Alt proof: $\{a^n b^n c^n \mid n \geq 0\}$ not CFL:

V & y can't have
both a & c

Thus case 1 not c
$u v^2 x y^2 z$  too few c's
Similarly
case 2 not a
... too few a's

Idea



Repeating shaded part i times
gives $uu^i x y^i z \in A$

$$S \Rightarrow^* u R z$$

$$R \Rightarrow^* v R y$$

$$R \Rightarrow^* x$$

$$R \rightarrow RE$$
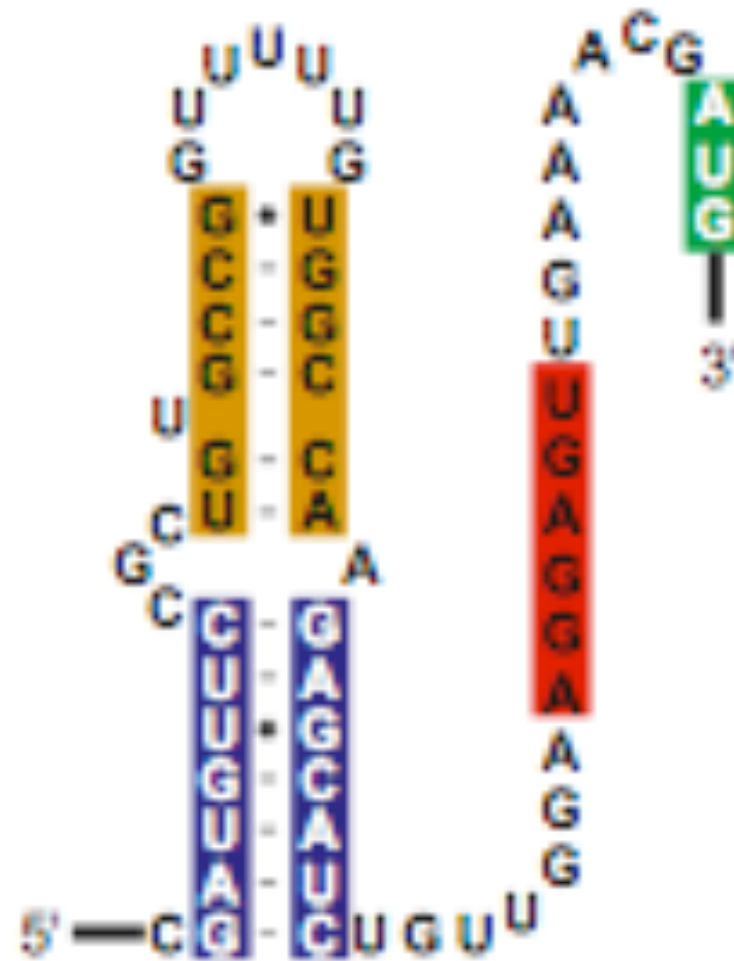$$E \rightarrow \varepsilon$$

