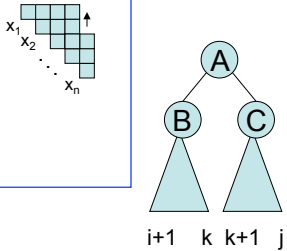# Cocke-Kasami-Younger Parser

Suppose all rules of form A → BC or A → a
(by mechanically transforming grammar)

Given $x = x_1...x_n$, want $M^A_{i,j}$ = {1 if ( $A \Rightarrow^* x_{i+1}...x_j$ ) else 0}

For j=2 to n
    $M^A$ [j-1,j] = {1 if (A → $x_j$ is a rule) else 0}
    for i = j-1 down to 1
        $M^A$ [i,j] = $\bigvee_{A \to BC, \ i < k < j} M^B[i,k] \wedge M^C[k,j]$

Time: O(n³)



i+1   k  k+1  j
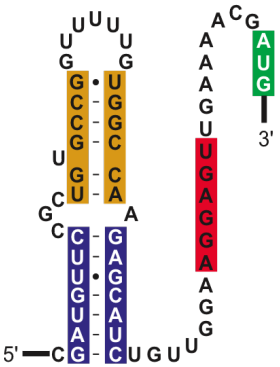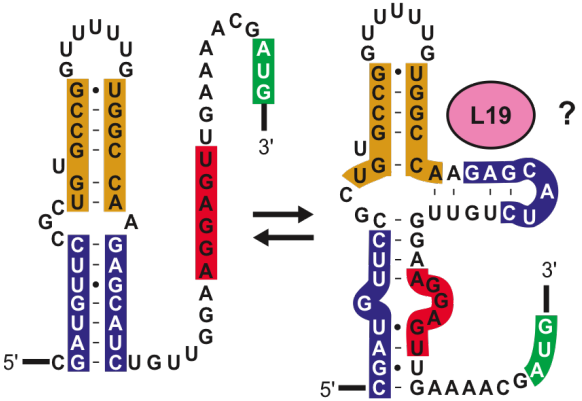
*And now for something completely different …*

# CFGs beyond compilers

# An RNA Structure



# An RNA Computer!
## Sensor & On/Off Switch



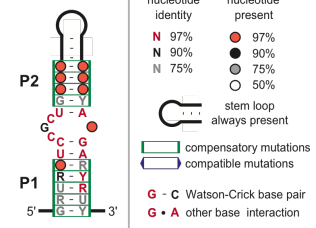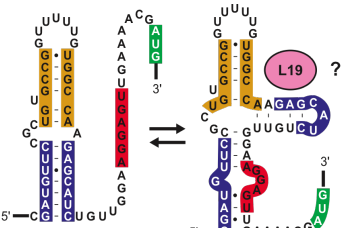L19 absent: Gene On    L19 present: Gene Off

## A  L19 (rplS) mRNA leader

**B**

nucleotide identity  nucleotide present

N 97%  ● 97%
N 90%  ● 90%
N 75%  ● 75%
       ○ 50%

stem loop always present

▭ compensatory mutations
▭ compatible mutations

G – C  Watson-Crick base pair
G • A  other base interaction

P2

P1

5'  3'

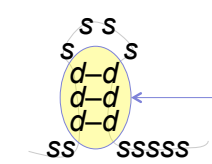**C**  *B. Subtilis* L19 mRNA leader switch?

L19  ?

---

# A CFG for RNA

$S \rightarrow LS \quad | \ L$

$L \rightarrow \text{“dFd”} \ | \ s$

$F \rightarrow \text{“dFd”} \ | \ LS$

"s" means unpaired;
"dFd" means paired
(Watson–Crick:

$aFu \ | \ uFa \ | \ gFc \ | \ cFg$

paren-like nesting)

$S \Rightarrow LS \Rightarrow^* LLLLLLLS$

$\Rightarrow LLLLLLLL$

$\Rightarrow^* ssLsssss$

$\Rightarrow ssdFdsssss$

$\Rightarrow ssddFddsssss$

$\Rightarrow ssdddFdddsssss$

$\Rightarrow \ \dots$

S S
S S S
d–d
d–d
d–d
ss   sssss

---

# Actually, a *Stochastic* CFG

Associate *probabilities* with rules, e.g.:

