

More on  
Stochastic Context Free Grammars  
for RNA Analysis

BMI/CS 776

[www.biostat.wisc.edu/bmi776/](http://www.biostat.wisc.edu/bmi776/)

Colin Dewey

[cdewey@biostat.wisc.edu](mailto:cdewey@biostat.wisc.edu)

Spring 2012

# Goals for Lecture

the key concepts to understand are the following

- approaches to determining the structure of an SCFG grammar
- the task of searching for sequences that “match” a given RNA structure

# Where do we get a grammar?

1. from a canonical secondary structure
2. through an iterative, refinement process
3. alternatively, use a simple, generic one

# 1. Deriving a Grammar from a Secondary Structure

given a structure



can construct a simple grammar characterizing it

$$s \rightarrow C s_1 G$$

$$s_1 \rightarrow A s_2 U$$

$$s_2 \rightarrow b_1 b_2 b_3$$

$$b_1 \rightarrow U$$

$$b_2 \rightarrow U$$

$$b_3 \rightarrow C$$

can add productions to allow for variation

$$s \rightarrow U s_1 A$$

$$s \rightarrow A s_1 U$$

$$s \rightarrow G s_1 C$$

base pair substitutions

$$s_1 \rightarrow s_1 A$$

insertions

$$b_2 \rightarrow A$$

$$b_2 \rightarrow C$$

$$b_2 \rightarrow G$$

single base substitutions

## 2. Deriving a Grammar Through an Iterative Process

- consider the approach used by Eddy & Durbin to learn an SCFG model of tRNAs

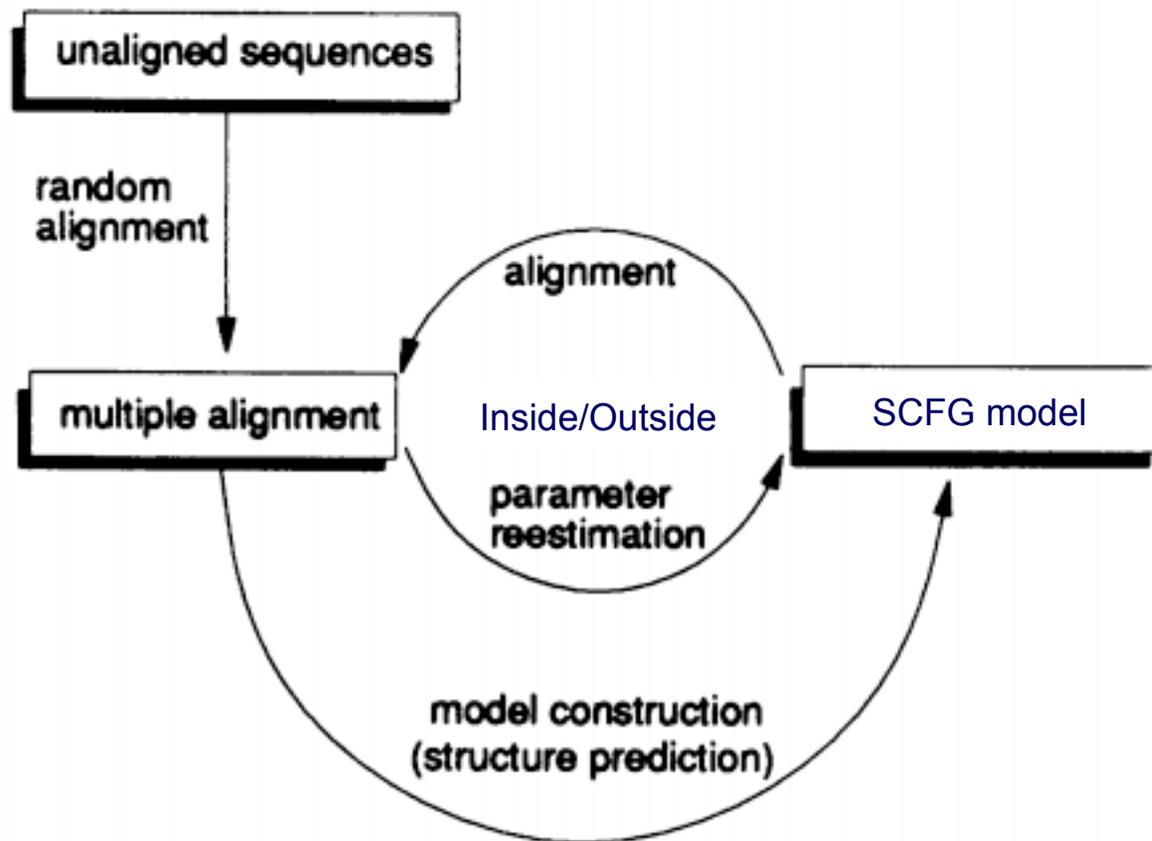
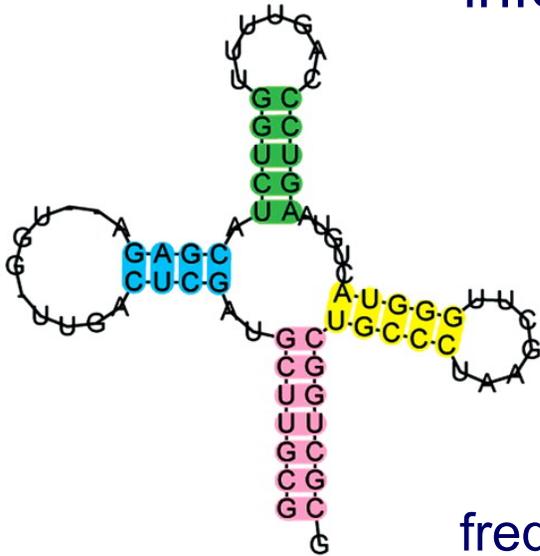