$$P = b^{|v|+1}$$

subtlety 1 : tree w/ fewest nodes $vy \neq \varepsilon$
subtlety 2 : pick rep. nearest leaves.
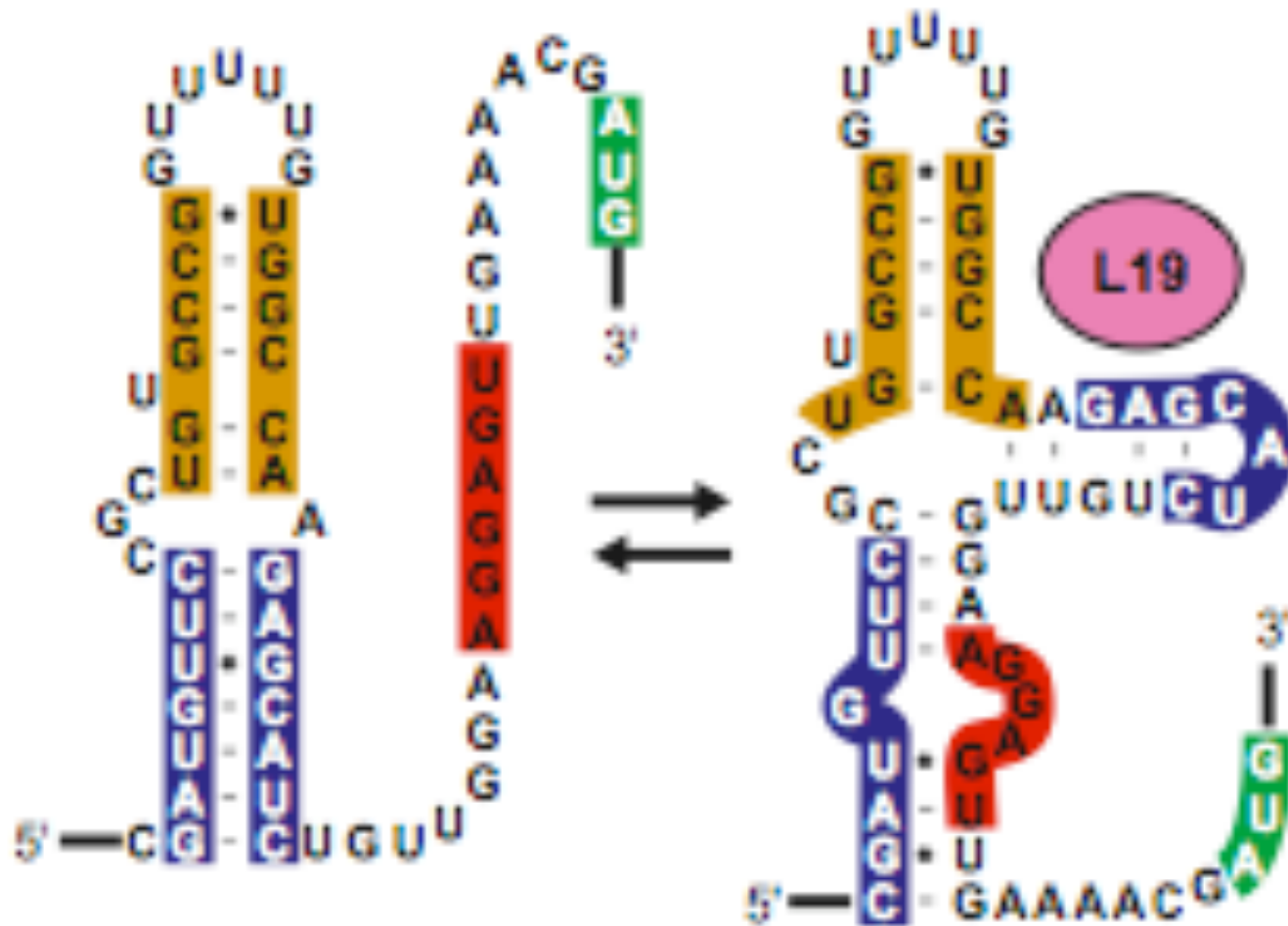
$$|vxy| \leq P \quad 31\text{-}2$$

