
CS681: Advanced Topics in Computational Biology

Week 10 Lecture 1

Can Alkan

EA224

calkan@cs.bilkent.edu.tr

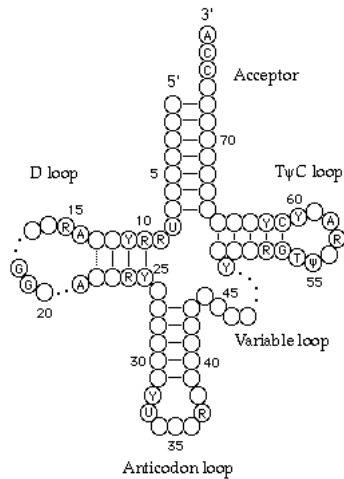
<http://www.cs.bilkent.edu.tr/~calkan/teaching/cs681/>

RNA folding

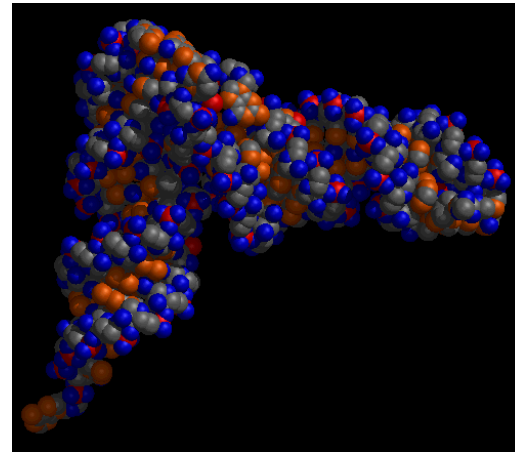
- Prediction of secondary structure of an RNA given its sequence
 - General problem is NP-hard due to “difficult” substructures, like pseudoknots
 - Most existing algorithms require too much memory ($\geq O(n^2)$), and run time ($\geq O(n^3)$) thus limited to smaller RNA sequences
-

