

# BMI/CS Lecture #22 - Stochastic Context Free Grammars for RNA Structure Modeling

Colin Dewey  
April 10, 2008

# Modeling RNA with Stochastic Context Free Grammars

- consider tRNA genes
  - 274 in yeast genome, ~1500 in human genome
  - get transcribed, like protein-coding genes
  - don't get translated, therefore base statistics much different than protein-coding genes
  - but secondary structure is conserved
- to recognize new tRNA genes, model known ones using stochastic context free grammars [Eddy & Durbin, 1994; Sakakibara et al. 1994]
- but what is a grammar?

# Transformational Grammars

- a transformational grammar characterizes a set of legal strings
- the grammar consists of
  - a set of abstract *nonterminal* symbols
$$\{S, c_1, c_2, c_3, c_4\}$$
  - a set of *terminal* symbols (those that actually appear in strings)
$$\{A, C, G, U\}$$
  - a set of *productions*

$$c_1 \rightarrow Uc_2 \qquad c_2 \rightarrow Ac_3 \qquad c_3 \rightarrow A$$

$$c_2 \rightarrow Gc_4 \qquad c_3 \rightarrow G$$

$$c_4 \rightarrow A$$

# A Grammar for Stop Codons

$$\begin{array}{llllll} s \rightarrow c_1 & c_1 \rightarrow Uc_2 & c_2 \rightarrow Ac_3 & c_3 \rightarrow A & c_4 \rightarrow A \\ & & c_2 \rightarrow Gc_4 & c_3 \rightarrow G & \end{array}$$

- this grammar can generate the 3 stop codons:  
UAA, UAG, UGA
- with a grammar we can ask questions like
  - what strings are derivable from the grammar?
  - can a particular string be derived from the grammar?

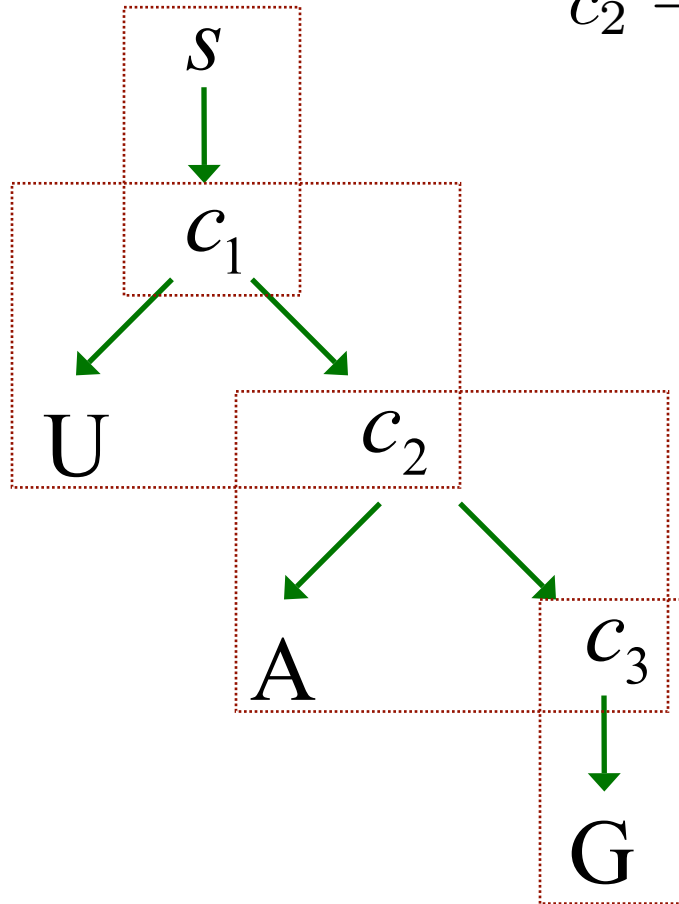
# The Derivation for UAG

$$\begin{array}{ccccccccc} s \rightarrow c_1 & c_1 \rightarrow Uc_2 & c_2 \rightarrow Ac_3 & c_3 \rightarrow A & c_4 \rightarrow A \\ & & c_2 \rightarrow Gc_4 & c_3 \rightarrow G & & & & & \end{array}$$

$$s \Rightarrow c_1 \Rightarrow Uc_2 \Rightarrow UAc_3 \Rightarrow UAG$$

# The Parse Tree for UAG

$$\begin{array}{ccccccccc} s \rightarrow c_1 & c_1 \rightarrow Uc_2 & c_2 \rightarrow Ac_3 & c_3 \rightarrow A & c_4 \rightarrow A \\ & & c_2 \rightarrow Gc_4 & c_3 \rightarrow G & & & & & \end{array}$$

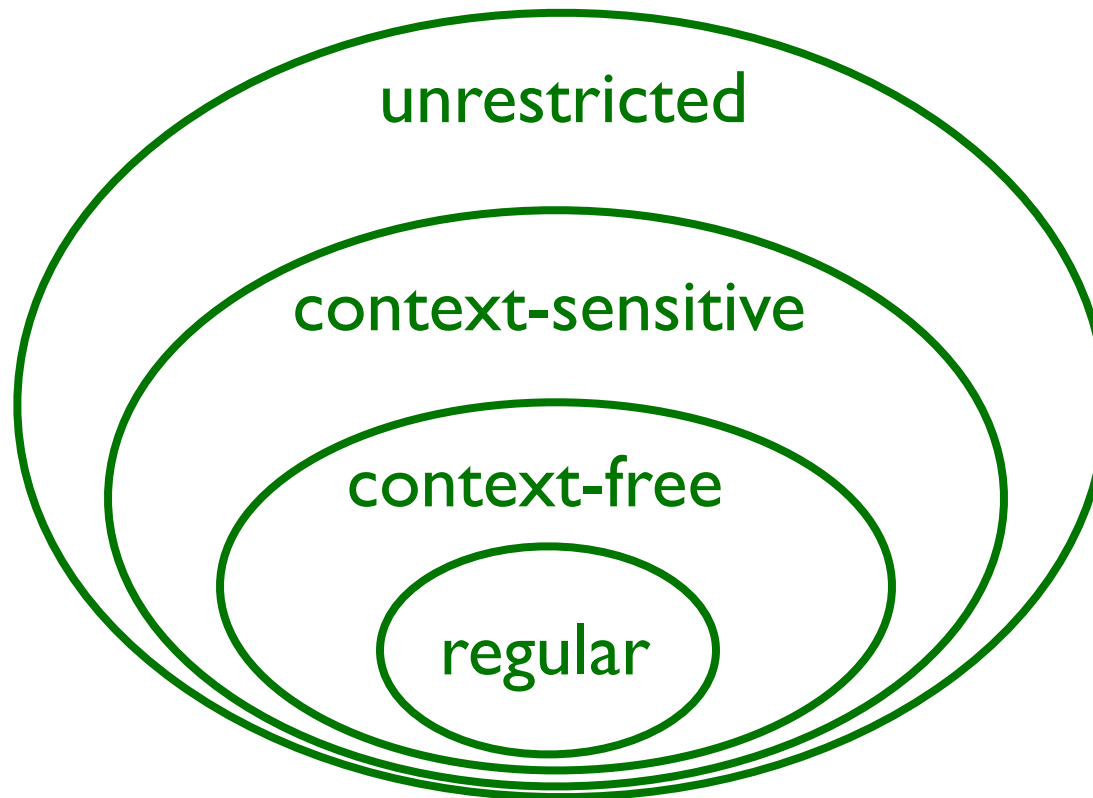


# A Probabilistic Version of the Grammar

$$\begin{array}{ccccc} \textcolor{red}{1.0} & \textcolor{red}{1.0} & \textcolor{red}{0.7} & \textcolor{red}{0.2} & \textcolor{red}{1.0} \\ s \rightarrow c_1 & c_1 \rightarrow Uc_2 & c_2 \rightarrow Ac_3 & c_3 \rightarrow A & c_4 \rightarrow A \\ & & \textcolor{red}{0.3} & \textcolor{red}{0.8} & \\ & & c_2 \rightarrow Gc_4 & c_3 \rightarrow G & \end{array}$$