# And now for something completely different

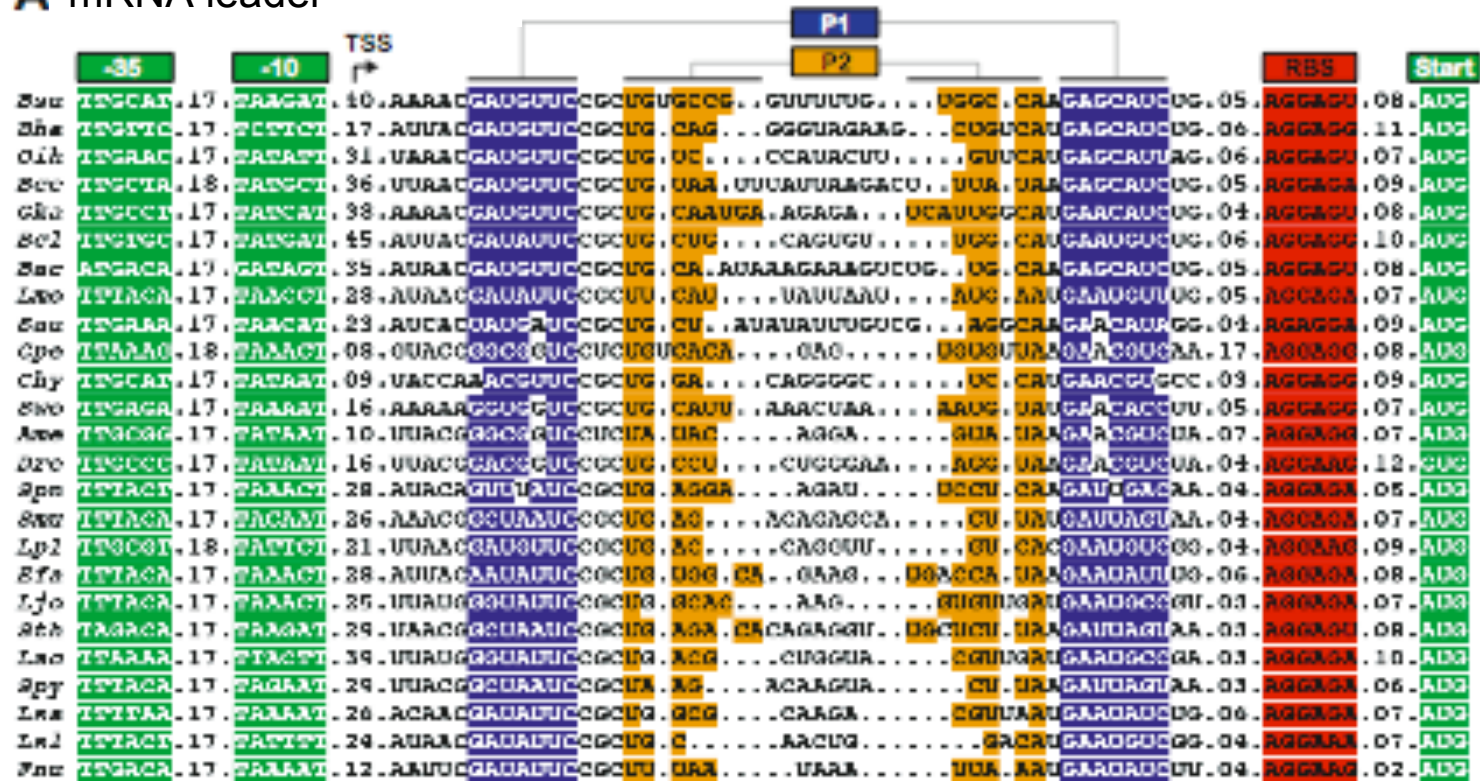CFG utility beyond compilers
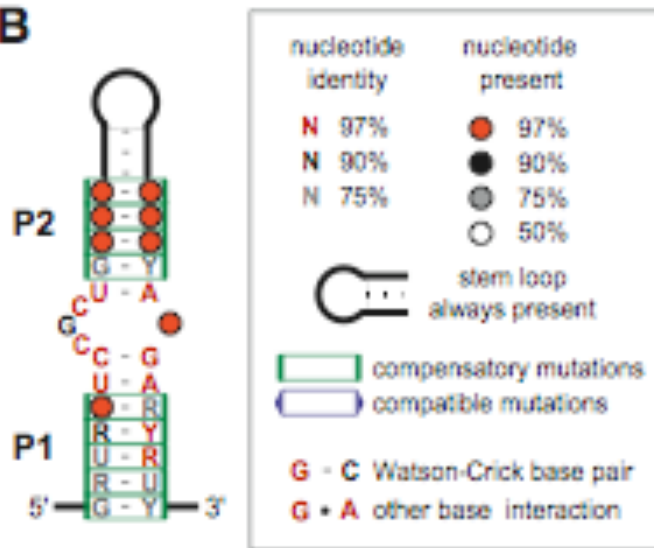
# An RNA Structure

