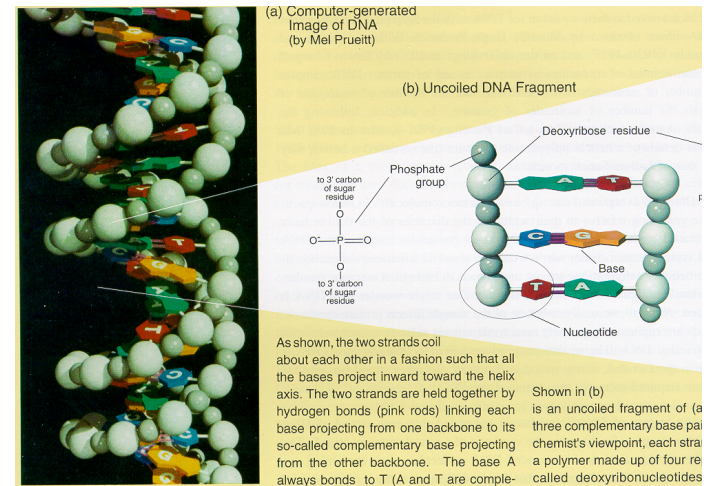


CSE 417: Algorithms and Computational Complexity

Winter 2009
W. L. Ruzzo

Dynamic Programming, II RNA Folding

The Double Helix



<http://www.rcsb.org/pdb/explore.do?structureId=1GAT>

NATURE VOL. 227 AUGUST 8 1970

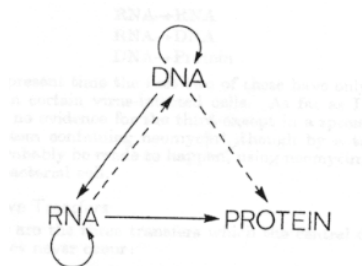
Central Dogma of Molecular Biology

by
FRANCIS CRICK
MRC Laboratory
Hills Road,
Cambridge CB2 2QH

The central dogma of molecular biology deals with the detailed residue-by-residue transfer of sequential information. It states that such information cannot be transferred from protein to either protein or nucleic acid.

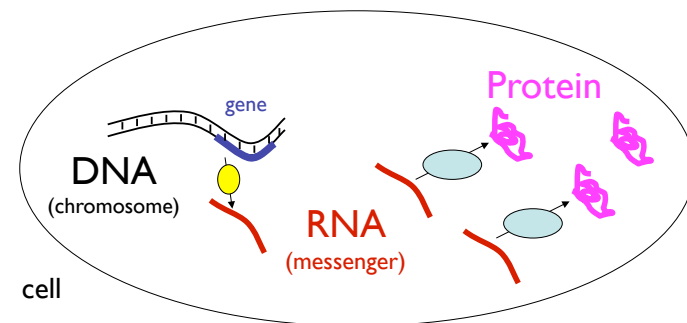
"The central dogma, enunciated by Crick in 1958 and the keystone of molecular biology ever since, is likely to prove a considerable over-simplification."

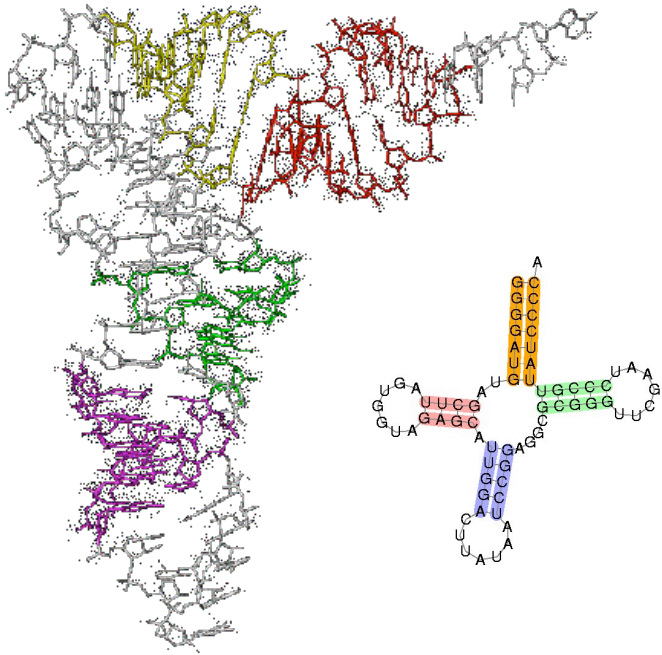
Fig. 2. The arrows show the situation as it seemed in 1958. Solid arrows represent probable transfers, dotted arrows possible transfers. The absent arrows (compare Fig. 1) represent the impossible transfers postulated by the central dogma. They are the three possible arrows starting from protein.