- each production has an associated probability
- the probabilities for productions with the same left-hand side sum to one
- *this* grammar has a corresponding Markov chain model

# The Chomsky Hierarchy



- a hierarchy of grammars defined by restrictions on productions



# The Chomsky Hierarchy

- regular grammars  $u \rightarrow Xv \quad u \rightarrow X$
- context-free grammars  $u \rightarrow \beta$
- context-sensitive grammars  $\alpha_1 u \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$
- unrestricted grammars  $\alpha_1 u \alpha_2 \rightarrow \gamma$
- where:
  - $u$  is a nonterminal
  - $X$  is a terminal
  - $\alpha, \gamma$  are any sequence of terminals/nonterminals
  - $\beta$  is any non-null sequence of terminals/nonterminals

# CFGs and RNA

- context free grammars are well suited to modeling RNA secondary structure because they can represent base pairing preferences
- a grammar for a 3-base stem with and a loop of either **GAAA** or **GCAA**

$$\begin{aligned}s &\rightarrow Aw_1U \mid Cw_1G \mid Gw_1C \mid Uw_1A \\w_1 &\rightarrow Aw_2U \mid Cw_2G \mid Gw_2C \mid Uw_2A \\w_2 &\rightarrow Aw_3U \mid Cw_3G \mid Gw_3C \mid Uw_3A \\w_3 &\rightarrow GAAA \mid GCAA\end{aligned}$$

# CFGs and RNA

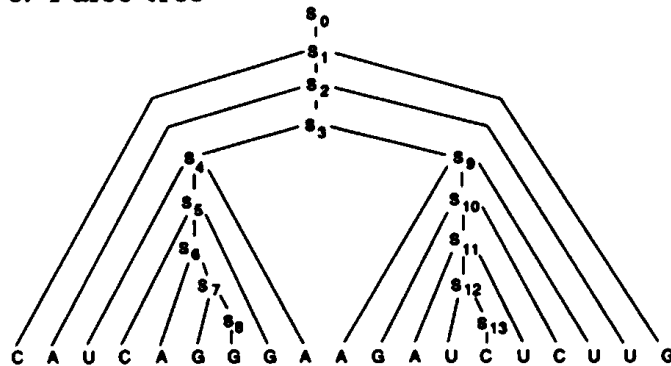
### a. Productions

$$P = \{ \begin{array}{ll} S_0 \rightarrow S_1, & S_7 \rightarrow G S_8, \\ S_1 \rightarrow C S_2 G, & S_8 \rightarrow G, \\ S_2 \rightarrow A S_3 U, & S_9 \rightarrow A S_{10} U, \\ S_3 \rightarrow S_4 S_9, & S_{10} \rightarrow G S_{11} C, \\ S_4 \rightarrow U S_5 A, & S_{11} \rightarrow A S_{12} U, \\ S_5 \rightarrow C S_6 G, & S_{12} \rightarrow U S_{13}, \\ S_6 \rightarrow A S_7, & S_{13} \rightarrow C \end{array} \}$$

### b. Derivation

$$\begin{aligned} S_0 &\Rightarrow S_1 \Rightarrow CS_2G \Rightarrow CAS_3UG \Rightarrow CAS_4S_9UG \\ &\Rightarrow CAUS_5AS_9UG \Rightarrow CAUCS_6GAS_9UG \\ &\Rightarrow CAUCAS_7GAS_9UG \Rightarrow CAUCAGS_8GAS_9UG \\ &\Rightarrow CAUCAGGAS_9UG \Rightarrow CAUCAGGGAAS_{10}UUG \\ &\Rightarrow CAUCAGGGAAGS_{11}CUUG \\ &\Rightarrow CAUCAGGGAAGAS_{12}UCUUG \\ &\Rightarrow CAUCAGGGAAGAU_{13}UCUUG \\ &\Rightarrow CAUCAGGGAAGAUUCUCUUG. \end{aligned}$$

c. Parse tree



#### d. Secondary Structure

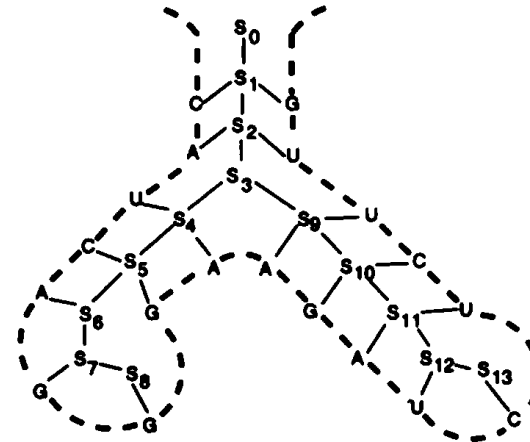
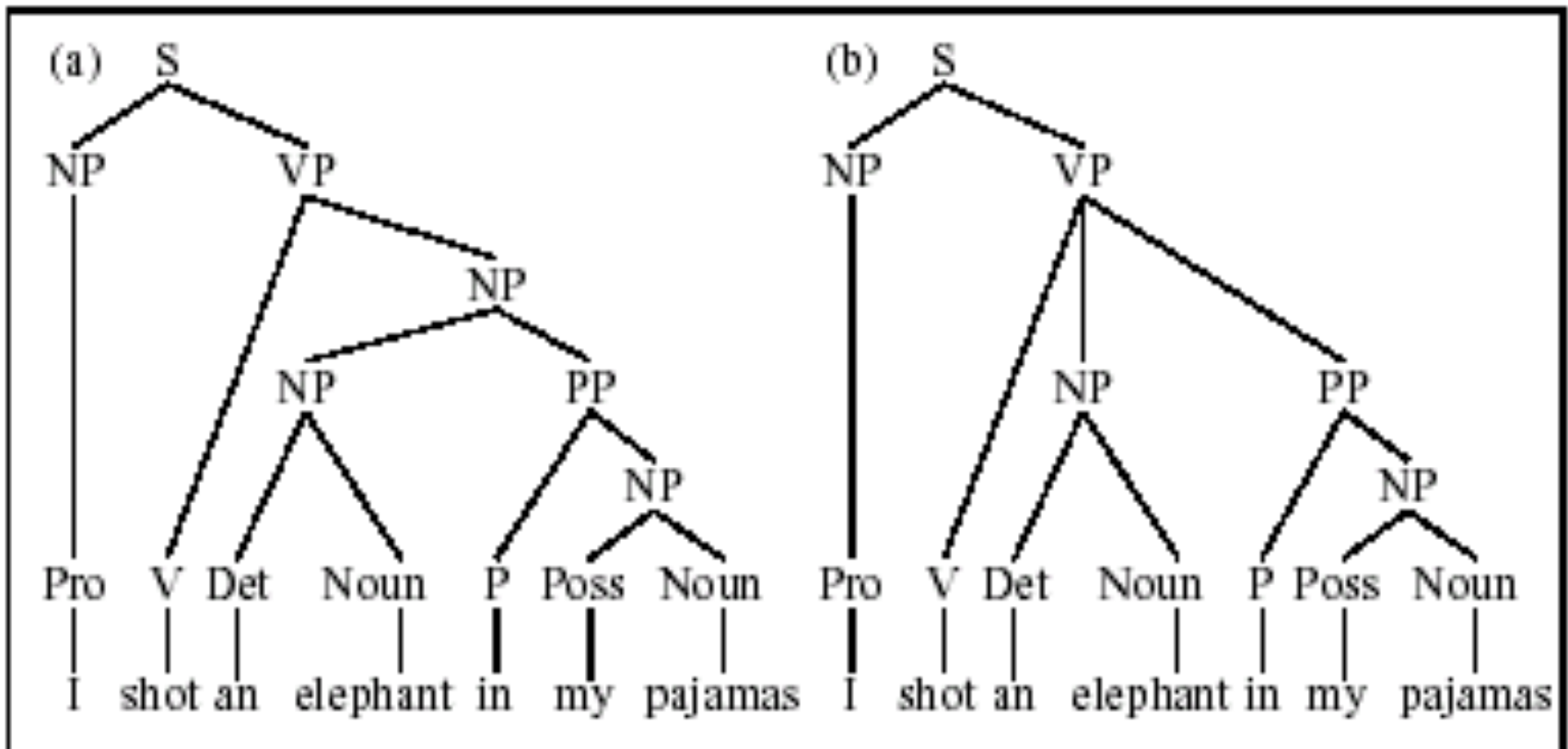