# An RNA Sensor & On/Off Switch



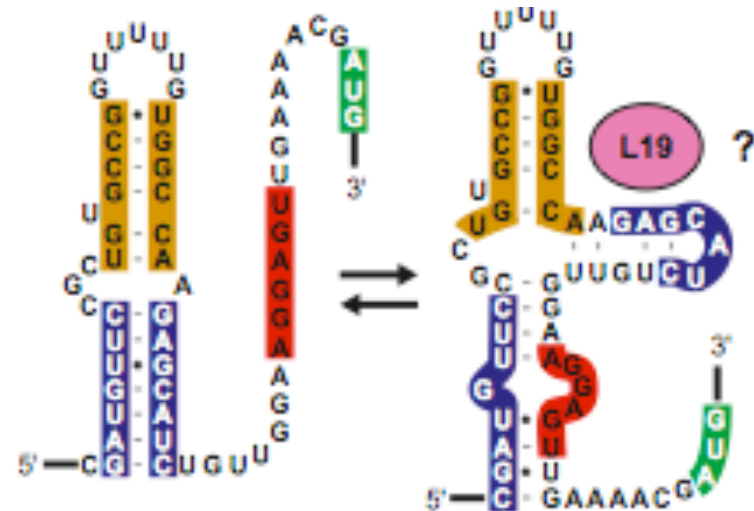L19 absent: Gene On     L19 present: Gene Off

**A** mRNA leader

**B**

**C** mRNA leader switch?

# An RNA Grammar

$S \rightarrow LS \mid L$

$L \rightarrow s \mid \text{``dFd''}$

$F \rightarrow LS \mid \text{``dFd''}$

"dFd" means
Watson-Crick
base pair:

$aFu \mid uFa \mid gFc \mid cFg$
paren-like nesting

a) 
$$S \rightarrow LS \rightarrow LLLLLLLS \rightarrow LLLLLLLL$$
$$\rightarrow ssLssssss \rightarrow ssdFdsssss$$
$$\rightarrow ssdddFdddsssss$$
$$\rightarrow ssdddLSdddsssss$$
$$\rightarrow ssdddLLLLdddsssss$$
$$\rightarrow ssdddssssdddsssss$$

b)
$$s^{\,ss}{}_{\,s}$$
$$d\text{-}d$$
$$d\text{-}d$$
$${}_{ss}d\text{-}d{}_{sssss}$$

c) 
$$F \rightarrow dFd \rightarrow ddFdd \rightarrow ddLSdd$$
$$\rightarrow ddLLdd \rightarrow ddLsdd \rightarrow dddFdsdd$$

# Actually, a _Stochastic_ CFG

Associate probabilities with rules:

$$S \rightarrow LS \quad (0.87) \qquad | \, L \qquad (0.13)$$
$$L \rightarrow S \quad (0.89*\text{p(s)}) \quad | \, dFd \quad (0.11*\text{p(dd)})$$
$$F \rightarrow LS \quad (0.21) \qquad | \, dFd \quad (0.79*\text{p(dd)})$$

Where p(s) & p(dd) are the probabilities of the specific single/paired nucleotides, perhaps from empirical data or a model of sequence evolution

# What SCFG Gives

"Prior" probabilities for

    frequencies of nucleotides/pairs

    fraction paired vs unpaired

    average lengths of each, etc.

Result: a probability distribution on sequences/structures

    E.g., is my sequence more likely to arise under this RNA model or a simple "background" model, say where A/C/G/T = 1/4?

# Cocke-Kasami-Younger Parser

Suppose all rules of form $A \to BC$ or $A \to a$
(by mechanically transforming grammar, or algorithm below…)