The "Central Dogma" of Molecular Biology

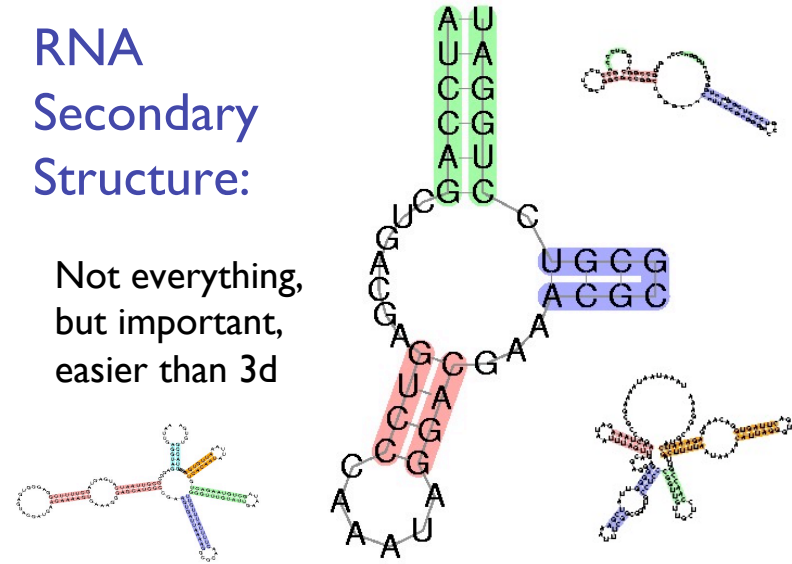
DNA → RNA → Protein





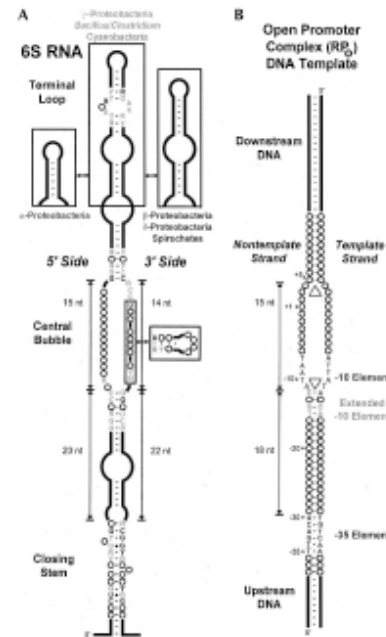
RNA Secondary Structure:

Not everything, but important, easier than 3d

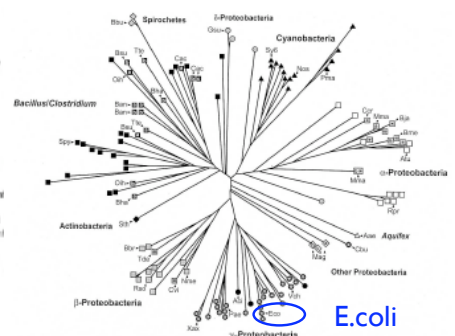
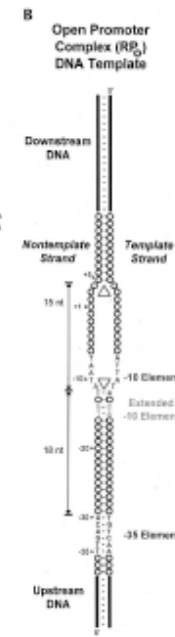


Why is structure important?

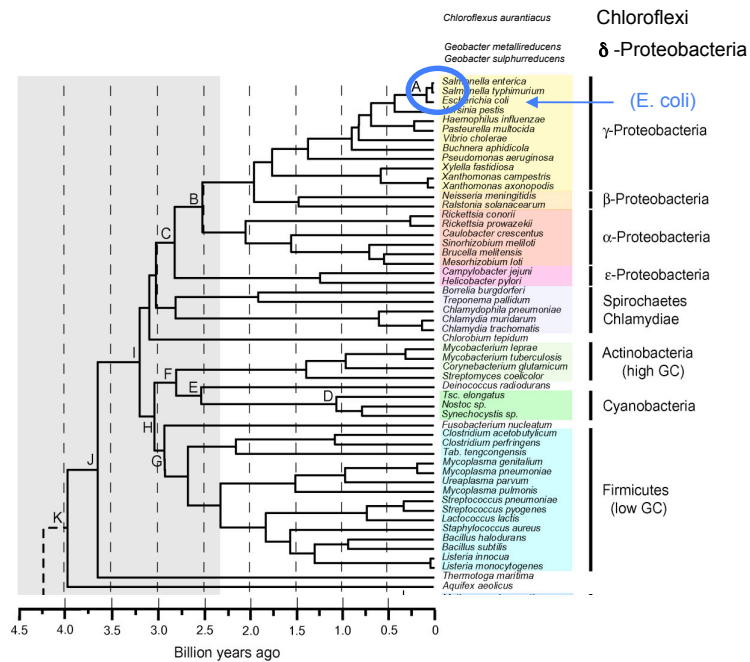
- For protein-coding, similarity in sequence is a powerful tool for finding related sequences
 - e.g. “hemoglobin” is easily recognized in all vertebrates
- For non-coding RNA, many different sequences have the same structure, and structure is most important for function.
 - So, using structure plus sequence, can find related sequences at much greater evolutionary distances