RNA Structural Levels

AAUCG...CUUCUCCA
Primary

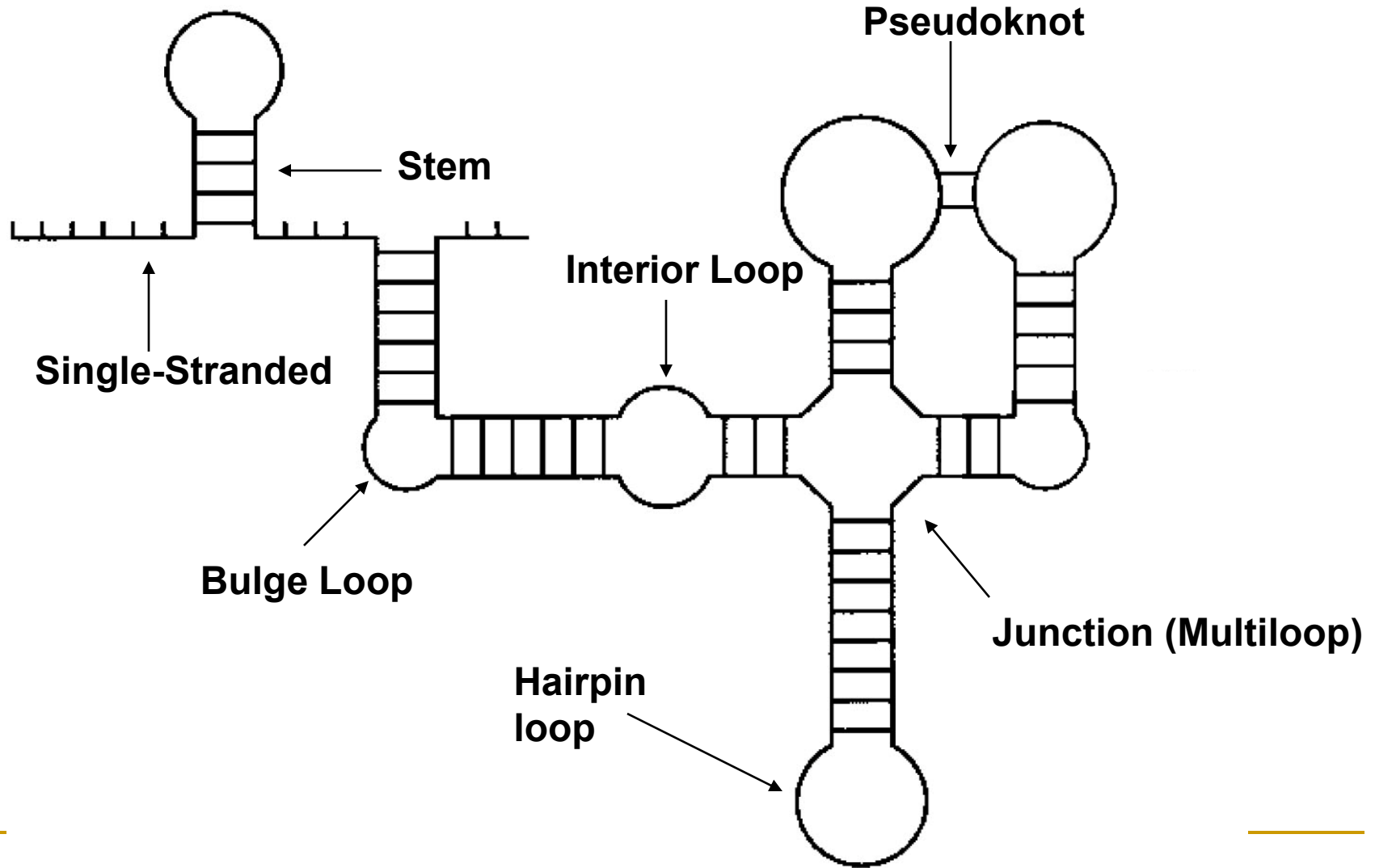


Secondary



Tertiary

RNA Secondary Structure



Predicting RNA secondary structure

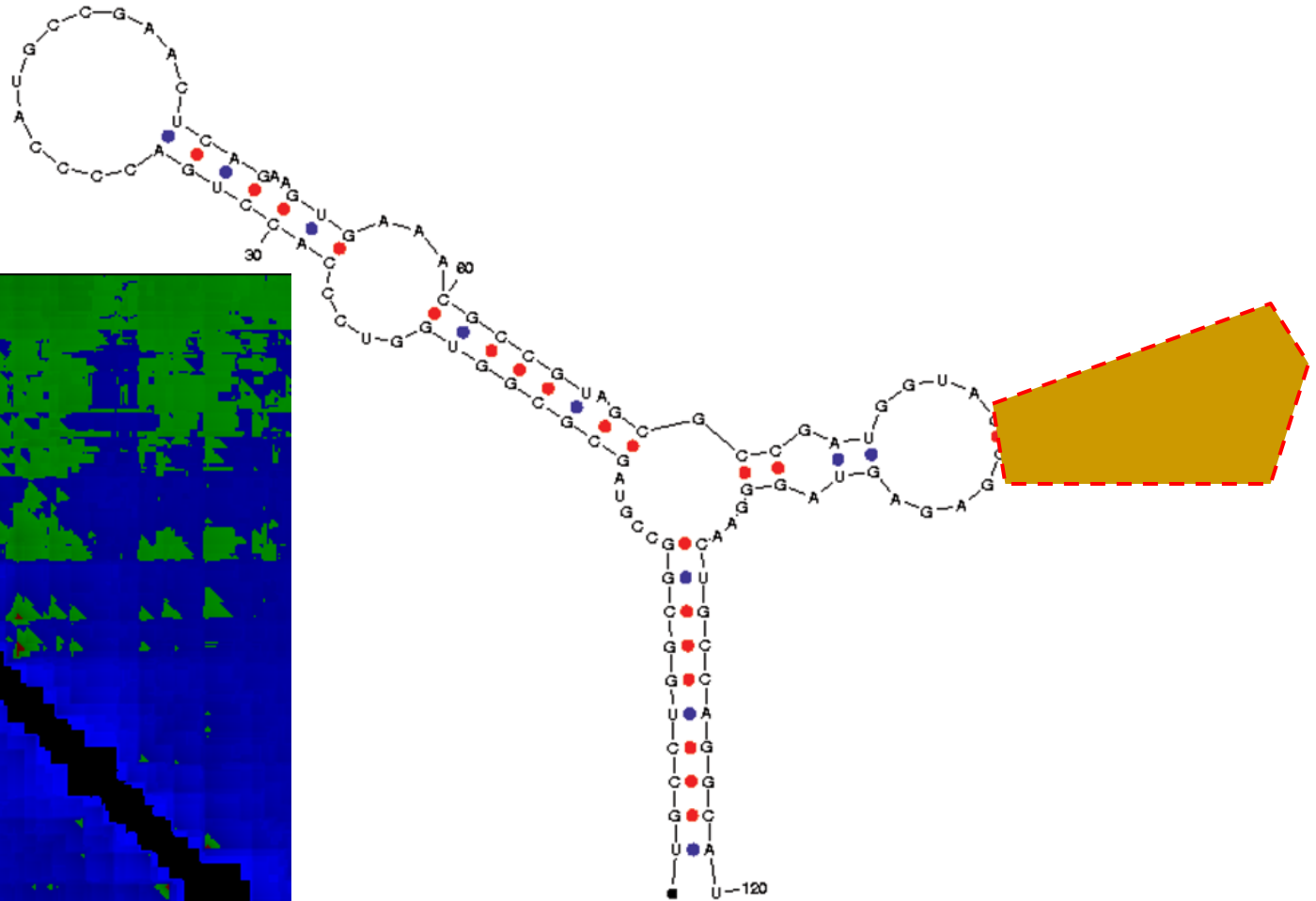
- Base pair maximization
 - Minimum free energy (most common)
 - Fold, Mfold (Zuker & Stiegler)
 - RNAfold (Hofacker)
 - Multiple sequence alignment
 - Use known structure of RNA with similar sequence
 - Covariance
 - Stochastic Context-Free Grammars
-

Alkan, Karakoç et al, RECOMB 2006

DENSITYFOLD

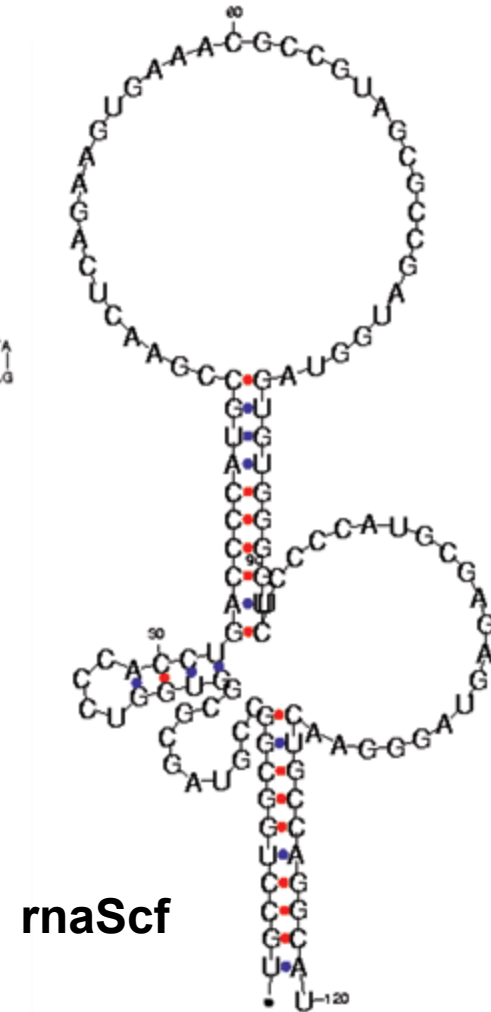
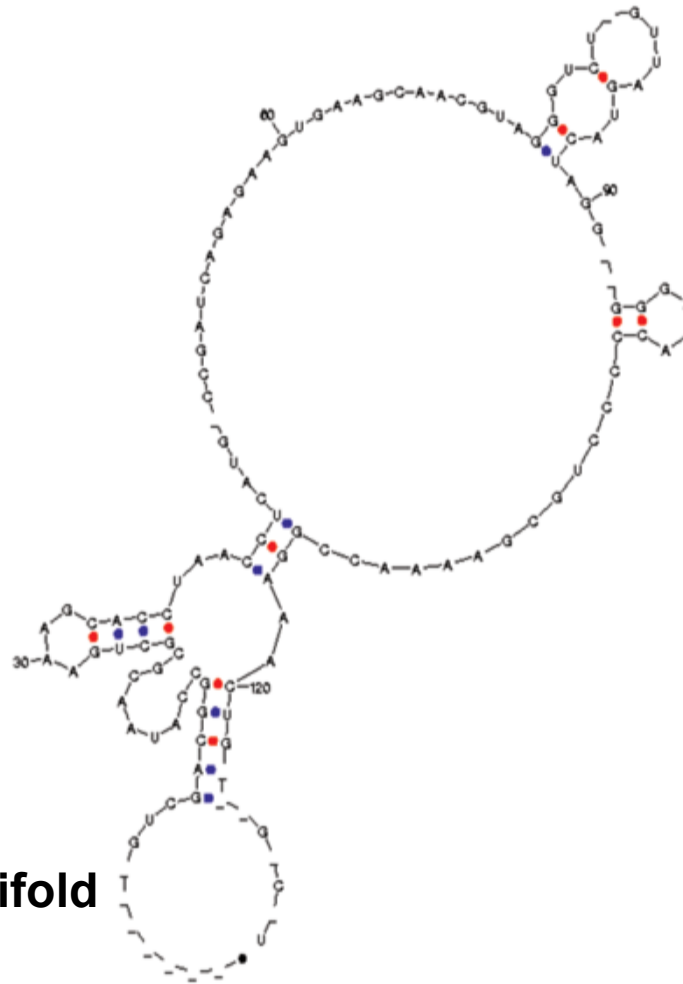
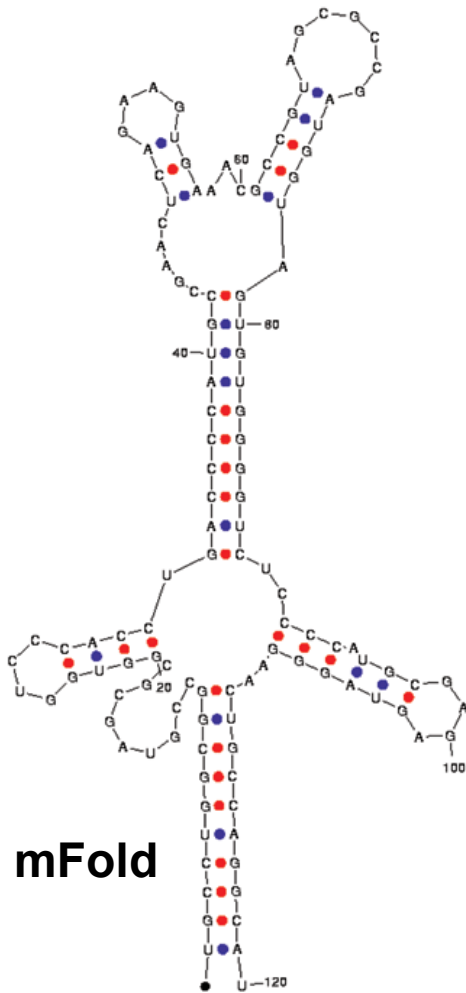


