Inside probabilities

Sentence probability derived from a PCFG  $P(w_{1m}|G) = P(N^1 \xRightarrow{*} w_{1m}|G)$

- Base case  $\beta_j(k,k) = P(w_k|N_{kk}^j, G) = P(w_{1m}|N_{1m}^1, G) = \beta_1(1, m)$   
 $\beta_j(k,k) = P(N^j \rightarrow w_k|G) \leftarrow \text{PCFG}$

- Induction due to CNF
  - CYK with probabilities



$$\begin{aligned}
 \beta_j(p, q) &= P(w_{pq}|N_{pq}^j, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} P(w_{pd}, N_{pd}^r, w_{(d+1)q}, N_{(d+1)q}^s | N_{pq}^j, G) \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} \underbrace{P(N_{pd}^r, N_{(d+1)q}^s | N_{pq}^j, G)}_{\text{PCFG}} \underbrace{P(w_{pd} | N_{pd}^r, N_{pd}^r, N_{(d+1)q}^s, G)}_{\text{PCFG}} \\
 &\quad \times \underbrace{P(w_{(d+1)q} | N_{pq}^j, N_{pd}^r, N_{(d+1)q}^s, w_{pd}, G)}_{\text{PCFG}} \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} \underbrace{P(N_{pd}^r, N_{(d+1)q}^s | N_{pq}^j, G)}_{\text{PCFG}} \underbrace{P(w_{pd} | N_{pd}^r, G)}_{\text{PCFG}} \\
 &\quad \times \underbrace{P(w_{(d+1)q} | N_{(d+1)q}^s, G)}_{\text{PCFG}} \\
 &= \sum_{r,s} \sum_{d=p}^{q-1} \underbrace{P(N^j \rightarrow N^r N^s)}_{\text{PCFG}} \underbrace{\beta_r(p, d)}_{\text{PCFG}} \underbrace{\beta_s(d+1, q)}_{\text{PCFG}}
 \end{aligned}$$

$$\beta^{NP}(3,5) = 0.18 \times 0.18 \times 0.4$$

# Inside probabilities

Compute inside probabilities using dynamic programming

S → NP VP	1.0	NP → NP PP	0.4
PP → P NP	1.0	NP → <i>astronomers</i>	0.1
VP → V NP	0.7	NP → <i>ears</i>	0.18
VP → VP PP	0.3	NP → <i>saw</i>	0.04
P → <i>with</i>	1.0	NP → <i>stars</i>	0.18
V → <i>saw</i>	1.0	NP → <i>telescopes</i>	0.1

$$C(3,5) = \{ N_{35}^S \rightarrow N_{3d}^r N_{(d+1)5}^s \mid G \}$$

$d=3$   $N_{35}^{NP} \rightarrow N_{33}^{NP} N_{45}^{PP}$   $\begin{pmatrix} s = NP \\ r = NP \\ s = PP \end{pmatrix}$   
 ~~$d=4$   $N_{35} \rightarrow N_{34} N_{55}$~~

	1	2	3	4	5
1	$\beta_{NP} = 0.1$		$\beta_S = 0.0126$		$\beta_S = 0.0015876$
2		$\beta_{NP} = 0.04$ $\beta_V = 1.0$	$\beta_{VP} = 0.126$		$\beta_{VP} = 0.015876$
3			$\beta_{NP} = 0.18$		$\beta_{NP} = 0.01296$
4				$\beta_P = 1.0$	$\beta_{PP} = 0.18$
5					$\beta_{NP} = 0.18$
	<i>astronomers</i>	<i>saw</i>	<i>stars</i>	<i>with</i>	<i>ears</i>

← Textbook  
 FSNLP  
 Lehigh Lib  
 ebook

HW5: compute cell (1,3)