Figure from: Sakakibara et al. *Nucleic Acids Research*, 1994

# Ambiguity in Parsing

“I shot an elephant in my pajamas. How he got in my pajamas, I’ll never know.” – Groucho Marx



# Ambiguity in Parsing

$$s \rightarrow Aw_1U \mid Cw_1G \mid Gw_1C \mid Uw_1A$$

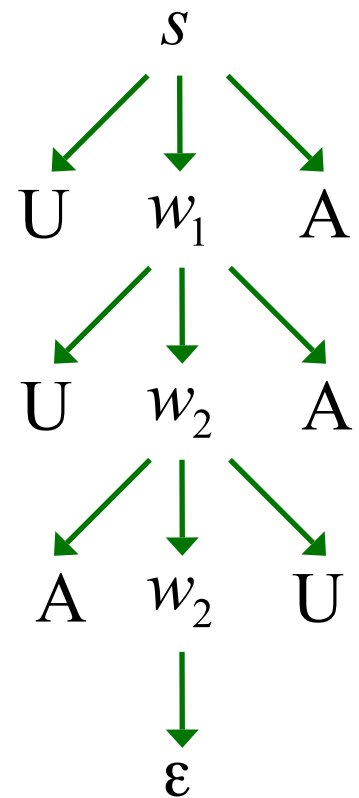
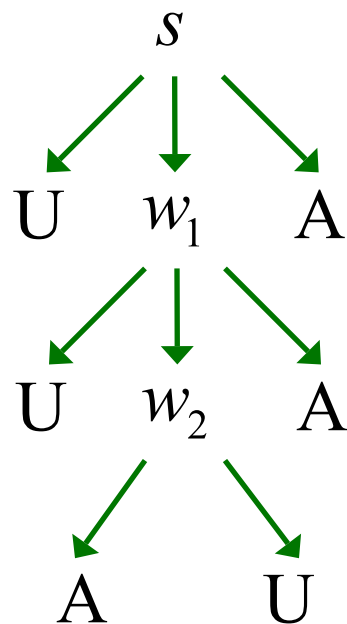
$$w_1 \rightarrow Aw_2U \mid Cw_2G \mid Gw_2C \mid Uw_2A$$

$$w_2 \rightarrow Aw_2U \mid Cw_2G \mid Gw_2C \mid Uw_2A$$

$$w_2 \rightarrow AA \mid AU$$

$$w_2 \rightarrow \epsilon$$

denotes the  
empty string



# Stochastic Context Free Grammars

$$\begin{array}{cccc} 0.25 & 0.25 & 0.25 & 0.25 \\ s \rightarrow Aw_1U & | & Cw_1G & | & Gw_1C & | & Uw_1A \end{array}$$

$$\begin{array}{cccc} 0.1 & 0.4 & 0.4 & 0.1 \\ w_1 \rightarrow Aw_2U & | & Cw_2G & | & Gw_2C & | & Uw_2A \end{array}$$

$$\begin{array}{cccc} 0.25 & 0.25 & 0.25 & 0.25 \\ w_2 \rightarrow Aw_3U & | & Cw_3G & | & Gw_3C & | & Uw_3A \end{array}$$

$$\begin{array}{cc} 0.8 & 0.2 \\ w_3 \rightarrow GAAA & | & GCAA \end{array}$$

# Stochastic Grammars?

*...the notion “probability of a sentence” is an entirely useless one, under any known interpretation of this term.*

— Noam Chomsky  
(famed linguist)

*Every time I fire a linguist, the performance of the recognizer improves.*

— Fred Jelinek  
(former head of IBM speech recognition group)

Credit for pairing these quotes goes to Dan Jurafsky and James Martin,  
*Speech and Language Processing*

# Three Key Questions

- How likely is a given sequence?

the Inside algorithm

- What is the most probable parse for a given sequence?

the Cocke-Younger-Kasami (CYK) algorithm

- How can we learn the SCFG parameters given a grammar and a set of sequences?

the Inside-Outside algorithm

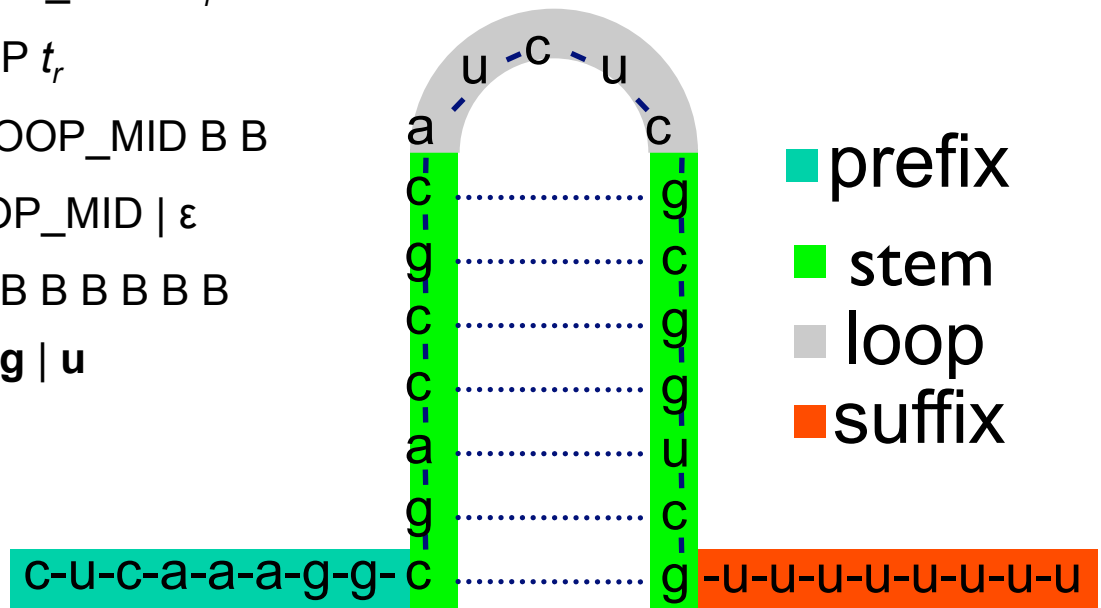


# A Stochastic Context Free Grammar for Terminators

START	→	PREFIX STEM_BOT1 SUFFIX
PREFIX	→	B B B B B B B B
STEM_BOT1	→	$t_l$ STEM_BOT2 $t_r$
STEM_BOT2	→	$t_l^*$ STEM_MID $t_r^*$   $t_l^*$ STEM_TOP2 $t_r^*$
STEM_MID	→	$t_l^*$ STEM_MID $t_r^*$   $t_l^*$ STEM_TOP2 $t_r^*$
STEM_TOP2	→	$t_l^*$ STEM_TOP1 $t_r^*$
STEM_TOP1	→	$t_l$ LOOP $t_r$
LOOP	→	B B LOOP_MID B B
LOOP_MID	→	B LOOP_MID   $\epsilon$
SUFFIX	→	B B B B B B B B
B	→	a   c   g   u