E.coli 5S rRNA



Energy Density Landscape

E.coli 5S rRNA predictions



Densityfold (alteRNA)

- Instead of finding minimum global free energy, find local minimum free energies
 - Emulate the folding process of RNA folding by aiming to keep locally stable substructures
 - *Energy density seen by a basepair*: the free energy of the “optimal substructure” normalized by distance
 - *Energy density of an unpaired base*: energy density of the nearest encapsulating basepair
 - Densityfold optimizes a linear combination of free energy and total energy density
 - For every potential basepair, compute the optimal contribution of the implied substructure
 - The optimization function is non linear
Hill climbing process for approximating the contributions of unpaired bases
-

Densityfold energy types

- $eH(i,j,.)$: free energy of a hairpin loop enclosed by the base pair $S[i].S[j]$
 - $eS(i,j,.)$: free energy of the base pair $S[i].S[j]$ provided that it forms a stacking pair with $S[i+1].S[j-1]$
 - $eBI(i,j,i',j')$: free energy of an internal loop or a bulge that starts with $S[i].S[j]$ and ends with $S[i'].S[j']$
-

Densityfold energy types

- $eM(i,j,i_1,j_1,\dots,i_k,j_k)$: free energy of multibranch loop that starts with $S[i].S[j]$ and branches out $S[i_1].S[j_1]$, $S[i_2].S[j_2]$, \dots , $S[i_k].S[j_k]$
- $eDA(j,j-1)$: free energy of an unpaired dangling base $S[j]$ when $S[j-1]$ forms a base pair with another base

Densityfold energy tables

- $ED(j)$: minimum total free energy density of a secondary structure for substring $S[1, j]$.
- $E(j)$: free energy of the energy density minimized secondary structure for substring $S[1, j]$.
- $ED_S(i, j)$: minimum total free energy density of a secondary structure for $S[i, j]$, provided that $S[i].S[j]$ is a base pair.
- $E_S(i, j)$: free energy of the energy density minimized secondary structure for the substring $S[i, j]$, provided that $S[i].S[j]$ is a base pair.

Densityfold energy tables

- $ED_{BI}(i, j)$: minimum total free energy density of a secondary structure for $S[i, j]$, provided that there is a bulge or an internal loop starting with base pair $S[i].S[j]$.
- $E_{BI}(i, j)$: free energy of an energy density minimized structure for $S[i, j]$, provided that a bulge or an internal loop starting with base pair $S[i].S[j]$.
- $ED_M(i, j)$: minimum total free energy density of a secondary structure for $S[i, j]$, such that there is a multibranch loop starting with base pair $S[i].S[j]$.
- $E_M(i, j)$: free energy of an energy density minimized structure for $S[i, j]$, provided there is a multibranch loop starting with base pair $S[i].S[j]$.

Calculating energy tables

$$ED(j) = \min \left\{ \begin{array}{l} ED(j-1) \\ \min_{1 \leq i \leq j-1} \{ED(i-1) + ED_S(i, j)\} \end{array} \right\}$$

$$ED_S(i, j) = \min \left\{ \begin{array}{l} +\infty, \\ eH(i, j), \\ 2 \frac{eS(i, j) + E_S(i+1, j-1)}{j-i+1} + ED_S(i+1, j-1), \\ ED_{BI}(i, j), \\ ED_M(i, j) \end{array} \right. \begin{array}{l} (i) \\ (ii) \\ (iii) \\ (iv) \\ (v) \end{array} \right\}$$

- Similar calculations for other tables
- $O(n^{k+2})$ time and $O(n^2)$ space

Linear combination of MFE and ED

For any $x \in \{S, BI, M\}$ let $ELC_x(i, j) = ED_x(i, j) + E_x(i, j)$.

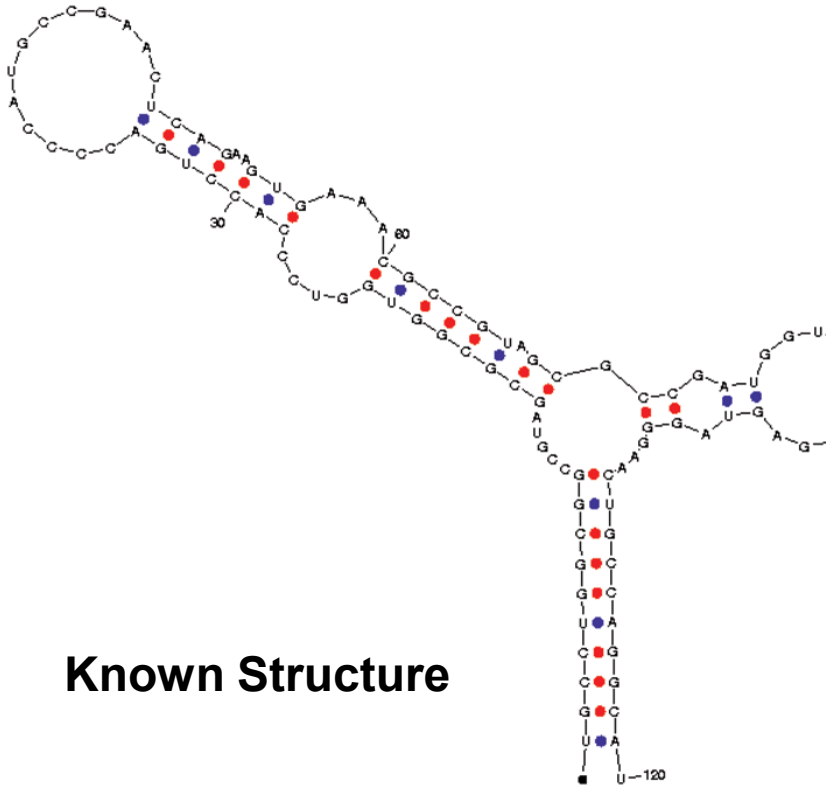
Optimize $ELC(n) = ED(n) + E(n)$.