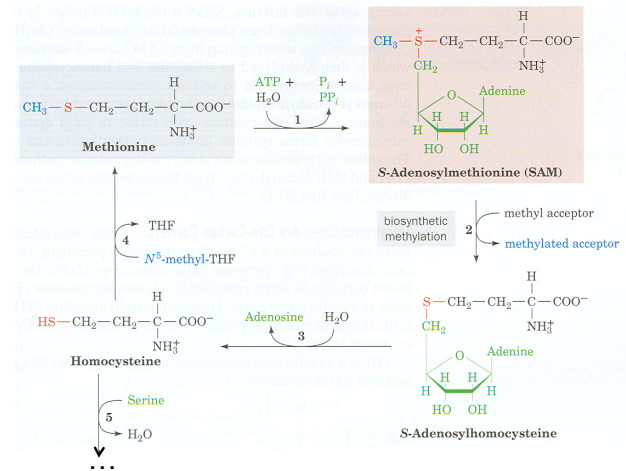
6S mimics an open promoter



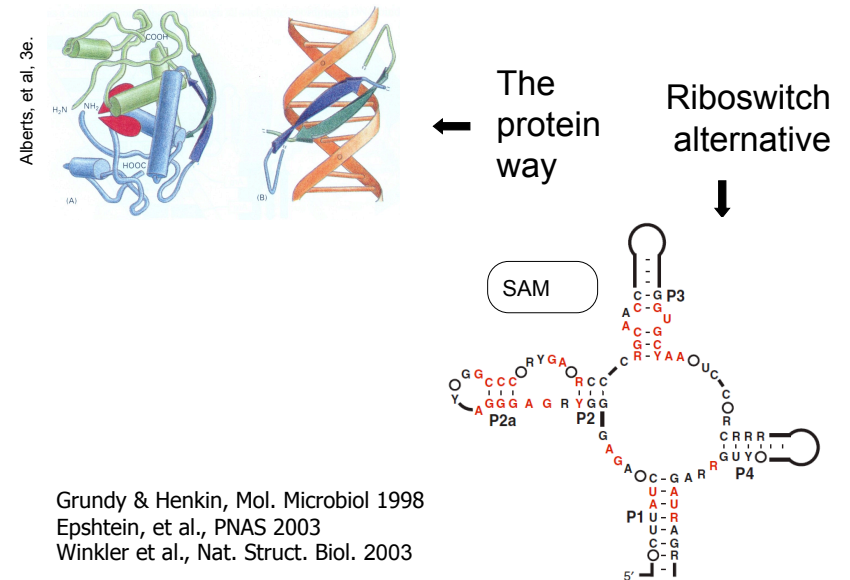
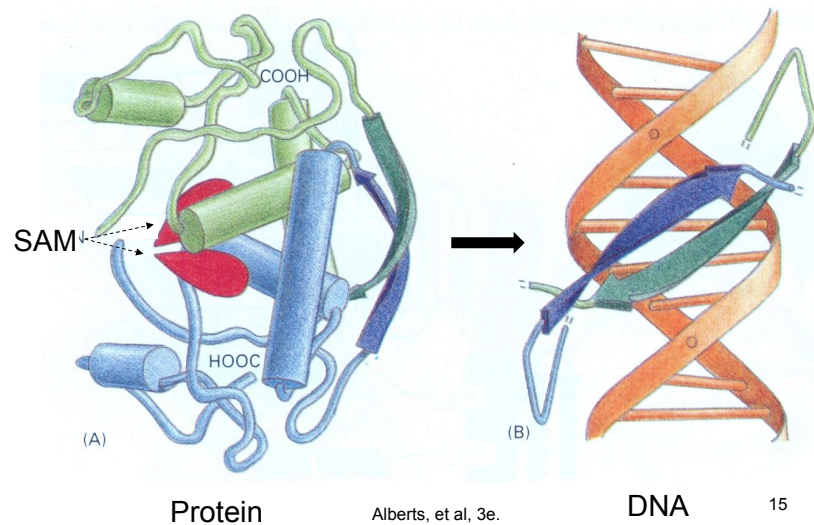
Barrick et al. RNA 2005
 Trotochaud et al. NSMB 2005
 Willkomm et al. NAR 2005



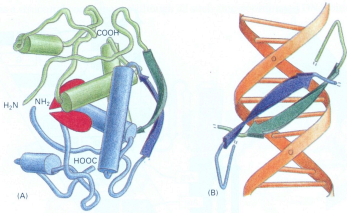
In Bacteria: A typical biosynthetic cycle around a critical metabolite ("SAM")



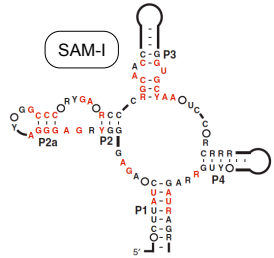
Gene Regulation: The MET Repressor



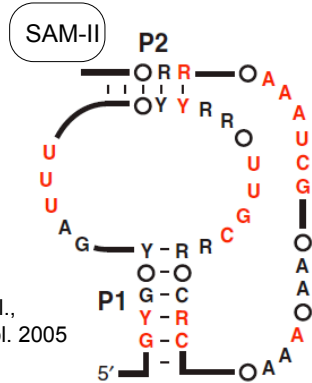
Alberts, et al. 3e.