$t = \{a, c, g, u\},$   
 $t^* = \{a, c, g, u, \epsilon\}$

STEM_BOT1	→	0.2	a STEM_BOT2 u
STEM_BOT1	→	0.2	u STEM_BOT2 a
STEM_BOT1	→	0.2	c STEM_BOT2 g
STEM_BOT1	→	0.2	g STEM_BOT2 c
STEM_BOT1	→	0.05	g STEM_BOT2 u
		⋮	



# Chomsky Normal Form

- it is convenient to assume that our grammar is in *Chomsky Normal Form*; i.e all productions are of the form:

$v \rightarrow yz$       right hand side consists of two nonterminals

$v \rightarrow A$       right hand side consists of a single terminal

- any CFG can be put into Chomsky Normal Form

# Parameter Notation

- for productions of the form  $v \rightarrow yz$ , we'll denote the associated probability parameters

$t_v(y, z)$       **transition**

- for productions of the form  $v \rightarrow A$ , we'll denote the associated probability parameters

$e_v(A)$       **emission**

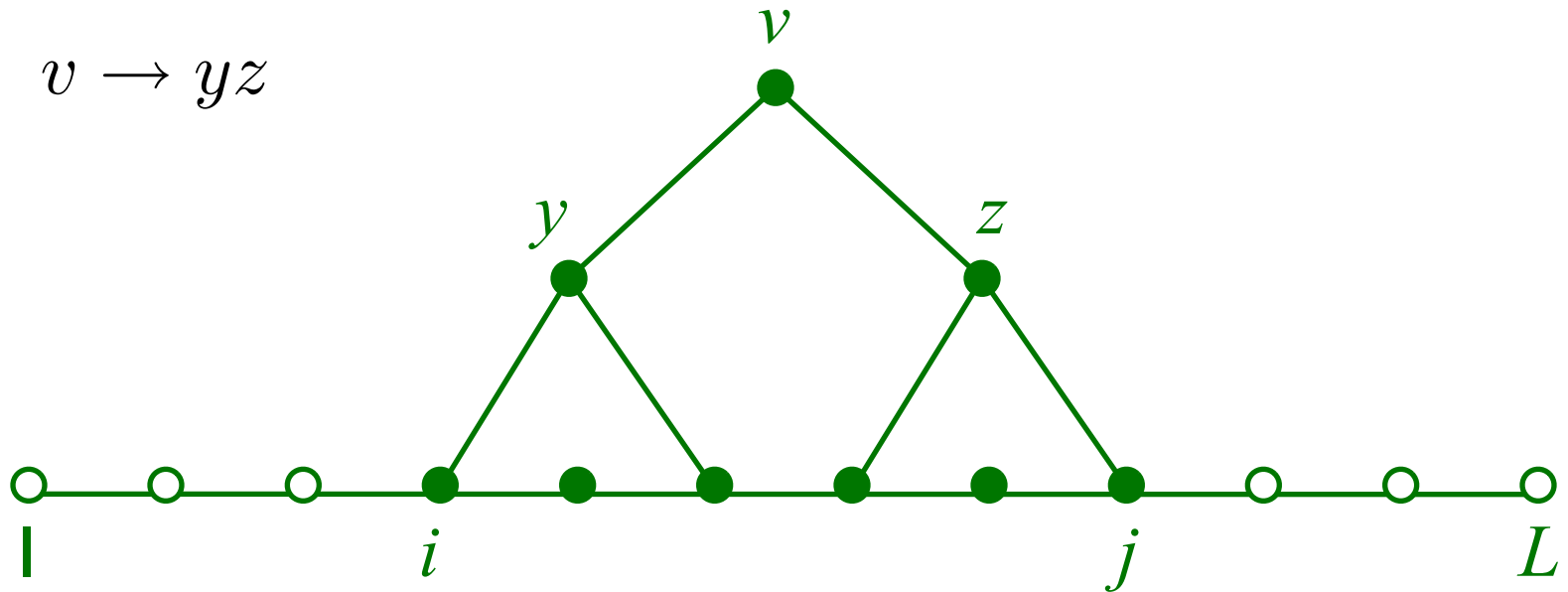
# Determining the Likelihood of a Sequence: The Inside Algorithm

- a dynamic programming method, analogous to the Forward algorithm
- involves filling in a 3D matrix

$$\alpha(i, j, v)$$

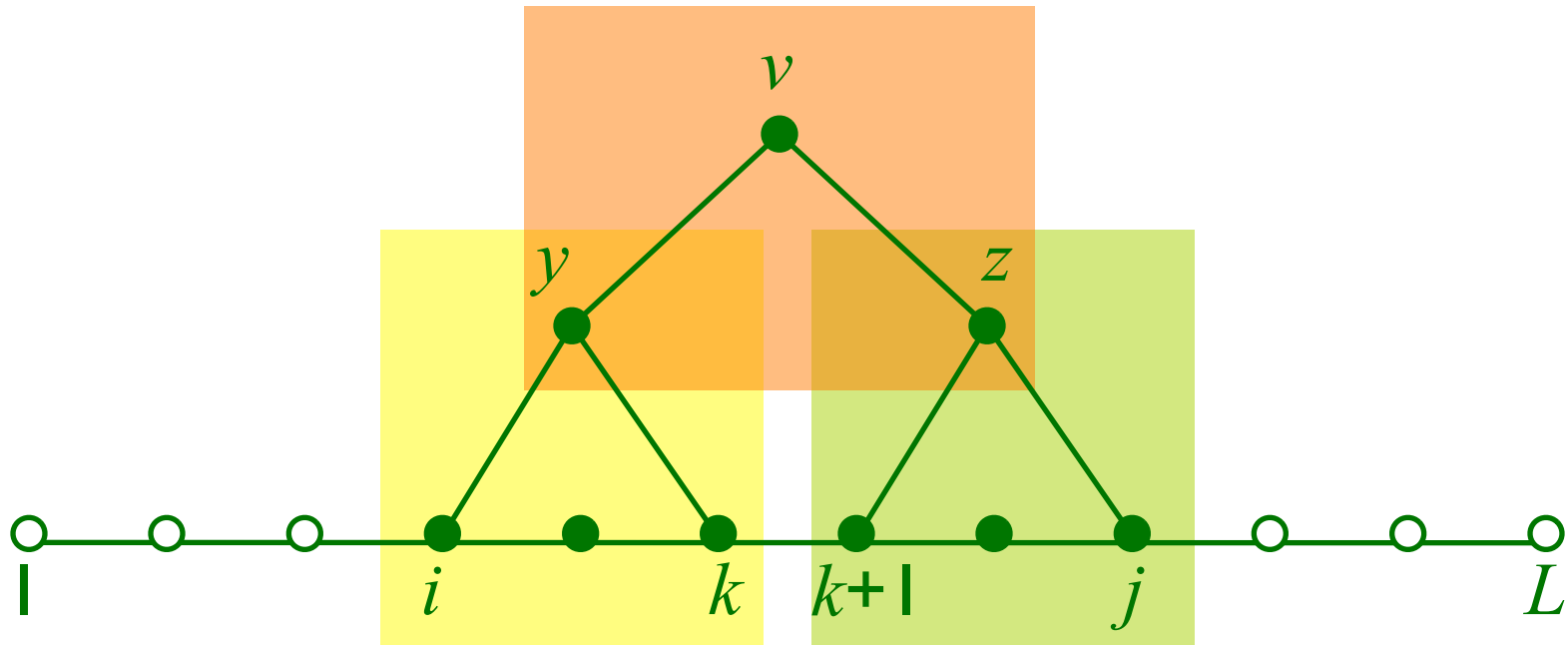
representing the probability of the all parse subtrees rooted at nonterminal  $v$  for the subsequence from  $i$  to  $j$