$$ELC(j) = \min \left\{ \begin{array}{l} ELC(j-1) \\ \min_{1 \leq i \leq j-1} \{ELC(i-1) + ELC_S(i, j)\} \end{array} \right\}$$

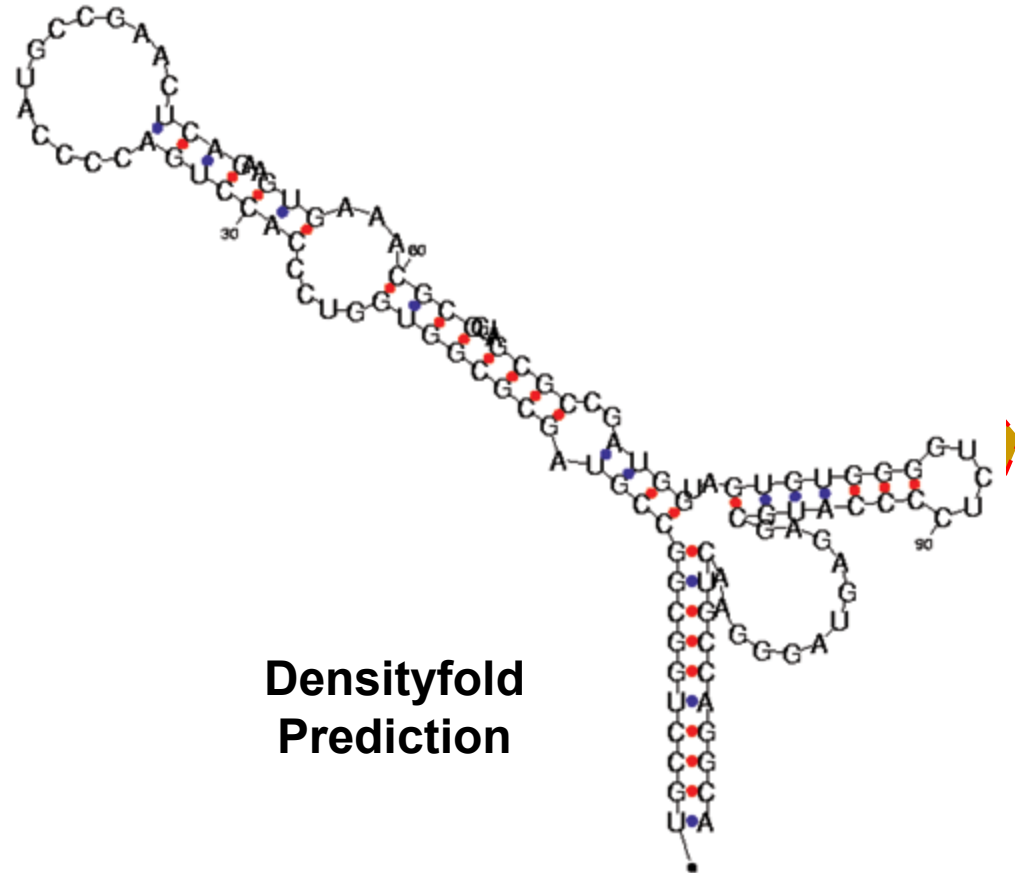
$$ELC_S(i, j) = \min \left\{ \begin{array}{ll} +\infty, & (i) \\ eH(i, j) \cdot (1 + \sigma), & (ii) \\ 2 \frac{eS(i, j) + E_S(i+1, j-1)}{j-i+1} + ELC_S(i+1, j-1) + \sigma \cdot eS(i, j), & (iii) \\ ELC_{BI}(i, j), & (iv) \\ ELC_M(i, j) & (v) \end{array} \right\}$$