← The protein way
Riboswitch alternatives

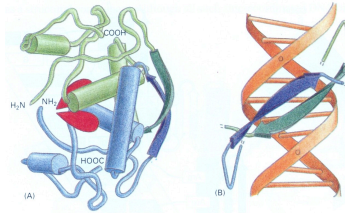


Grundy, Epshtein, Winkler et al., 1998, 2003

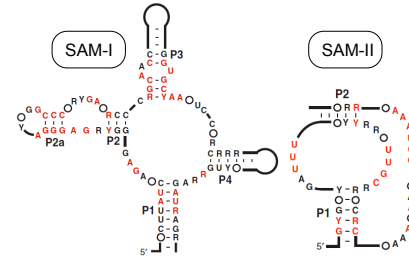


Corbino et al.,
Genome Biol. 2005

Alberts, et al. 3e.

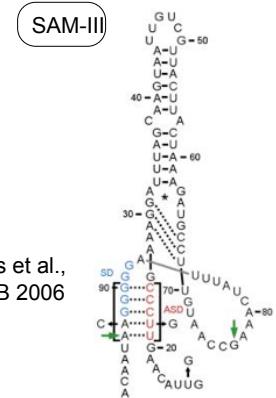


← The protein way
Riboswitch alternatives



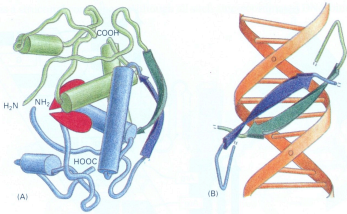
Grundy, Epshtein, Winkler et al., 1998, 2003

Corbino et al.,
Genome Biol. 2005

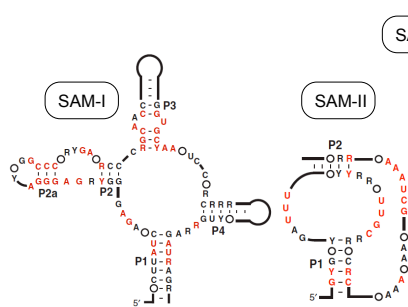


Fuchs et al.,
NSMB 2006

Alberts, et al. 3e.



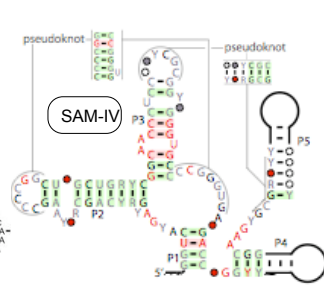
← The protein way
Riboswitch alternatives



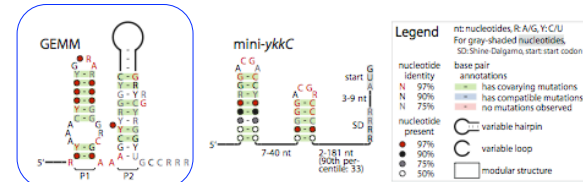
Grundy, Epshtein, Winkler et al., 1998, 2003

Corbino et al.,
Genome Biol. 2005

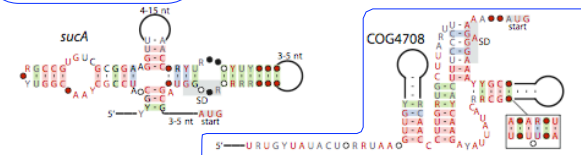
Fuchs et al.,
NSMB 2006



Weinberg et al.,
RNA 2008



And many other examples. Widespread, deeply conserved, structurally sophisticated, functionally diverse, biologically important uses for ncRNA throughout prokaryotic world.



Weinberg, et al. Nucl. Acids Res., July 2007 35: 4809-4819.

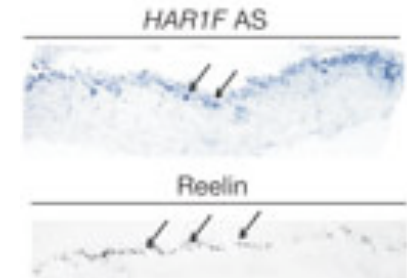
Vertebrates

- Bigger, more complex genomes
- <2% coding
- But >5% conserved in sequence?
- And 50-90% transcribed?
- And *structural* conservation, if any, invisible (without proper alignments, etc.)

• What's going on?

21

Fastest Human Gene?