Given $x = x_1...x_n$, want $M_{i,j} = \{ A \mid A \to x_{i+1}...x_j \}$

For j=2 to n

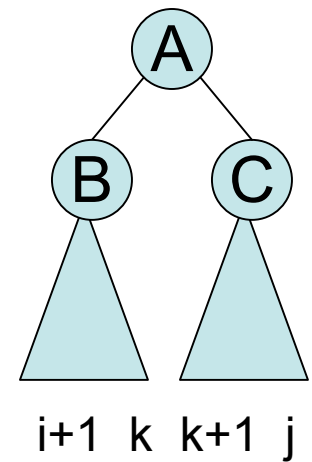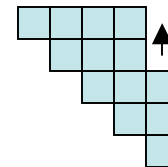    $M[j-1,j] = \{A \mid A \to x_j$ is a rule$\}$

    for i = j-1 down to 1

        $M[i,j] = \cup_{i < k < j} M[i,k] \otimes M[k,j]$

Where $X \otimes Y = \{A \mid A \to BC , B \in X,$ and $C \in Y \}$

Time: $O(n^3)$

# "Inside" Algorithm for SCFG

Just like CKY, but instead of just recording *possibility* of A in M[i,j], record its *probability*:

For each A, do sum instead of union, over all possible k and all possible A $\rightarrow$ BC rules, of products of their respective probabilities.

Result: for each i, j, A, have $\Pr(A \Rightarrow^* x_{i+1}\ldots x_j )$

# The SCFG "Viterbi" algorithm

Like inside, but use max instead of sum;

Gives probability of the *single* parse tree having max probability; (inside sums probability over *all* legal trees)

# ncRNA Discovery in Bacteria

**Cmfinder--A Covariance Model Based RNA Motif  Finding Algorithm**, Yao, Weinberg, Ruzzo,
*Bioinformatics*, 2006, 22(4): 445-452,

**A Computational Pipeline for High Throughput Discovery of cis-Regulatory Noncoding RNA in Prokaryotes**. Yao, Barrick, Weinberg, Neph, Breaker, Tompa and Ruzzo.
*PLoS Comput Biol*. 3(7): e126, July 6, 2007.