- Similar formulations for ELC_{BI} and ELC_M
- $O(n^4)$ running time

Densityfold prediction: E.coli 5S rRNA



Known Structure



Densityfold Prediction

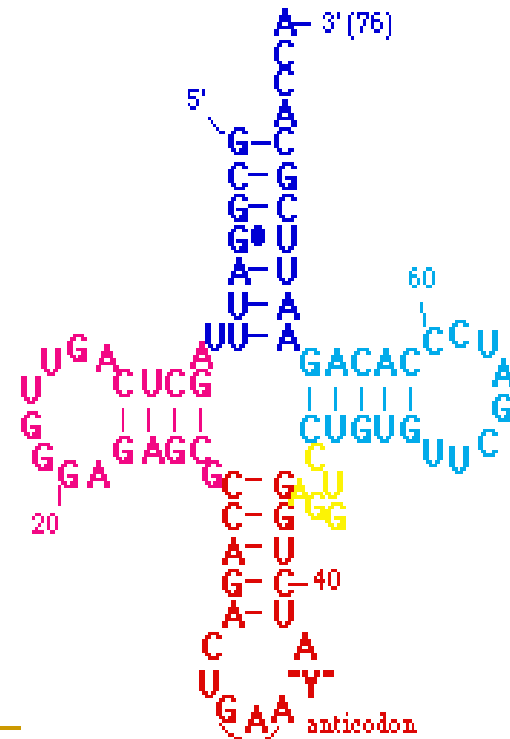
CONTRAFOLD

CONTRAFold

Probabilistic RNA folding algorithm

Problem: Given an RNA sequence, predict the most likely secondary structure

AUCCCGUAUCGAUC
AAAUCCAUGGGUAC
CCUAGUGAAAGUGUA
UAUACGUGCUCUGAU
UCUUUACUGAGGAGU
CAGUGAACGAACUGA

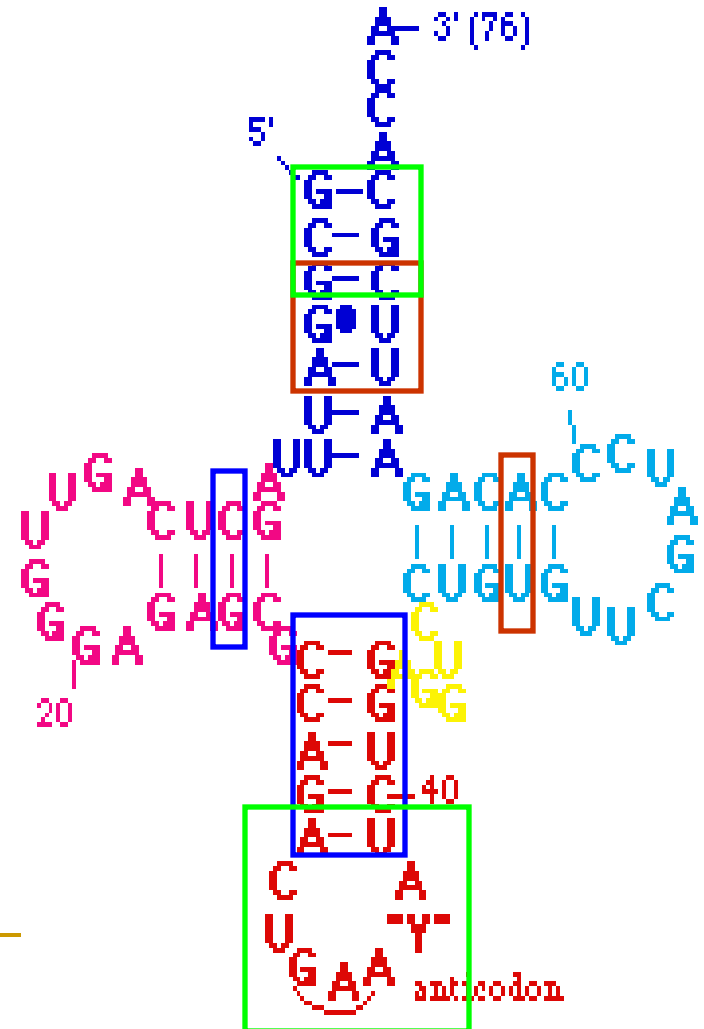


CONTRAFold

- CONTRAFold looks at *features* that indicate a good structure

For example:

- C-G base pairings
- A-U base pairings
- Helices of length 5
- Hairpin loops of size 9
- Bulge loops of size 2
- CG/GC Base-pair stacking interactions



Choosing a structure

- Every *feature* f_i is associated with a *weight* w_i .
- The probability of a structure y , given a sequence x , is determined by the following relationship:

of occurrences of feature i , in structure y generated from sequence x

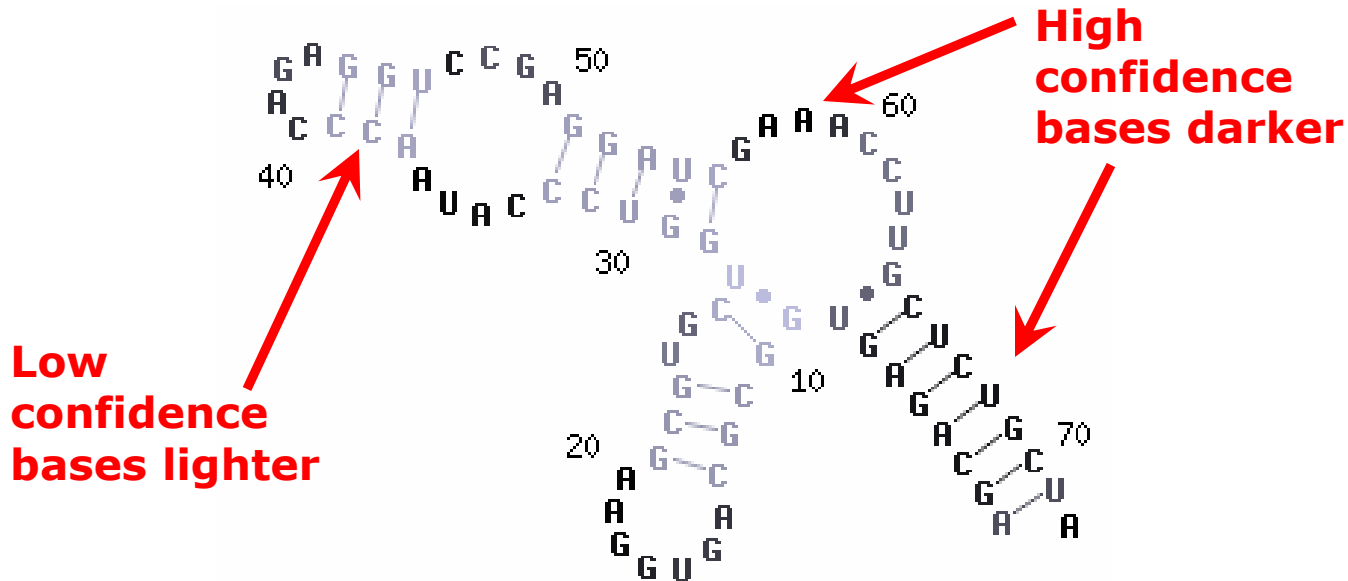
$$P(\mathbf{y} \mid \mathbf{x}) \propto \exp\left(\sum_i w_i \cdot F_i(\mathbf{x}, \mathbf{y})\right)$$

structure sequence weight of Feature i

$$F(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} F_1(\mathbf{x}, \mathbf{y}) \\ F_2(\mathbf{x}, \mathbf{y}) \\ \vdots \end{bmatrix} = \begin{bmatrix} \# \text{ of hairpin loops in } y \text{ of length } 4 \\ \# \text{ of GC/GC stacking interactions in } y \\ \vdots \end{bmatrix}.$$

Choosing a structure

- Considers all structures and finds optimal structure via *dynamic programming* in $O(n^3)$
- Added bonus: probability associated with each base



Maximum Expected Accuracy

For a candidate structure $\hat{\mathbf{y}}$ with true structure \mathbf{y}

$$\hat{\mathbf{y}}_{\text{mea}} = \operatorname{argmax}_{\hat{\mathbf{y}}} E_{\mathbf{y}} [\text{accuracy}_{\gamma}(\hat{\mathbf{y}}, \mathbf{y})]$$

$$M_{1,L} = \max_{\mathbf{y}} E_{\mathbf{y}} [\text{accuracy}_{\gamma}(\hat{\mathbf{y}}_{\text{mea}}, \mathbf{y})]$$