Figure from Eddy & Durbin, *Nucleic Acids Research*, 22(11):2079-2088, 1994.

# Eddy & Durbin: Model Construction Step

- given multiple alignment, compute mutual information between pairs of positions

A



$$M(i, j) = \sum_{x_i, x_j} f_{x_i x_j} \log_2 \frac{f_{x_i x_j}}{f_{x_i} f_{x_j}}$$

frequency of  $x_i$  in column  $i$

joint frequency of  $x_i$  in column  $i$  and  $x_j$  in column  $j$

B

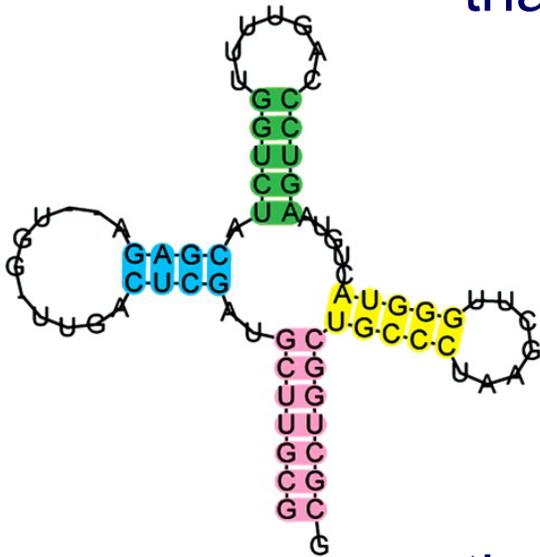
X51423.1/1802-1874	GGGCUUGUAGCUCAGCU--GGU--AGAGCGCCGCCUUUGCAAAGCGGAGGCCUGGGUCCGAAUCCCAGCAAGUCCA
AF008220.1/6713-6785	GCGGUUGUGGCGAAGU--GGUU-AACGCAACAGAUUGUGGCUCUGGCAUUCGUGGGUUCGAUUCCAUCAUUCGCC
X54408.1/1-73	GCCCCAUCGUCUAGA--GGCCUAGGACACUCCCUUUCACGGAGAAAA-CGCGGAUUCGAAUCCCGUGGGGUA
J05395.1/2325-2252	GGUUUCGUGUCUAGUC--GGUU-AUGGCAUCUGUUAACACGCAGAACGUCCCCAGUUCGAUCCUGGGCGAAAUCG
X12977.1/475-548	CGUGAUUAGCGCAGCCCGGU--AGCGCAUCUGGUUUGGGACAGAGGGUCAAAGGUUCGAUCCUUUAUCACCGA
X03126.1/364-434	AAGAGUAUAGUUUAAA--GGU--AAAACAGAAAGCUUCAACCUUUAUUU-UCUUAGUUCGAGUCUAAGUGCUCUUG
X52070.1/777-849	UCCUCCGUAGCUCAAUU--GGC--AGAGCAGCCGGCUGUUAACCGGAGGUUACUGGUUCGAGUCCAGUCGGGGGAG
X13558.1/186-115	UGGGCGUAGCCAAGC--GGU--AAGGCAAACGGUUUUUGGUCCCGCUAUUCGGAGGUUCGAUCCUUCCGUCCCAG
X52392.1/3967-4040	GCUAGCGUGGCAGAGCUCGGCA-AAUGCAAAGGCUUAAGCCUUUAUC-CAGAGGUUCAAAUCCUCUCCUAGCU
K00287.1/1-73	UCCUUGUUAAGCUCAGUU--GGU--AGAGCGUUCGGCUUUUAACCGAAUUGUCAGGGGUUCGAGCCCCUAUGAGGAG
Consensus (-31.49)	(((((.....))))(((((.....)))).....((((.....))))))))))