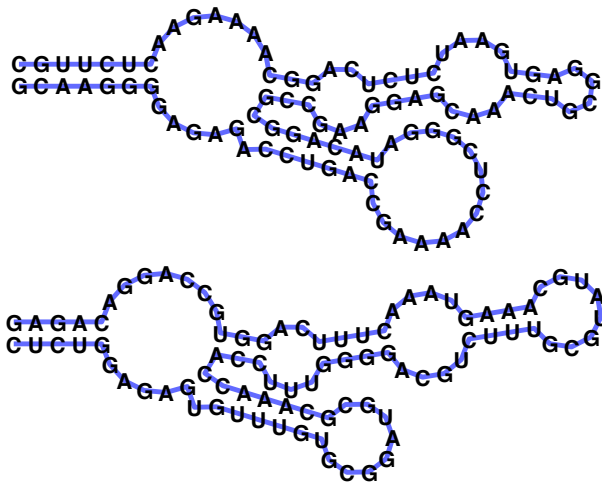


a

Position	20	30	40	50
Human	AGACGTTACAGCAACG	TGTCAGCTGAAATGATGGGGTAGACGCACGT		
Chimpanzee	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Gorilla	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Orang-utan	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Macaque	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Mouse	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Dog	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Cow	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Platypus	ATAAATTACAGCAATTTATCAAAATGAAATATAGGTGTAGACACATGT			
Opossum	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Chicken	AGAAATTACAGCAATTTATCAACTGAAATATAGGTGTAGACACATGT			
Fold	(((((.....)))).....)[][][(((.....)))]			
Pair symbol	l m n o p q r	r q p o n	m l	r s t u v w x x w v u t s r



Q: What's so hard?



A: Structure often more important than sequence

Origin of Life?

Life needs

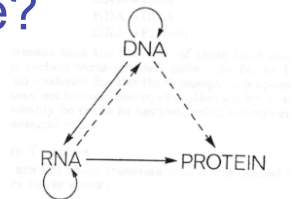
information carrier: DNA

molecular machines, like enzymes: Protein

making proteins needs DNA + RNA + proteins

making (duplicating) DNA needs proteins

Horrible circularities! How could it have arisen in an abiotic environment?



Origin of Life?

RNA can carry information, too

RNA double helix; RNA-directed RNA polymerase

RNA can form complex structures

RNA enzymes exist (ribozymes)

RNA can control, do logic (riboswitches)

The “RNA world” hypothesis:
1st life was RNA-based

6.5 RNA Secondary Structure

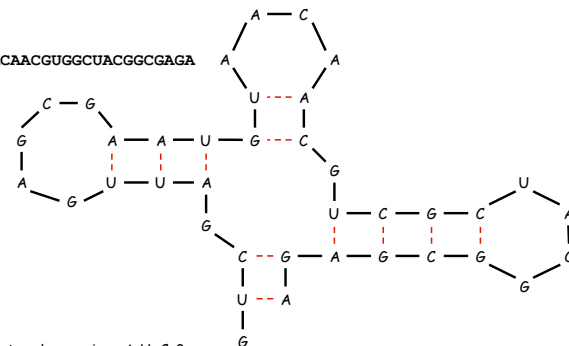
Nussinov’s Algorithm – core technology
for RNA structure prediction

RNA Secondary Structure

RNA. String $B = b_1b_2\dots b_n$ over alphabet $\{A, C, G, U\}$.

Secondary structure. RNA is usually single-stranded, and tends to loop back and form base pairs with itself. This structure is essential for understanding behavior of molecule.

Ex: GUCGAUUGAGCGAAUGUAACAACGUGGCCUACGGCGAGA



RNA Secondary Structure (somewhat oversimplified)

Secondary structure. A set of pairs $S = \{(b_i, b_j)\}$ that satisfy:

- [Watson-Crick.]
 - S is a *matching*, i.e. each base pairs with at most one other, and
 - each pair in S is a Watson-Crick pair: A-U, U-A, C-G, or G-C.
- [No sharp turns.] The ends of each pair are separated by at least 4 intervening bases. If $(b_i, b_j) \in S$, then $i < j - 4$.
- [Non-crossing.] If (b_i, b_j) and (b_k, b_l) are two pairs in S , then we cannot have $i < k < j < l$. (Violation of this is called a *pseudoknot*.)

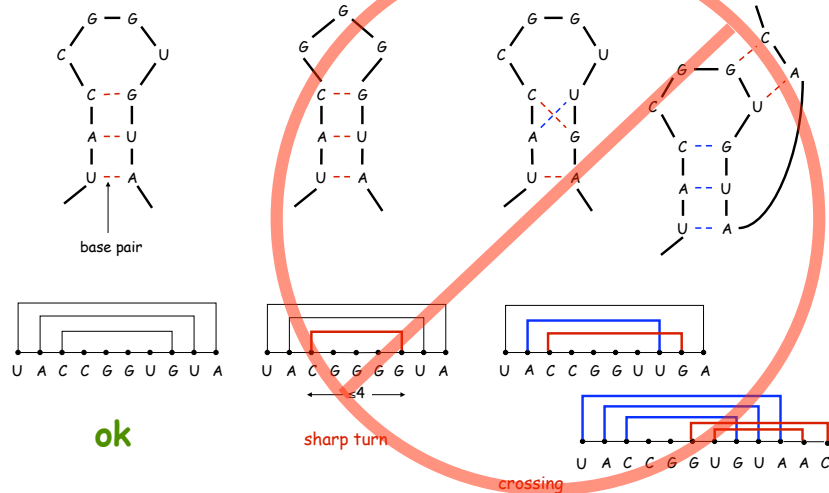
Free energy. Usual hypothesis is that an RNA molecule will form the secondary structure with the optimum total free energy.

approximate by number of base pairs

Goal. Given an RNA molecule $B = b_1b_2\dots b_n$, find a secondary structure S that maximizes the number of base pairs.