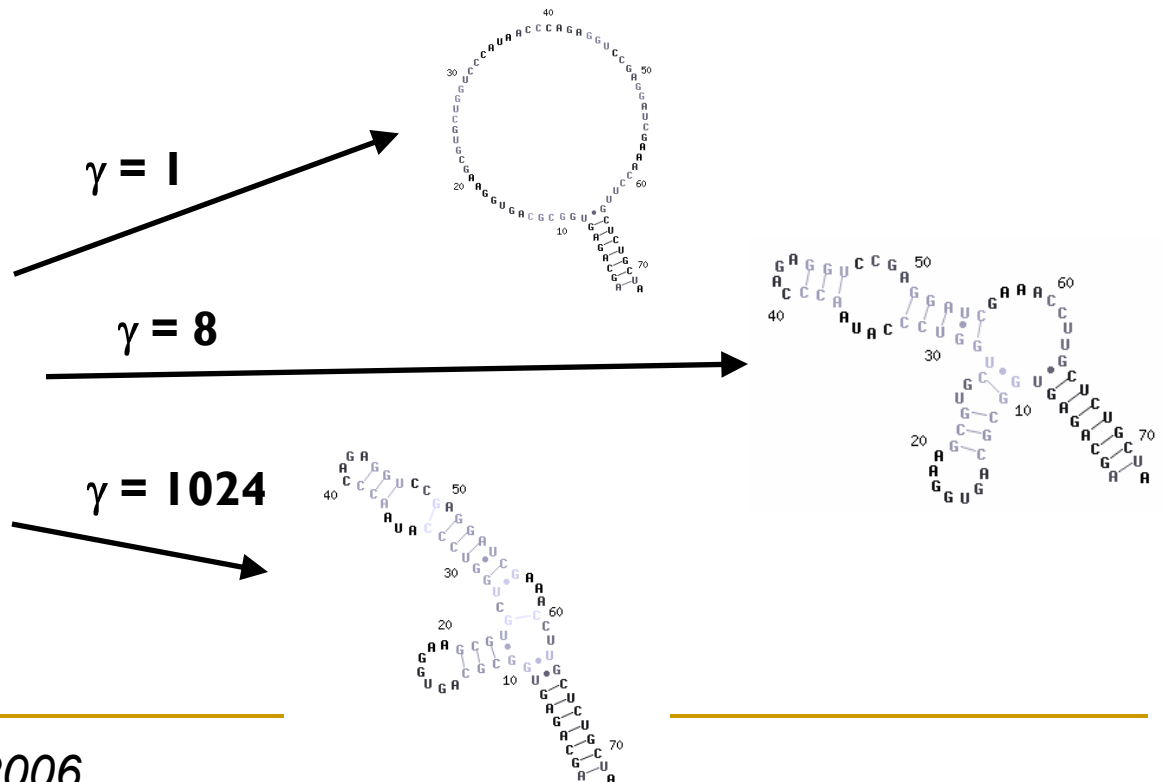
$$M_{i,j} = \max \left\{ \begin{array}{ll} q_i & \text{if } i=j \\ q_i + M_{i+1,j} & \text{if } i < j \\ q_j + M_{i,j-1} & \text{if } i < j \\ \gamma \cdot 2p_{ij} + M_{i+1,j+1} & \text{if } i+2 < j \\ M_{i,k} + M_{k+1,j} & \text{if } i \leq k < j \end{array} \right.$$

Sensitivity vs Specificity: γ

$$\text{Sensitivity} = \frac{\# \text{ correct base pairings}}{\# \text{ true base pairings}}$$

$$\text{Specificity} = \frac{\# \text{ correct base pairings}}{\# \text{ predicted base pairings}}$$

AUCCCGUAUCGAUC
AAAUAUCCAUGGGUAC
CCUAGUGAAAGUGUA
UAUACGUGCUCUGAU
UCUUUACUGAGGAGU
CAGUGAACGAACUGA

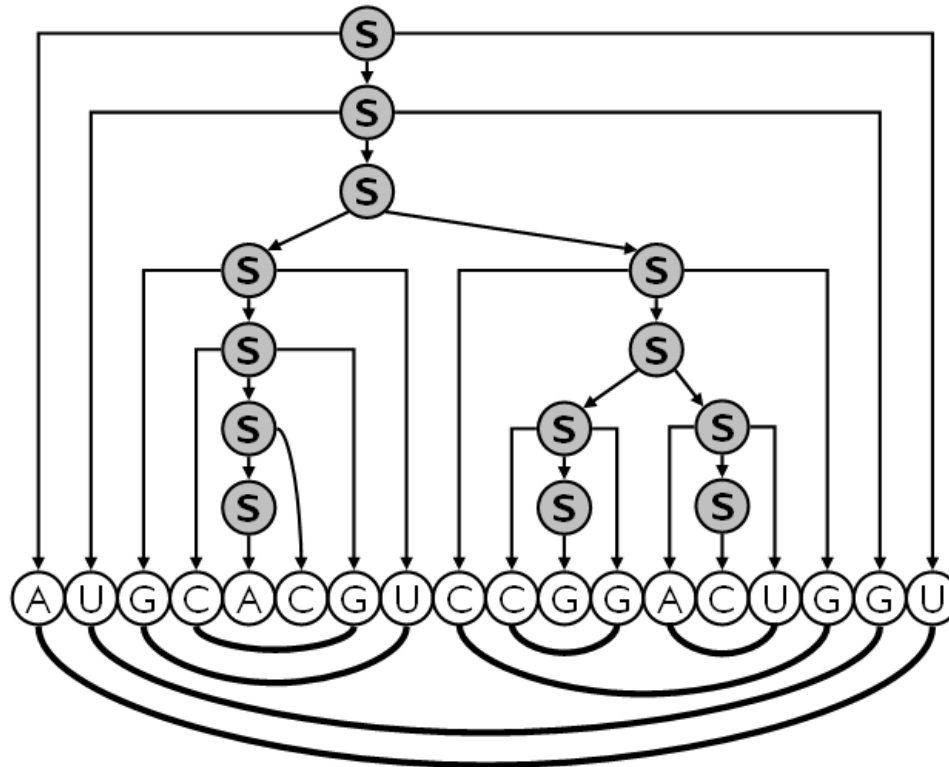


Learning to predict good structures

- CONTRAfold learns the relative value, or *weight*, of each of its *features*
- A *training set* is a collection of known correct solutions that a program learns from.
- CONTRAfold trains on set of published examples of known RNA structures taken from a database called Rfam (RNA families)
- CONTRAfold determines the *weight* for each *feature* that maximizes its performance on the training set.