Figure from Voß, *Nucleic Acids Research*, 34:5471-5481, 2006.

# Eddy & Durbin: Model Construction Step

- use a DP like Nussinov to find folded structure that maximizes mutual information

A



$$\gamma(i, j) = \max \begin{cases} \gamma(i + 1, j) \\ \gamma(i, j - 1) \\ \gamma(i + 1, j - 1) + M(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k + 1, j)] \end{cases}$$

- then derive grammar from this structure

B

X51423.1/1802-1874	GGGCUUGUAGCUCAGCU--GGU--AGAGCGCCGCCUUUGCAAAGCGGAGGCCUGGGUCCGAAUCCCAGCAAGUCCA
AF008220.1/6713-6785	GCGGUUGUGGCGAAGU--GGUU-AACGCAACAGAUUGUGGCUCUGGCAUUCGUGGGUUCGAUUCCAUCAUUCGCC
X54408.1/1-73	GCCCCAUCGUCUAGA--GGCCUAGGACACUCUCCUUUCACGGAGAAA-CGCGGAUUCGAAUCCGCUGGGGUA
J05395.1/2325-2252	GGUUUCGUGUCUAGUC--GGUU-AUGGCAUCUGUUUAACACGCAGAACGUCCCCAGUUCGAUCCUGGGCGAAAUCG
X12977.1/475-548	CGUGAUUAGCGCAGCCCGU--AGCGCAUCUGUUUGGGACCAGAGGGUCAAGGUUCGAUCCUUUAUCACCGA
X03126.1/364-434	AAGAGUAUAGUUUAAA--GGU--AAAACAGAAAGCUUCAACCUUUAUU-UCUUAGUUCGAGUCUAAGUGCUCUUG
X52070.1/777-849	UCCUCCGUAGCUCAAUU--GGC--AGAGCAGCCGGCUGUUAACCCGGCAGGUUACUGGUUCGAGUCCAGUCGGGGGAG
X13558.1/186-115	UGGGCGUA GCCAAGC--GGU--AAGGCAA CCGGUUUUGGUCCGCUAUUCGGAGGUUCGAUCCUUCCGUCCCAG
X52392.1/3967-4040	GCUAGCGUGGCAGAGCUCGGCA-AAUGCAAAGGCUUAAGCCUUUAUC-CAGAGGUUCAAAUCCUCUCCUAGCU
K00287.1/1-73	UCCUUGUUA GCUCAGUU--GGU--AGAGCGUUCGGCUUUUAACCGAAUUGUCAGGGGUUCGAGCCCCUAUGAGGAG
Consensus (-31.49)	(((((.....))))).(((.....))).....((((.....)))))).....