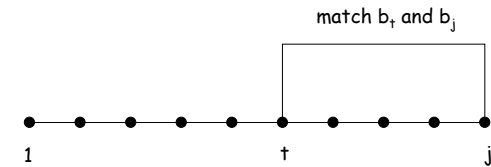
RNA Secondary Structure: Examples

Examples.



RNA Secondary Structure: Subproblems

First attempt. $OPT[j] =$ maximum number of base pairs in a secondary structure of the substring $b_1b_2\dots b_j$.



Difficulty. Results in two sub-problems.

- Finding secondary structure in: $b_1b_2\dots b_{t-1}$. ← $OPT(t-1)$
- Finding secondary structure in: $b_{t+1}b_{t+2}\dots b_{j-1}$. ← not OPT of anything; need more sub-problems

Dynamic Programming Over Intervals: (R. Nussinov's algorithm)

Notation. $OPT[i, j] =$ maximum number of base pairs in a secondary structure of the substring $b_i b_{i+1} \dots b_j$.

- Case 1. If $i \geq j - 4$.
 - $OPT[i, j] = 0$ by no-sharp turns condition.
- Case 2. Base b_j is not involved in a pair.
 - $OPT[i, j] = OPT[i, j-1]$
- Case 3. Base b_j pairs with b_t for some $i \leq t < j - 4$.
 - non-crossing constraint decouples resulting sub-problems
 - $OPT[i, j] = 1 + \max_t \{ OPT[i, t-1] + OPT[t+1, j-1] \}$
 - ↑
take max over t such that $i \leq t < j-4$ and b_t and b_j are Watson-Crick complements

Key point:
Either last base is unpaired (case 1,2) or paired (case 3)

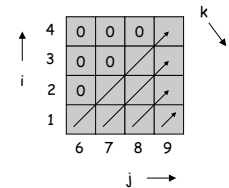
Remark. Same core idea in CKY algorithm to parse context-free grammars.

Bottom Up Dynamic Programming Over Intervals

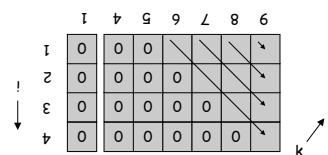
- Q. What order to solve the sub-problems?
A. Do shortest intervals first.

```

RNA( $b_1, \dots, b_n$ ) {
  for  $k = 5, 6, \dots, n-1$ 
    for  $i = 1, 2, \dots, n-k$ 
       $j = i + k$ 
      Compute  $OPT[i, j]$ 
  return  $OPT[1, n]$  using recurrence
}
    
```

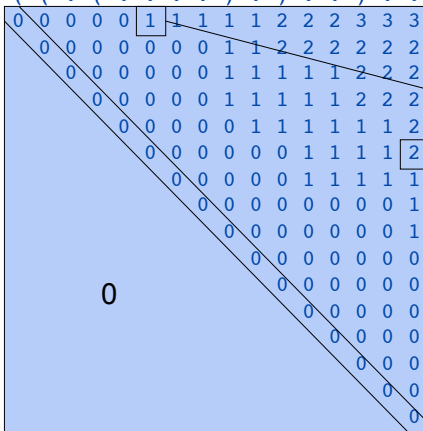


Running time. $O(n^3)$.



C U C C G G U U G C A A U G U C

n = 16



E.g.:
OPT[1,6] = 1:
C U C C G G
(.....)

E.g.:
OPT[6,16] = 2:
G U U G C A A U G U C
((.....).....)

Computing one cell: OPT[2,18] = ?

G G G A A A A C C C A A A G G G G U U U n = 20

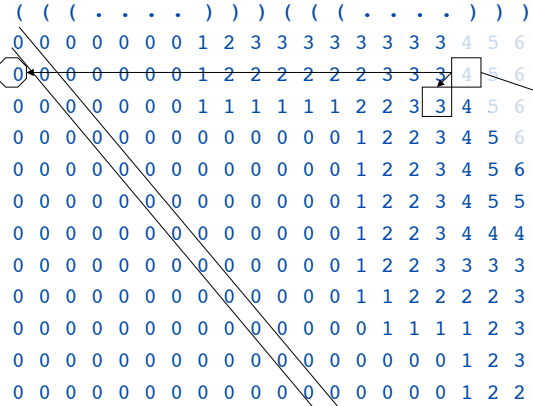


Case 1:
2 ≥ 18-4? no.
Case 2:
B₁₈ unpaired?
Always a possibility;
then OPT[2,18] ≥ 3
G G A A A C C C A A A G G G G U
((.....))(.....)

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$

Computing one cell: OPT[2,18] = ?