$S \rightarrow LS$  (p = 0.87)

$S \rightarrow L$  (p = 0.13)

$\dots$

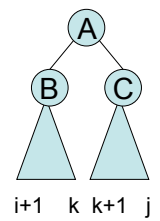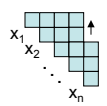Now we can ask, not only
"Does S generate w?"
But also
*"How likely is it?"*

---

# Cocke-Kasami-Younger Parser

Suppose all rules of form A → BC or A → a
(by mechanically transforming grammar)

Given $x = x_1 \dots x_n$, want $M^A_{i,j} = \{1$ if $(A \Rightarrow^* x_{i+1} \dots x_j)$ else $0\}$

For j=2 to n
  $M^A[j\text{-}1,j] = \{1$ if $(A \rightarrow x_j$ is a rule) else $0\}$
  for i = j-1 down to 1
    $M^A[i,j] = \bigvee_{A \rightarrow BC, \ i < k < j} M^B[i,k] \wedge M^C[k,j]$

Time: $O(n^3)$

$x_1$
$x_2$
$x_n$

A
B   C
i+1  k  k+1  j

# "Inside" Algorithm for SCFG

Suppose all rules of form A → BC or A → a
(by mechanically transforming grammar)

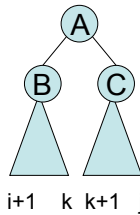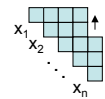Given $x = x_1 \ldots x_n$, want $M^A_{i,j} = p(A \Rightarrow^* x_{i+1} \ldots x_j)$

For j=2 to n

$\quad M^A[j-1,j] = p(\text{rule } A \to x_j)$

$\quad$ for i = j-1 down to 1

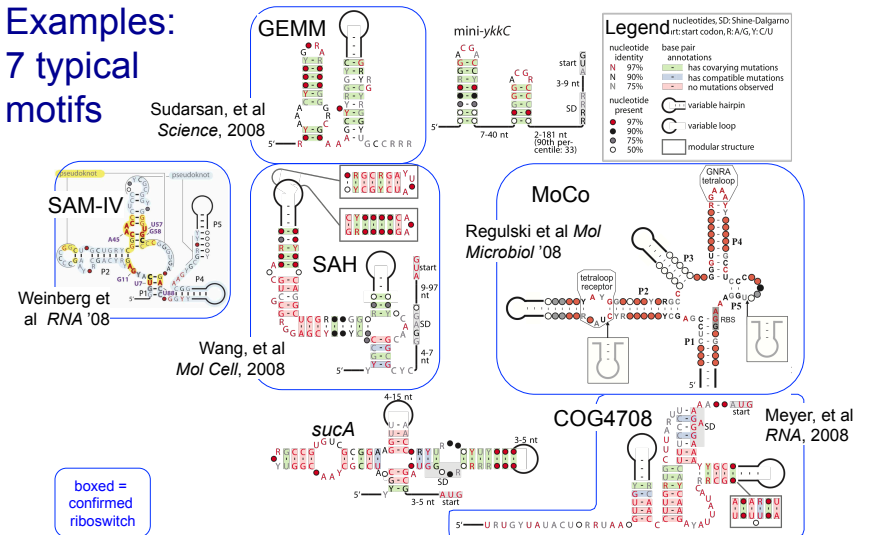$\qquad M^A[i,j] = \sum_{A \to BC, \, i < k < j} M^B[i,k] \times M^C[k,j] \times p(A \to BC)$

I.e., *probability* of A in M[i,j], instead of its *possibility*

Time: $O(n^3)$



---

## Examples: 7 typical motifs



GEMM — Sudarsan, et al *Science*, 2008

mini-*ykkC*

SAM-IV — Weinberg et al *RNA* '08

SAH — Wang, et al *Mol Cell*, 2008

MoCo — Regulski et al *Mol Microbiol* '08

*sucA*

COG4708 — Meyer, et al *RNA*, 2008

boxed = confirmed riboswitch

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

---

## 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; smart algorithms have a big impact

(Recall the $O(n^3)$…)

---

## More?

Check out CSE 427/428: "Comp Bio"