# Determining the Likelihood of a Sequence: The Inside Algorithm



- $\alpha(i, j, v)$  : the probability of all parse subtrees rooted at nonterminal  $v$  for the subsequence from  $i$  to  $j$

# Determining the Likelihood of a Sequence: The Inside Algorithm



$$\alpha(i, j, v) = \sum_{y=1}^M \sum_{z=1}^M \sum_{k=i}^{j-1} \alpha(i, k, y) \alpha(k+1, j, z) t_v(y, z)$$

$M$  is the number of nonterminals in the grammar

# The Inside Algorithm

- initialization (for  $i = 1$  to  $L$ ,  $v = 1$  to  $M$ )

$$\alpha(i, i, v) = e_v(x_i)$$

- iteration (for  $i = L-1$  to  $1$ ,  $j = i+1$  to  $L$ ,  $v = 1$  to  $M$ )

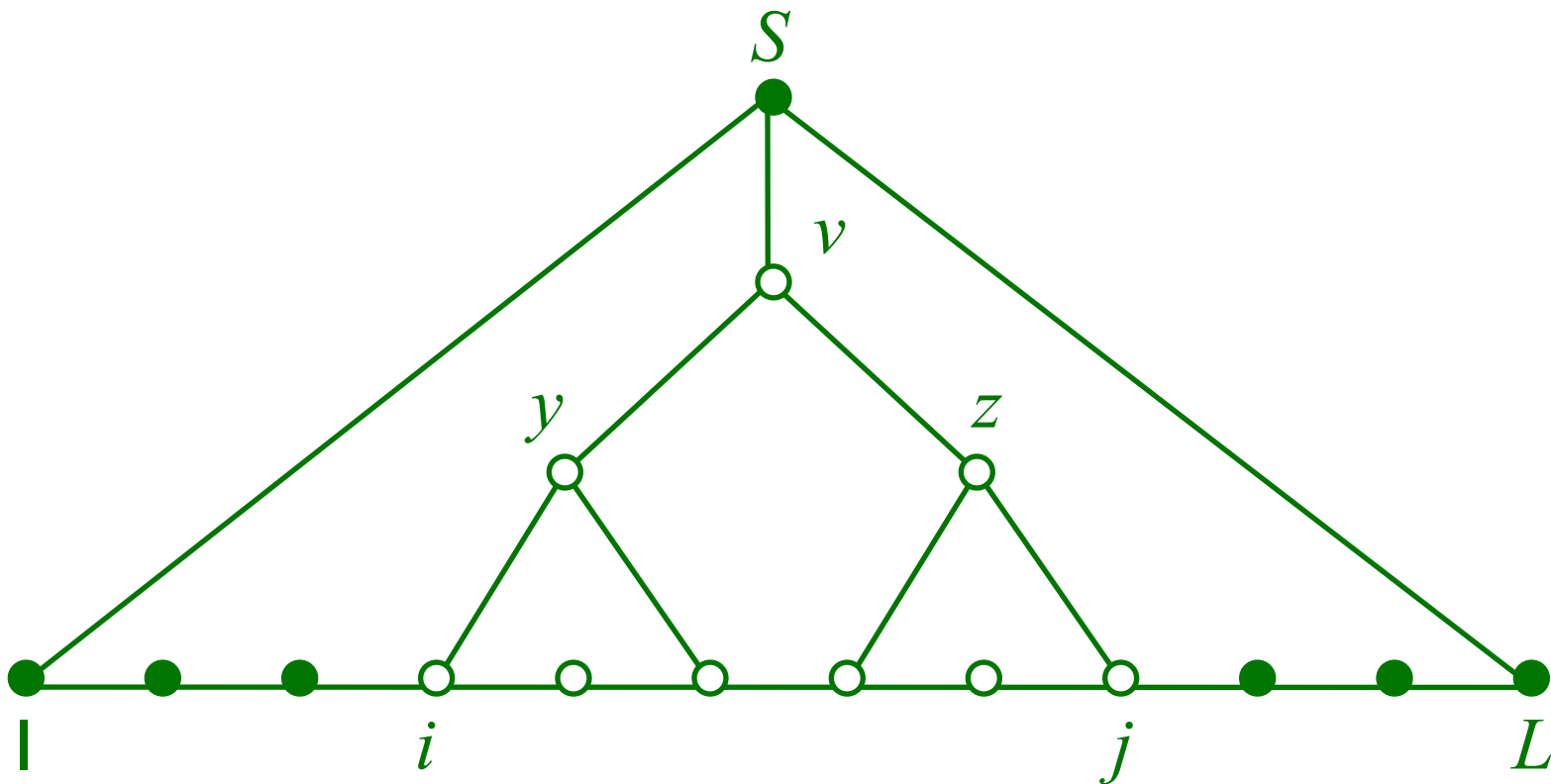
$$\alpha(i, j, v) = \sum_{y=1}^M \sum_{z=1}^M \sum_{k=i}^{j-1} \alpha(i, k, y) \alpha(k+1, j, z) t_v(y, z)$$

- termination

$$\Pr(x) = \alpha(1, L, 1)$$

↑  
start nonterminal

# The Outside Algorithm

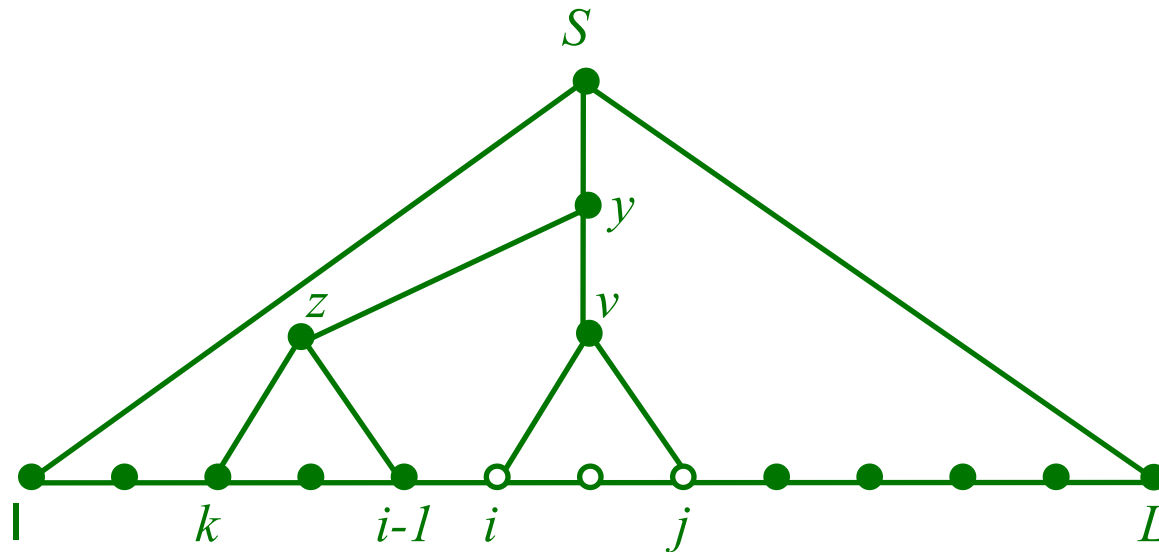


- $\beta(i, j, v)$  : the probability of parse trees rooted at the start nonterminal, excluding the probability of all subtrees rooted at nonterminal  $v$  covering the subsequence from  $i$  to  $j$