Figure from Voß, *Nucleic Acids Research*, 34:5471-5481, 2006.

# 3. Using a Simple, Generic Grammar

- a grammar that could characterize almost any RNA structure

$$s \rightarrow C s G \mid G s C \mid A s U \mid U s A$$

$$s \rightarrow C s \mid G s \mid A s \mid U s$$

$$s \rightarrow s G \mid s C \mid s U \mid s A$$

$$s \rightarrow G \mid C \mid U \mid A$$

$$s \rightarrow s s$$

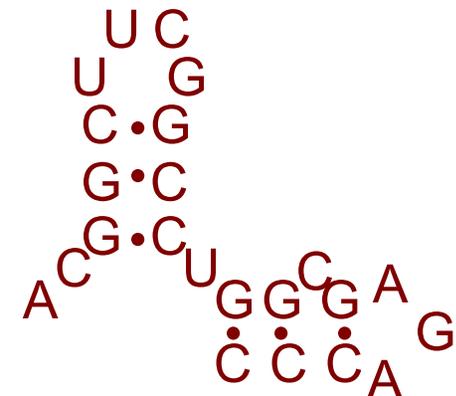
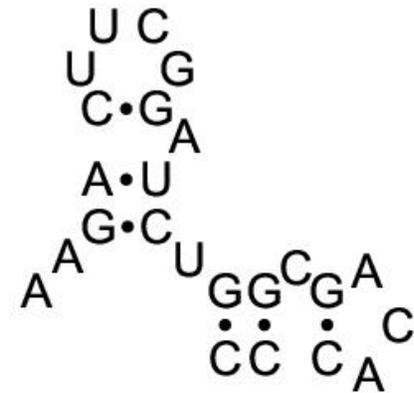
# Searching Sequence for a Secondary Structure

Given

- a single RNA sequence with its secondary structure
- another RNA query sequence