G G G A A A A C C C A A A G G G G U U U n = 20

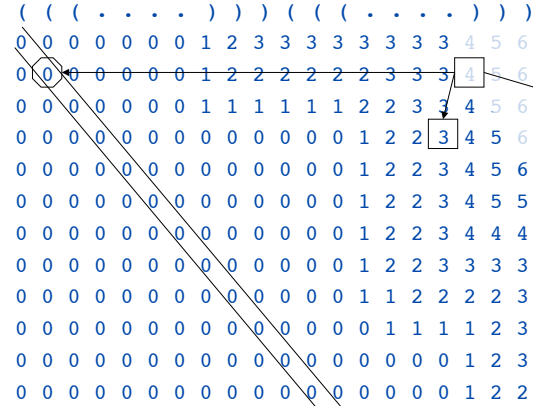


Case 3, 2 ≤ t < 18-4:
t = 2: no pair

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$

Computing one cell: OPT[2,18] = ?

G G G A A A A C C C A A A G G G G U U U n = 20

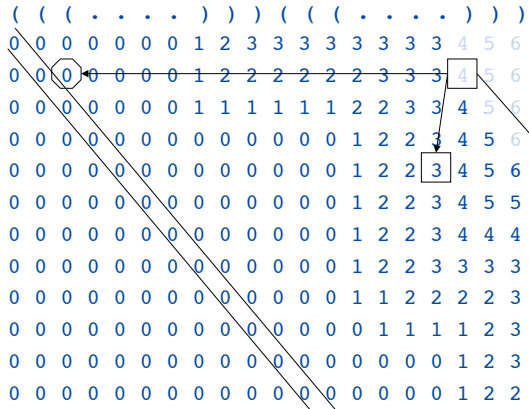


Case 3, 2 ≤ t < 18-4:
t = 3: no pair

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$

Computing one cell: $OPT[2,18] = ?$

G G G A A A A C C C A A A G G G G U U U $n = 20$



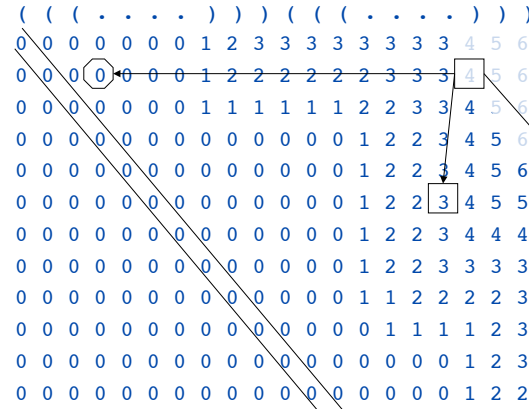
Case 3, $2 \leq t < 18-4$:
 $t = 4$: yes pair
 $OPT[2,18] \geq 1+0+3$
 GGAAAACCCAAAGGGGU
 ..(...(((.....))))

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$



Computing one cell: $OPT[2,18] = ?$

G G G A A A A C C C A A A G G G G U U U $n = 20$



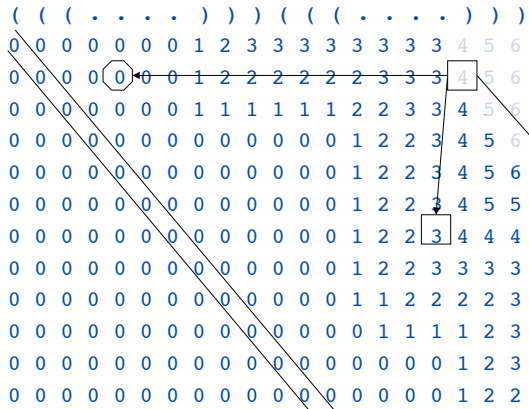
Case 3, $2 \leq t < 18-4$:
 $t = 5$: yes pair
 $OPT[2,18] \geq 1+0+3$
 GGAAAACCCAAAGGGGU
 ..(...(((.....))))

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$



Computing one cell: $OPT[2,18] = ?$

G G G A A A A C C C A A A G G G G U U U $n = 20$



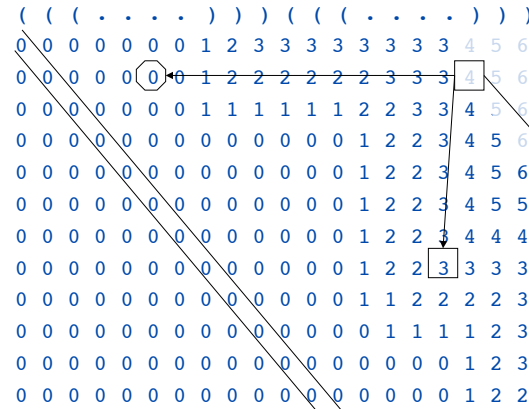
Case 3, $2 \leq t < 18-4$:
 $t = 6$: yes pair
 $OPT[2,18] \geq 1+0+3$
 GGAAAACCCAAAGGGGU
(((.....))))