# The Outside Algorithm

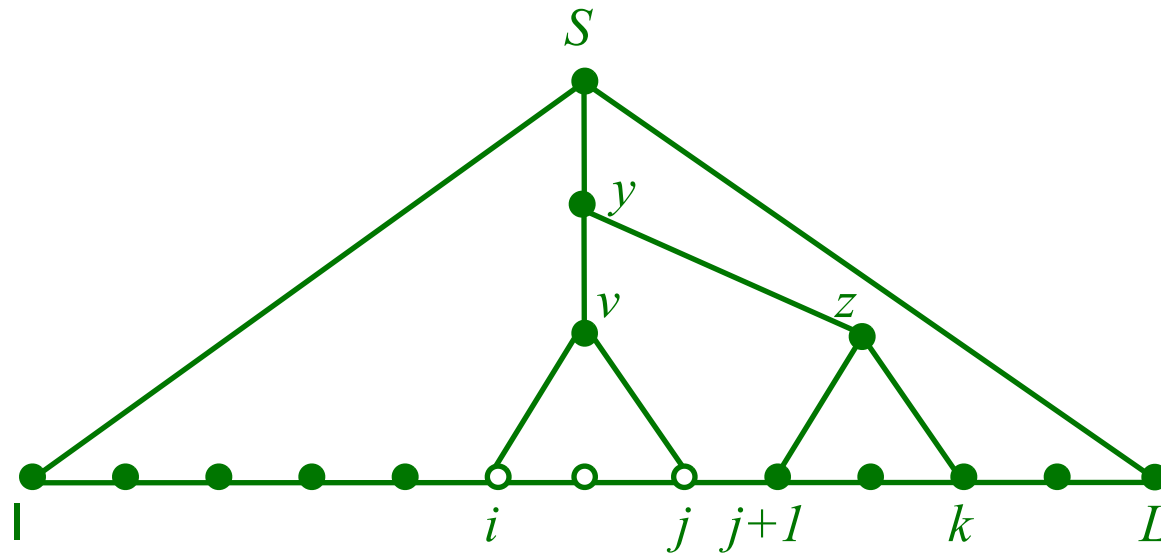
- we can recursively calculate  $\beta(i, j, v)$  from  $\beta$  values we've calculated for  $y$
- the first case we consider is where  $v$  is used in productions of the form:  $y \rightarrow zv$



$$\sum_{y=1}^M \sum_{z=1}^M \sum_{k=1}^{i-1} \alpha(k, i-1, z) \beta(k, j, y) t_y(z, v)$$

# The Outside Algorithm

- the second case we consider is where  $v$  is used in productions of the form:  $y \rightarrow vz$



$$\sum_{y=1}^M \sum_{z=1}^M \sum_{k=j+1}^L \alpha(j+1, k, z) \beta(i, k, y) t_y(v, z)$$

# The Outside Algorithm

- initialization

$$\beta(1, L, 1) = 1 \quad (\text{the } \textit{start} \text{ nonterminal})$$

$$\beta(1, L, v) = 0 \quad \text{for } v = 2 \text{ to } M$$

- iteration (for  $i = 1$  to  $L, j = L$  to  $i, v = 1$  to  $M$ )

$$\begin{aligned} \beta(i, j, v) = & \sum_{y=1}^M \sum_{z=1}^M \sum_{k=1}^{i-1} \alpha(k, i-1, z) \beta(k, j, y) t_y(z, v) + \\ & \sum_{y=1}^M \sum_{z=1}^M \sum_{k=j+1}^L \alpha(j+1, k, z) \beta(i, k, y) t_y(v, z) \end{aligned}$$

# Learning SCFG Parameters

- if we know the parse tree for each training sequence, learning the SCFG parameters is simple
  - no hidden state during training
  - count how often each parameter (i.e. production) is used
  - normalize/smooth to get probabilities
- more commonly, there are many possible parse trees per sequence – we don't know which one is correct
  - thus, use an EM approach (Inside-Outside)
  - iteratively
    - determine expected # times each production is used
    - consider all parses
    - weight each by its probability
    - set parameters to maximize these counts

# The Inside-Outside Algorithm

- we can learn the parameters of an SCFG from training sequences using an EM approach called Inside-Outside
- in the E-step, we determine  $c(v)$ 
  - the expected number of times each *nonterminal* is used in parses  $c(v \rightarrow yz)$
  - the expected number of times each *production* is used in parses  $c(v \rightarrow A)$
- in the M-step, we update our production probabilities

# The Inside-Outside Algorithm

- the EM re-estimation equations (for 1 sequence) are:

$$\hat{e}_v(A) = \frac{c(v \rightarrow A)}{c(v)} = \frac{\sum_{i|x_i=A} \beta(i, i, v) e_v(A)}{\sum_{i=1}^L \sum_{j=i}^L \beta(i, j, v) \alpha(i, j, v)}$$

← cases where  $v$  used to generate  $A$

← cases where  $v$  used to generate any subsequence

$$\hat{t}_v(y, z) = \frac{c(v \rightarrow yz)}{c(v)}$$

$$= \frac{\sum_{i=1}^{L-1} \sum_{j=i+1}^L \sum_{k=i}^{j-1} \beta(i, j, v) t_v(y, z) \alpha(i, k, y) \alpha(k+1, j, z)}{\sum_{i=1}^L \sum_{j=i}^L \beta(i, j, v) \alpha(i, j, v)}$$

# The CYK Algorithm

- analogous to Viterbi algorithm
- like Inside algorithm but
  - max operations instead of sums
  - retain traceback pointers
- traceback is a little more involved than Viterbi
  - need to reconstruct parse tree instead of recovering simple path