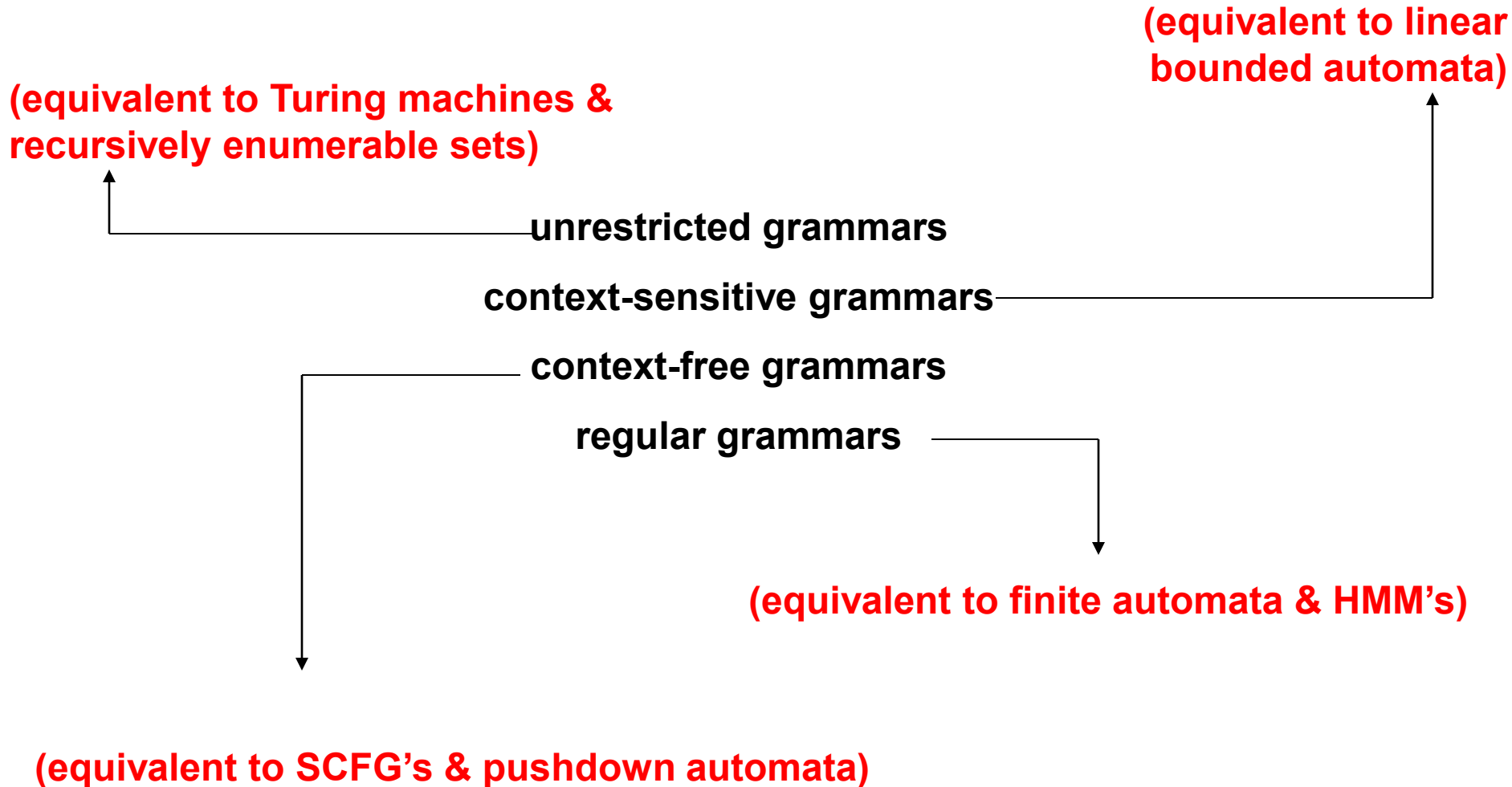
STOCHASTIC CONTEXT-FREE GRAMMARS

SCFG



- RNA folding can be represented as context-free grammars

Chomsky hierarchy



Context-free grammars

A *context-free grammar* is a generative model denoted by a 4-tuple:

$$G = (V, \alpha, S, R)$$

where:

α is a *terminal alphabet*, (e.g., $\{a, c, g, t\}$)

V is a *nonterminal alphabet*, (e.g., $\{A, B, C, D, E, \dots\}$)

$S \in V$ is a special *start symbol*, and

R is a set of rewriting rules called *productions*.

Productions in R are rules of the form:

$$X \rightarrow \lambda$$

where $X \in V$, $\lambda \in (V \cup \alpha)^*$

Context “freeness”

The “*context-freeness*” is imposed by the requirement that the l.h.s of each production rule may contain only a **single** symbol, and that symbol must be a **nonterminal**:

$$X \rightarrow \lambda$$

Thus, a CFG **cannot** specify *context-sensitive* rules such as:

$$wXz \rightarrow w\lambda z$$

Derivations

Suppose a CFG G has generated a *terminal string* $x \in \alpha^*$. A *derivation* $S \Rightarrow^* x$ denotes a possible for generating x .

A *derivation* (or *parse*) consists of a series of applications of productions from R , beginning with the *start symbol* S and ending with the *terminal string* x :

$$S \Rightarrow s_1 \Rightarrow s_2 \Rightarrow s_3 \Rightarrow \cdots \Rightarrow x$$

where $s_i \in (V \cup \alpha)^*$.

We'll concentrate on leftmost derivations where the leftmost nonterminal is always replaced first.

Context-free vs. regular

The advantage of CFG's over HMM's lies in their ability to model arbitrary runs of matching pairs of elements, such as matching pairs of parentheses:

$$\dots(((((((\dots)))))))))\dots$$

When the number of matching pairs is unbounded, a finite-state model such as a DFA or an HMM is inadequate to enforce the constraint that all left elements must have a matching right element.

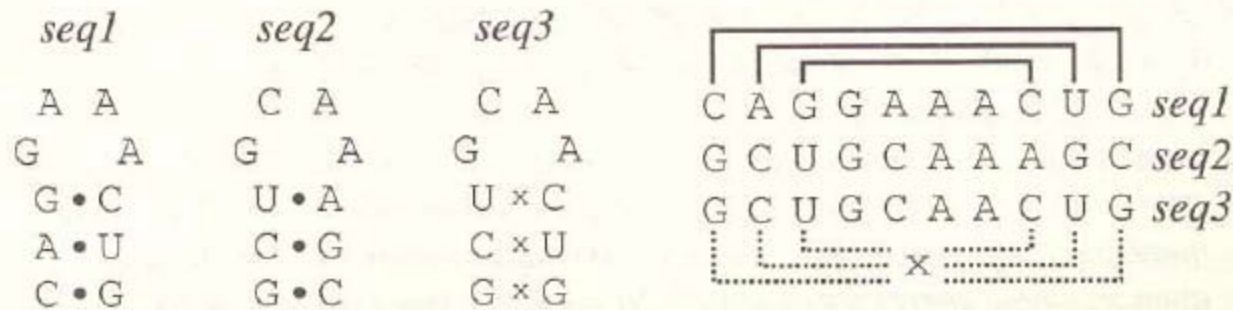
In contrast, in a CFG we can use rules such as $X \rightarrow (X)$. A sample derivation using such a rule is:

$$X \Rightarrow (X) \Rightarrow ((X)) \Rightarrow (((X))) \Rightarrow ((((X))))$$

An additional rule such as $X \rightarrow \epsilon$ is necessary to terminate the recursion.