**Identification of 22 candidate structured RNAs in bacteria using the CMfinder comparative genomics pipeline**. Weinberg, Barrick, Yao, Roth, Kim, Gore, Wang, Lee, Block, Sudarsan, Neph, Tompa, Ruzzo and Breaker.
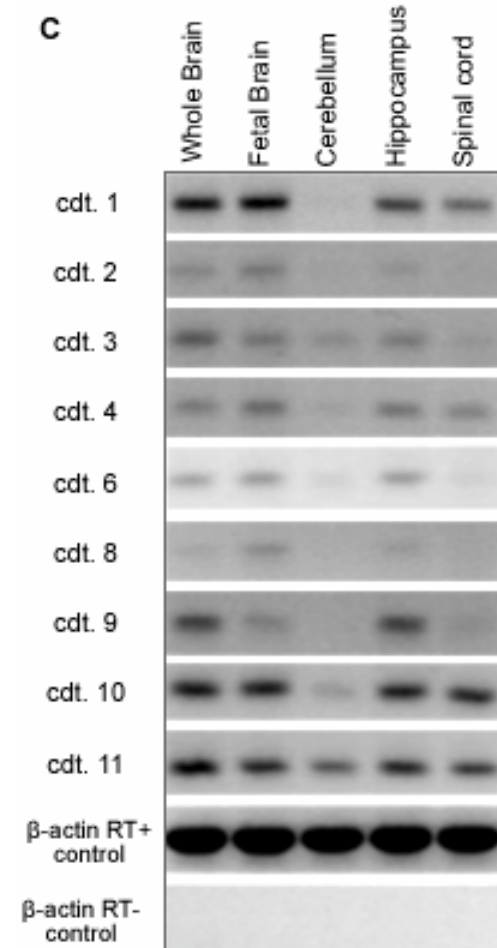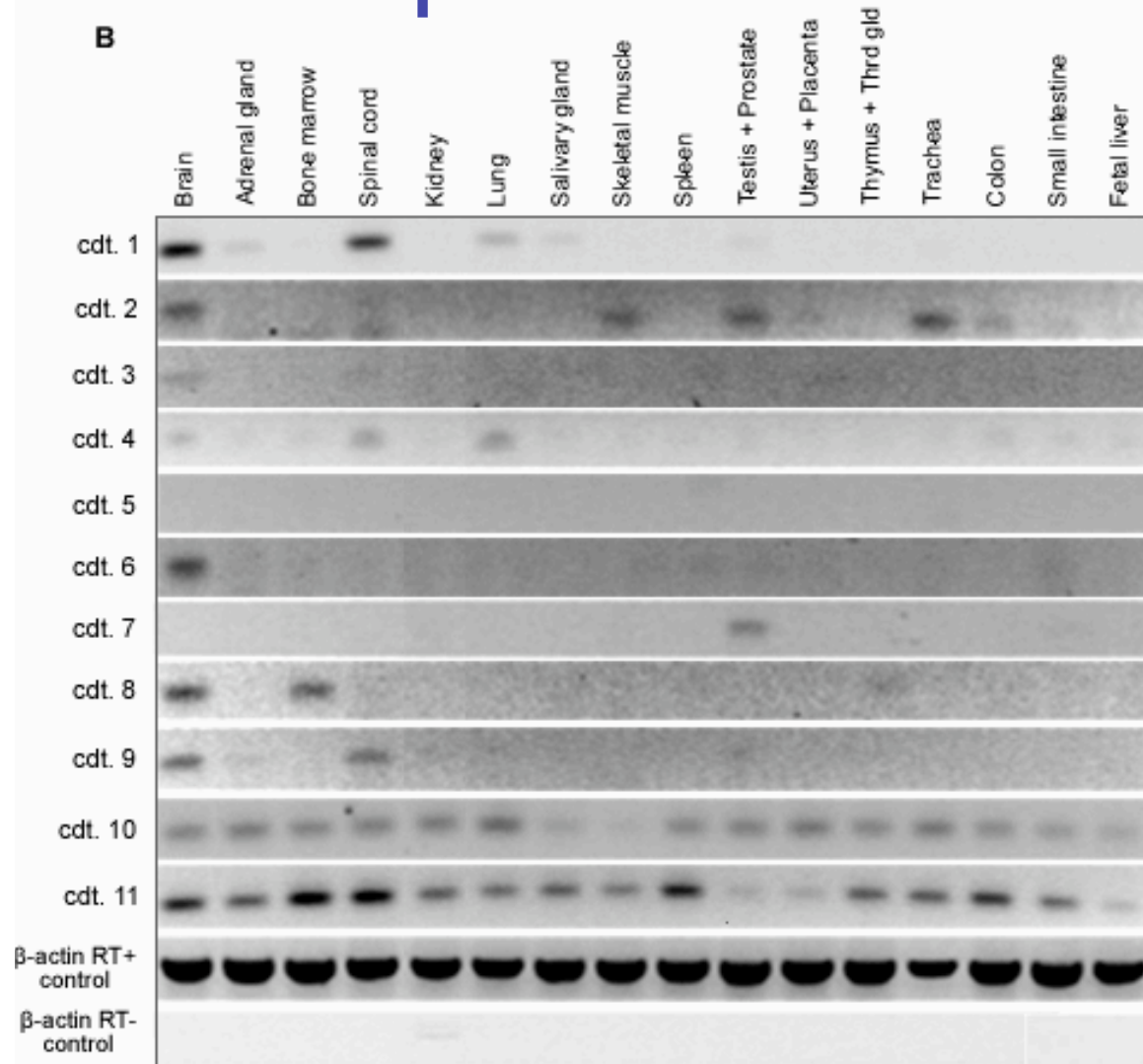*Nucl. Acids Res.,* July 2007 35: 4809-4819.

# ncRNA Discovery in Vertebrates

**Comparative genomics beyond sequence based alignments: RNA structures in the ENCODE regions**

Torarinsson, Yao, Wiklund, Bramsen , Hansen, Kjems, Tommerup, Ruzzo and Gorodkin

Genome Research, to appear

# Experimental Validation

# Bottom Line

CFG technology is a key tool for RNA description, discovery and search

A very active research area. (Some call RNA the "dark matter" of the genome.)

Huge compute hog: results above represent hundreds of CPU-years, and smart algorithms can have a big impact

# More?

Check out CSE 427