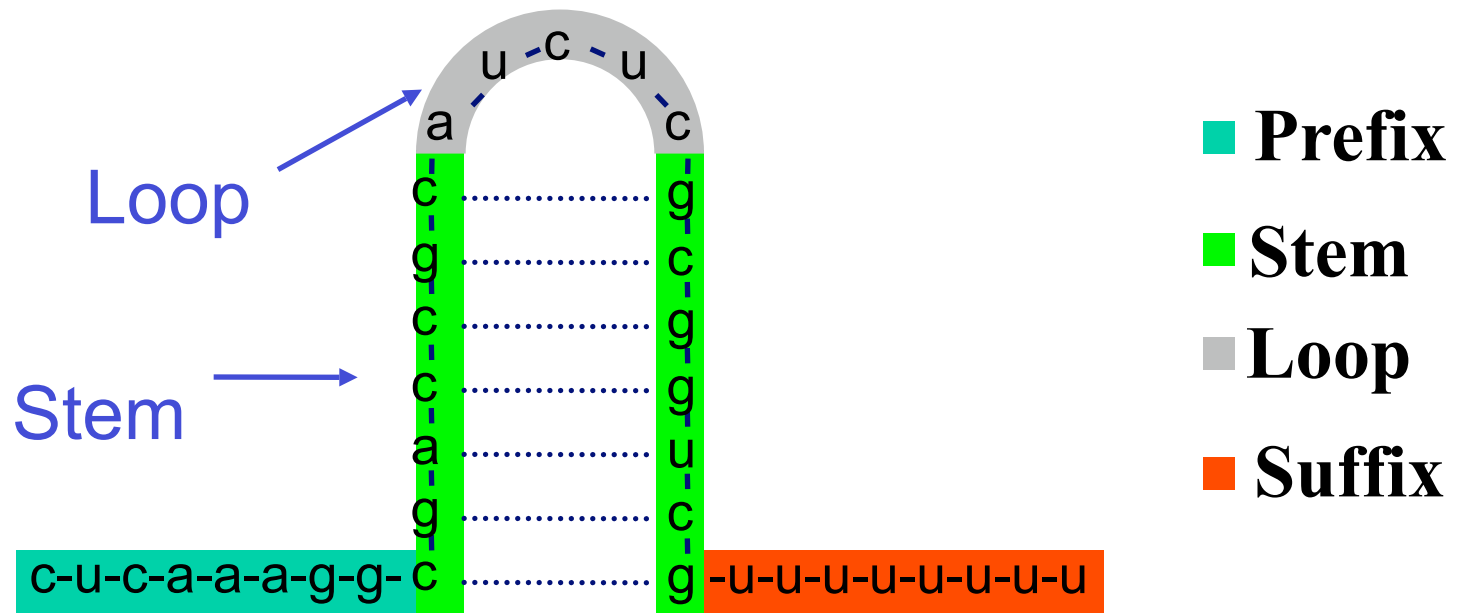
# Comparison of SCFG Algorithms to HMM Algorithms

	HMM algorithm	SCFG algorithm
optimal alignment	Viterbi	CYK
probability of sequence	forward	inside
EM parameter estimation	forward-backward	inside-outside
memory complexity	$O(LM)$	$O(L^2M)$
time complexity	$O(LM^2)$	$O(L^3M^3)$



# Recognizing Terminators with SCFGs

- [Bockhorst & Craven, *IJCAI* 2001]



- a prototypical terminator has the structure above
- the lengths and base compositions of the elements can vary a fair amount

# Our Initial Terminator Grammar

START	→	PREFIX STEM_BOT1 SUFFIX
PREFIX	→	B B B B B B B B B
STEM_BOT1	→	$t_l$ STEM_BOT2 $t_r$
STEM_BOT2	→	$t_l^*$ STEM_MID $t_r^*$   $t_l^*$ STEM_TOP2 $t_r^*$
STEM_MID	→	$t_l^*$ STEM_MID $t_r^*$   $t_l^*$ STEM_TOP2 $t_r^*$
STEM_TOP2	→	$t_l^*$ STEM_TOP1 $t_r^*$
STEM_TOP1	→	$t_l$ LOOP $t_r$
LOOP	→	B B LOOP_MID B B
LOOP_MID	→	B LOOP_MID   $\epsilon$
SUFFIX	→	B B B B B B B B B
B	→	a   c   g   u

Nonterminals are uppercase,  
terminals are lowercase

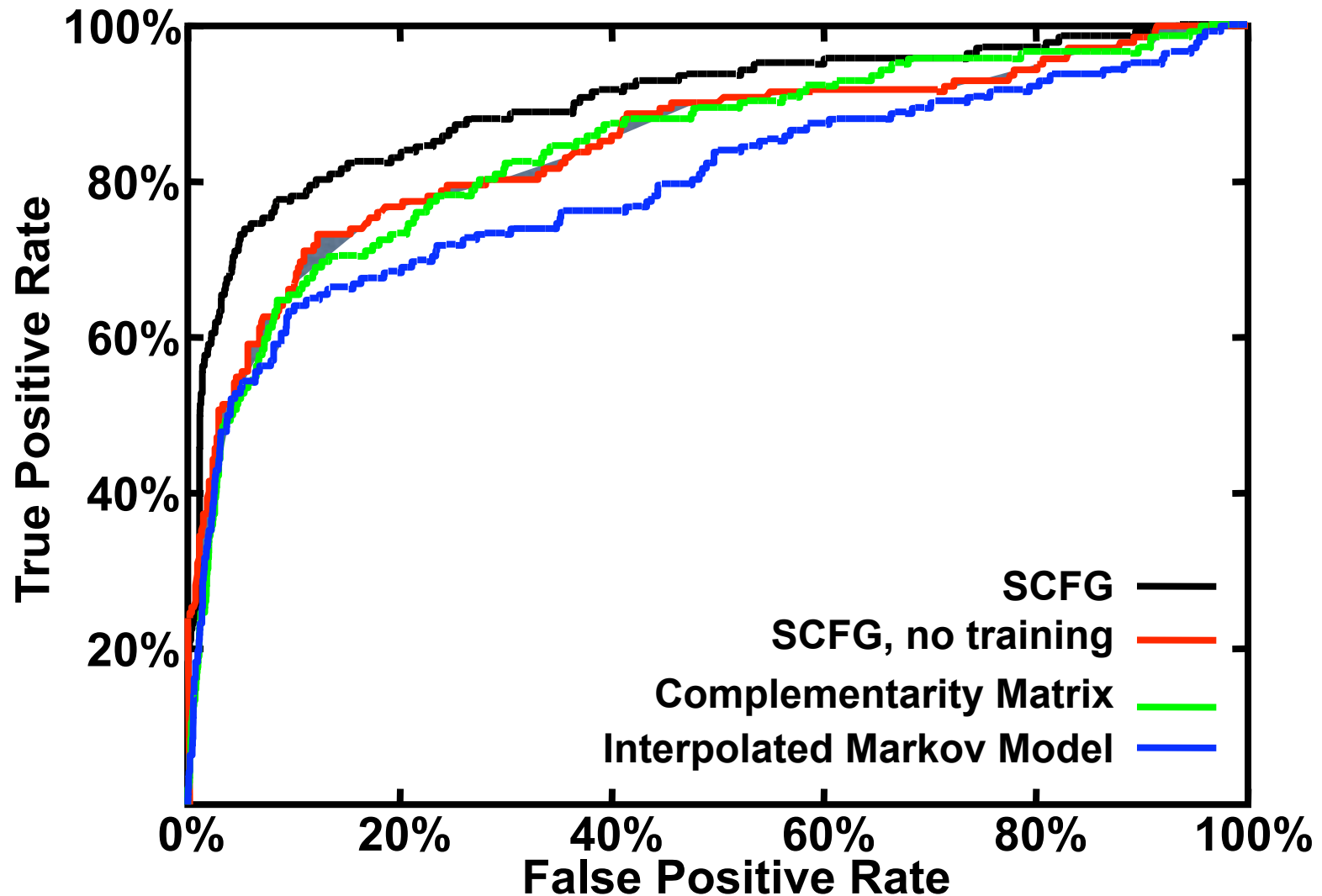
$$t = \{a, c, g, u\},$$

$$t^* = \{a, c, g, u, \epsilon\}$$

# Terminator SCFG Experiments

- compare predictive accuracy of
  - SCFG with learned parameters
  - SCFG without learning (but parameters initialized using domain knowledge)
  - interpolated Markov models (IMMs)
    - can represent distribution of bases at each position
      - \* cannot easily encode base pair dependencies
  - complementarity matrices
    - Brendel et al., *J Biom Struct and Dyn* 1986
    - ad hoc way of considering base pairings
      - \* cannot favor specific base pairs by position