A CFG for an RNA

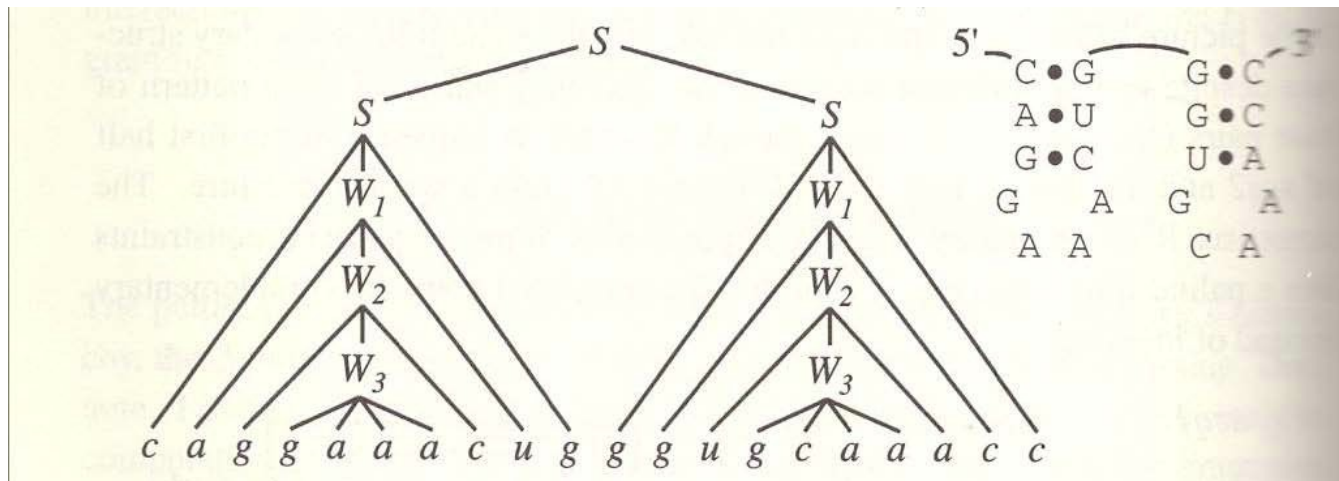
- RNA hairpin with 3 bp stem and a 4-base loop (GAAA or GCAA)



S-> aXu | cXg | gXc | uXa
 X-> aYu | cYg | gYc | uYa
 Y-> aZu | cZg | gZc | uZa
 Z->gaaa | gcaa

Parse trees

- A representation of a parse of a string by a CFG
- **Root** – start nonterminal S
- **Leaves** – terminal symbols in the given string
- **Internal nodes** - nonterminals
- The children of an internal node are the productions of that nonterminal (left-to-right order)



Stochastic CFG

A *stochastic context-free grammar* (SCFG) is a CFG plus a probability distribution on productions:

$$G = (V, \alpha, S, R, P_p)$$

where $P_p : R \rightarrow [0, 1]$, and probabilities are normalized at the level of each l.h.s. symbol X :

$$\forall X \in V \left[\sum_{X \rightarrow \lambda} P_p(X \rightarrow \lambda) = 1 \right]$$

Thus, we can compute the probability of a single derivation $S \Rightarrow^* x$ by multiplying the probabilities for all productions used in the derivation:

$$\prod_i P(X_i \rightarrow \lambda_i)$$

We can sum over all possible (leftmost) derivations of a given string x to get the probability that G will generate x at random:

$$P(x | G) = \sum_j P(S \Rightarrow_j^* x | G).$$

An example

As an example, consider $G=(V_G, \alpha, S, R_G, P_G)$, for $V_G=\{S, L, N\}$, $\alpha=\{a, c, g, t\}$, and R_G the set consisting of:

$$S \rightarrow a S t \mid t S a \mid c S g \mid g S c \mid L \quad (\mathbf{P=0.2})$$

$$L \rightarrow N N N N \quad (\mathbf{P=1.0})$$

$$N \rightarrow a \mid c \mid g \mid t \quad (\mathbf{P=0.25})$$

Then the probability of the sequence $acgtacgtacgt$ is given by:

$$P(acgtacgtacgt) =$$

$$P(S \Rightarrow aSt \Rightarrow acSgt \Rightarrow acgScgt \Rightarrow acgtSacgt \Rightarrow \\ acgtLacgt \Rightarrow acgtNNNNacgt \Rightarrow acgtaNNNacgt \Rightarrow \\ acgtacNNacgt \Rightarrow acgtacgNacgt \Rightarrow acgtacgtacgt) =$$

$$0.2 \quad 0.2 \quad 0.2 \quad 0.2 \quad 0.2 \quad 1 \quad 0.25 \quad 0.25 \quad 0.25 \quad 0.25 = 1.25 \cdot 10^{-6}$$

because this sequence has only one possible (leftmost) derivation under grammar G .

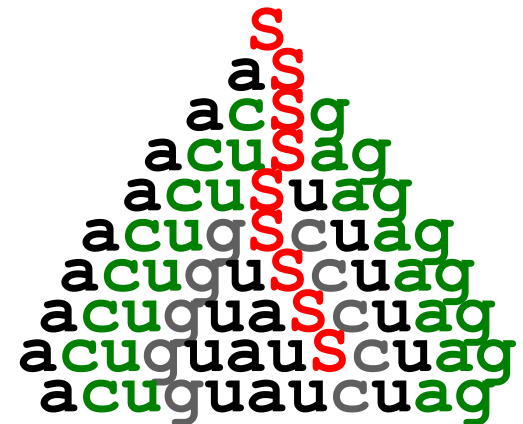
Structure using SFCG

- Grammar rules with associated probabilities

$S \rightarrow$	aSu	cSg	aS	uS	\dots	Su	SS	ϵ
P	.21	.15	.11	.08		.03	.22	.02

- We select the set of transformations that highest probability of generating the input sequence. This set gives us our structure.
- Let's generate a structure for the sequence acuuauuag

acuuauuag
- (((()))))



Chomsky Normal Form

A CNF grammar is one in which all productions are of the form:

$$X \rightarrow YZ$$

or:

$$X \rightarrow a$$

Non-CNF:

$$S \rightarrow aSt | tSa | cSg | gSc | L$$

$$L \rightarrow NNNN$$

$$N \rightarrow a | c | g | u$$

CNF:

$$S \rightarrow AS_T | TS_A | CS_G | GS_C | NL_1$$

$$S_A \rightarrow SA$$

$$S_T \rightarrow ST$$

$$S_C \rightarrow SC$$

$$S_G \rightarrow SG$$

$$L_1 \rightarrow NL_2$$

$$L_2 \rightarrow NN$$

$$N \rightarrow a | c | g | u$$

$$A \rightarrow a$$

$$C \rightarrow c$$

$$G \rightarrow g$$

$$T \rightarrow u$$

Parsing CFG

Two questions for a CFG:

- 1) Can a grammar G derive string x ?
- 2) If so, what series of productions would be used during the derivation? (*there may be multiple answers!*)

Additional questions for an SCFG:

- 1) What is the *probability* that G derives string x ?
- 2) What is the *most probable* derivation of x via G ?

Parsing CFG

- CYK Algorithm (Cocke-Younger-Kasami)
 - Dynamic Programming method
 - Modified CYK for SCFG
 - “Inside algorithm”
 - Training similar to HMM
 - If parses are known for training data sequences, simply count the number of times for each production, calculate probabilities (labeled sequence training for HMM)
 - If parses are not known, apply an EM algorithm called “Inside-Outside” (“forward-backward” for HMM)
-