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$



Computing one cell: $OPT[2,18] = ?$

G G G A A A A C C C A A A G G G G U U U $n = 20$

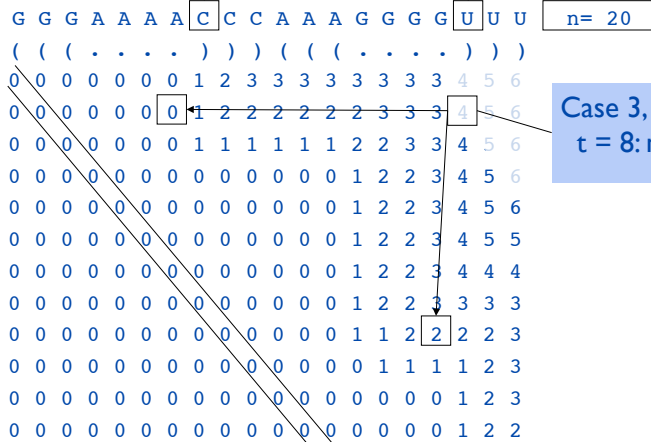


Case 3, $2 \leq t < 18-4$:
 $t = 7$: yes pair
 $OPT[2,18] \geq 1+0+3$
 GGAAAACCCAAAGGGGU
(((.....))))

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$



Computing one cell: OPT[2,18] = ?

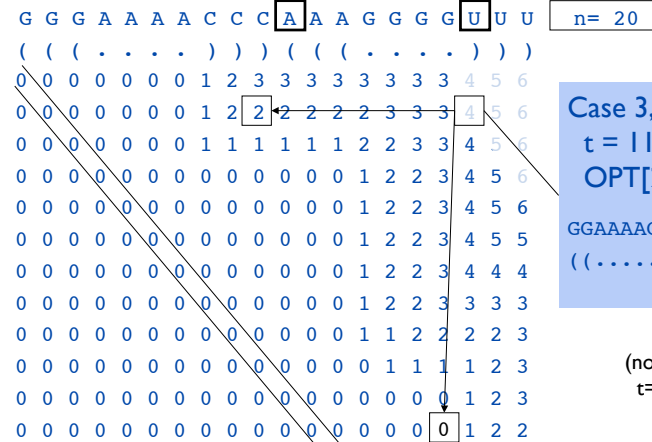


Case 3, $2 \leq t < 18-4$:
 $t = 8$: no pair

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$



Computing one cell: OPT[2,18] = ?



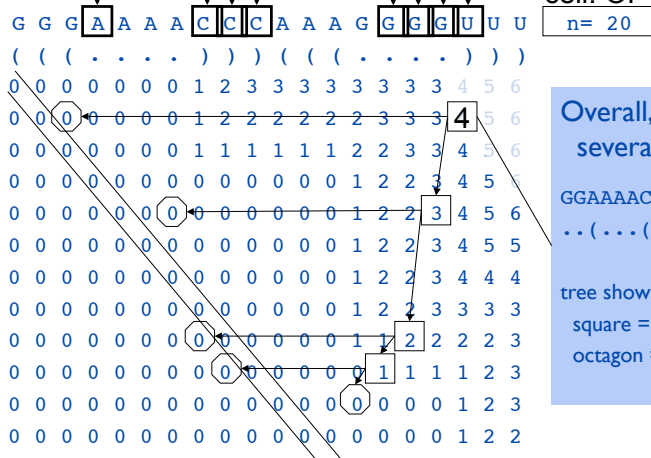
Case 3, $2 \leq t < 18-4$:
 $t = 11$: yes pair
 $OPT[2,18] \geq 1+2+0$
GGAAAACCCAAAGGGUU
((.....))((.....))

(not shown:
 $t=9,10,12,13$)

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$



Computing one cell: OPT[2,18] = 4



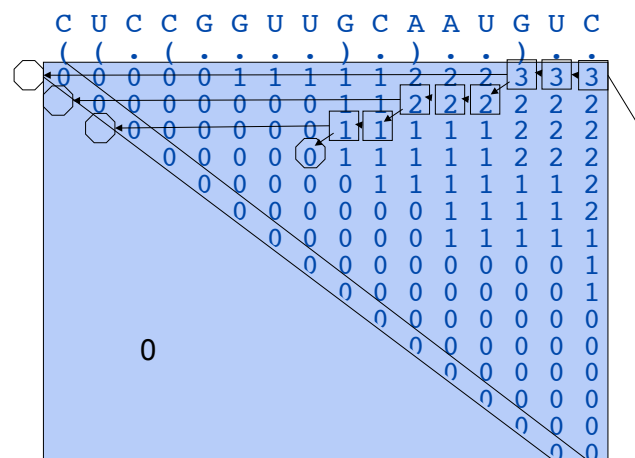
Overall, Max = 4
several ways, e.g.:
GGAAAACCCAAAGGGUU
..((.....))((.....))

tree shows trace back:
square = case 3
octagon = case 1

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$



Another Trace Back Example



E.g.:
 $OPT[1,16] = 3$:
CUCCGGUUGCA AUGUC
((.....))((.....))

$$OPT(i,j) = \begin{cases} 0 & \text{if } i \geq j-4 \\ \max \left\{ \begin{array}{l} OPT[i,j-1] \\ 1 + \max_t (OPT[i,t-1] + OPT[t+1,j-1]) \end{array} \right\} & \text{otherwise} \end{cases}$$