# SCFGs vs. Related Methods



# Learning SCFG Structure

- given the productions of a grammar, can learn the probabilities using the Inside-Outside algorithm
- we have developed an algorithm that can add new nonterminals & productions to a grammar during learning [Bockhorst & Craven, *IJCAI 01*]
- basic idea:
  - identify nonterminals that seem to be “overloaded”
  - split these nonterminals into two; allow each to specialize

# Refinement Algorithm Overview

Given:

- set of sequences
- initial grammar structure hypothesis

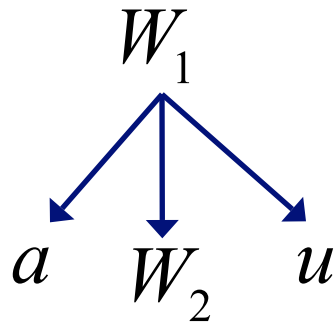
Do:

repeat

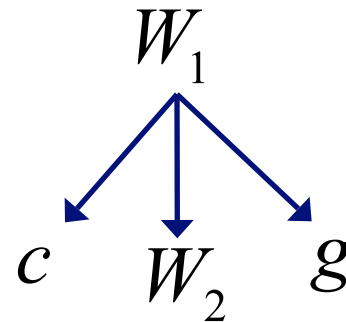
- 1) find MAP estimates for probabilities
- 2) refine grammar structure
  - 2.1) diagnostically identify overloaded nonterminal
  - 2.2) apply **EXPAND** operator

# Refining the Grammar in a SCFG

- there are various “contexts” in which each grammar nonterminal may be used
- consider two contexts for the nonterminal  $W_2$



$$\begin{array}{l}
 W_2 \rightarrow aW_3u \mid 0.4 \\
 \quad cW_3g \mid 0.1 \\
 \quad gW_3c \mid 0.1 \\
 \quad uW_3a \quad 0.4
 \end{array}$$



$$\begin{array}{l}
 W_2 \rightarrow aW_3u \mid 0.1 \\
 \quad cW_3g \mid 0.4 \\
 \quad gW_3c \mid 0.4 \\
 \quad uW_3a \quad 0.1
 \end{array}$$

- if the probabilities for  $W_2$  look very different, depending on its context, we add a new nonterminal and specialize

# Refining the Grammar in a SCFG

- we can compare two probability distributions  $P$  and  $Q$  using Kullback-Leibler divergence

$$H(P||Q) = \sum_i P(x_i) \log_2 \frac{P(x_i)}{Q(x_i)}$$

	$P$		$Q$
$W_2 \rightarrow aW_3u$	0.4	$W_2 \rightarrow aW_3u$	0.1
$cW_3g$	0.1	$cW_3g$	0.4
$gW_3c$	0.1	$gW_3c$	0.4
$uW_3a$	0.4	$uW_3a$	0.1

- or we can compare expected number of times each production is used (over training set) using  $\chi^2$



# EXPAND Operator

EXPAND(P,N) /\* N on RHS of production P \*/

1) create new nonterminal N'

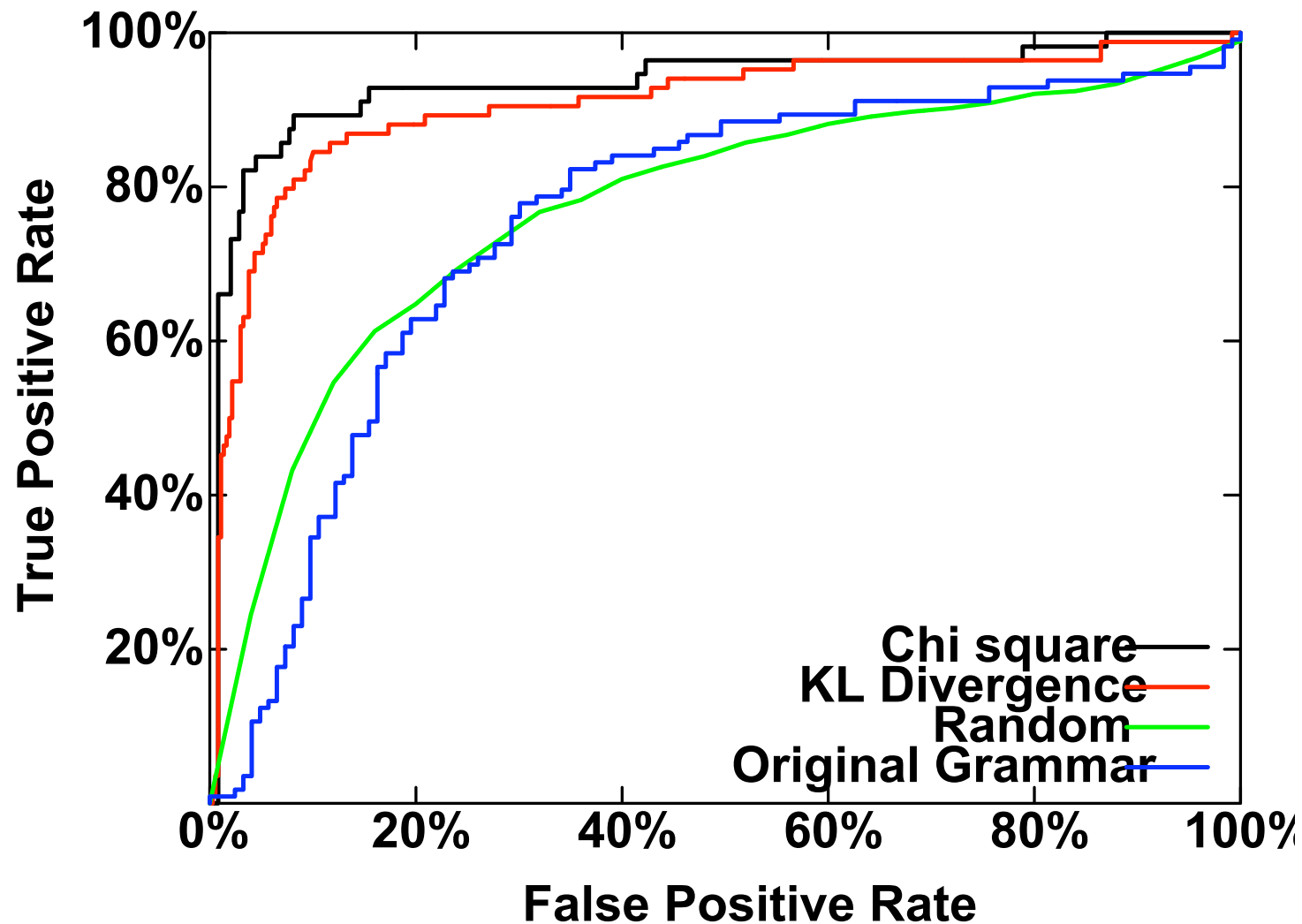
2) replace N with N' in P.

3) for each production  $N \rightarrow \bullet$ , create a production  $N' \rightarrow \bullet$

# Learning Terminator SCFGs

- extracted grammar from the literature (~ 120 productions)
- data set consists of 142 known *E. coli* terminators, 125 sequences that do not contain terminators
- learn parameters using Inside-Outside algorithm (an EM algorithm)
- consider adding nonterminals guided by three heuristics
  - KL divergence
  - chi-squared
  - random

# SCFG Accuracy After Adding 25 New Nonterminals



# SCFG Accuracy vs. Nonterminals Added