ACGGCUUCGGCCUGGCGAGACCC

Determine if the query sequence has the “same” secondary structure





# Deriving a Grammar from Secondary Structure

given a structure



can construct a simple grammar characterizing it

$$s \rightarrow C s_1 G$$

$$s_1 \rightarrow A s_2 U$$

$$s_2 \rightarrow b_1 l_1$$

$$l_1 \rightarrow b_2 b_3$$

$$b_1 \rightarrow U$$

$$b_2 \rightarrow U$$

$$b_3 \rightarrow C$$

can add productions to allow for variation

$$s \rightarrow U s_1 A$$

$$s \rightarrow A s_1 U$$

$$s \rightarrow G s_1 C$$

base pair substitutions

$$s_1 \rightarrow s_1 A$$

insertions

$$b_2 \rightarrow A$$

$$b_2 \rightarrow C$$

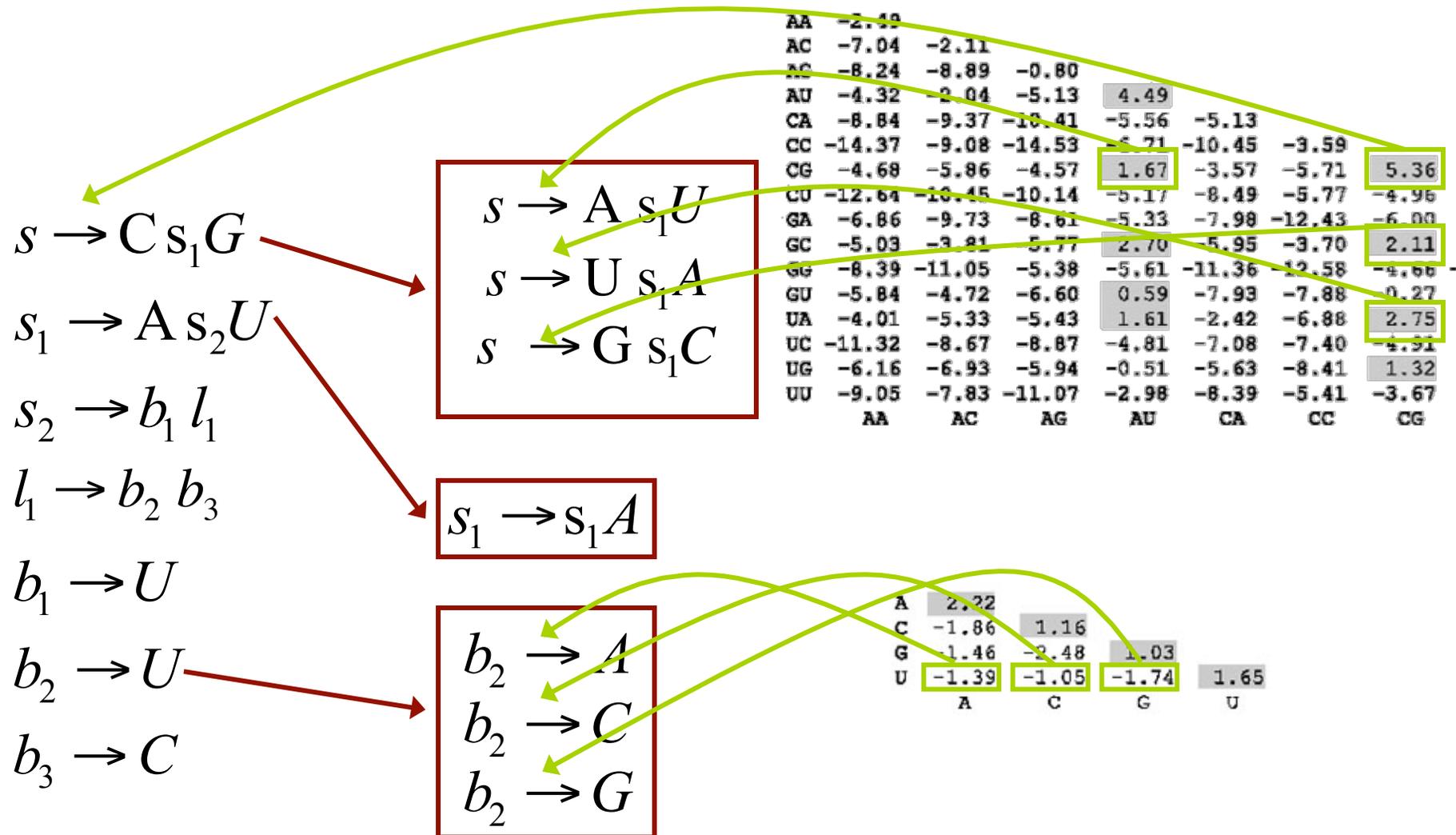
$$b_2 \rightarrow G$$

single base substitutions



# Setting the Parameters in the Grammar

- Infer parameters from the RIBOSUM matrices (taking into account the latter are log-odds scores)

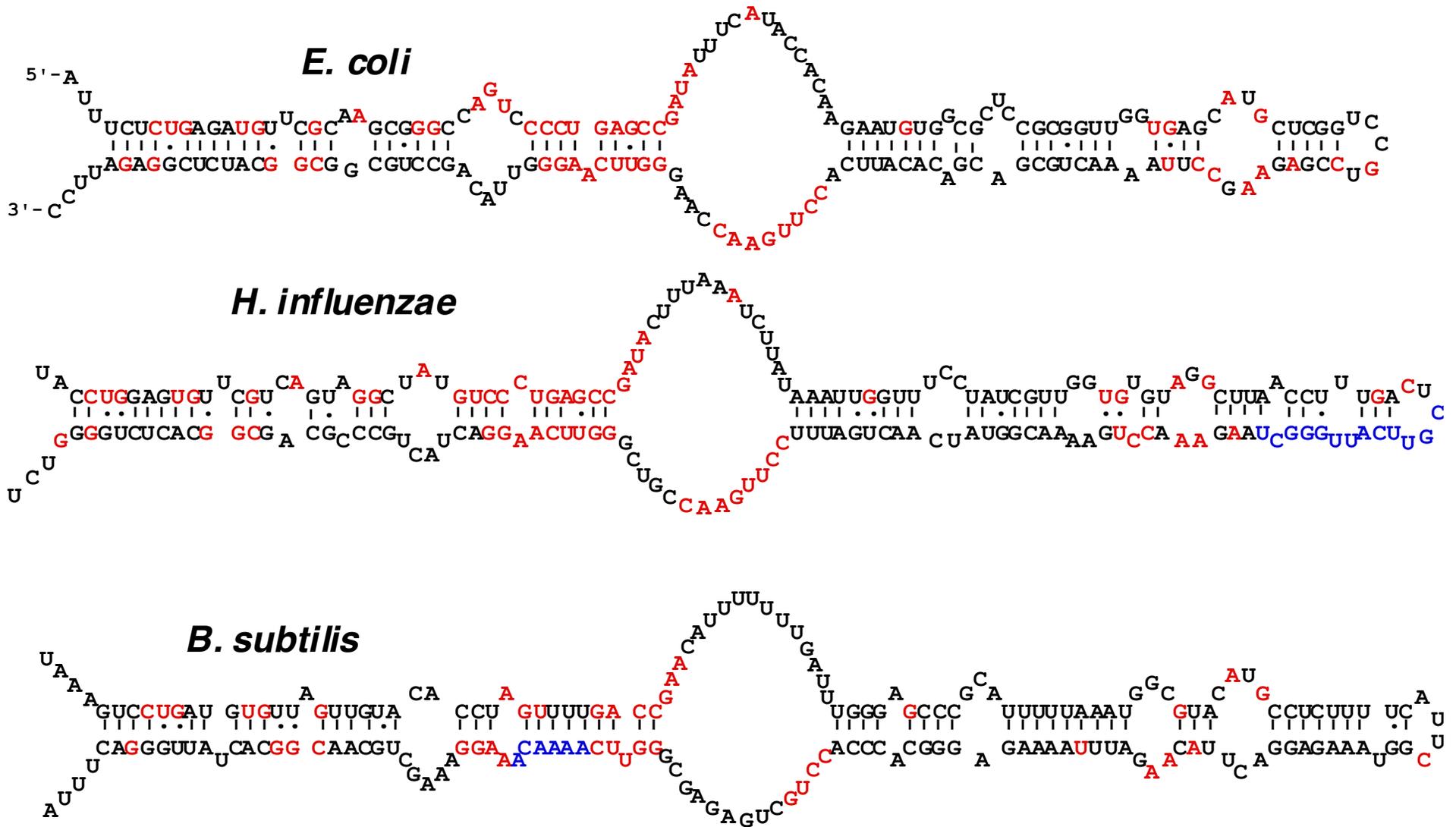


# RSEARCH: Searching Sequence for a Secondary Structure

[Klein & Eddy, *BMC Bioinformatics* 2003]

- the RSEARCH algorithm implements this idea
- but uses a somewhat different SCFG formulation – covariance models (see section 10.3 in Durbin et al.)

# 6S RNA Secondary Structure



# An RSEARCH Case Study

- finding 6S genes in bacterial genomes
  - E. coli 6S as the query structure
  - searched 14 other genomes with known 6S genes
    - ~ 5,000 intergenic sequences on average
  - the top-scoring RSEARCH hit in all 14 genomes was the known 6S gene

# Summary of RNA Analysis Tasks

- given a sequence, predict its secondary structure
- given a set of related RNA sequences, construct a model of the set
  - parameter learning (Inside-Outside)
  - structure refinement
- given a model of an RNA class, find sequences that belong to the class (Inside or CYK)
- given a sequence/structure, find other sequences with similar structure
- others not discussed