



ALGORITHMS OF BIOINFORMATICS

8 RNA Structure Prediction

29 January 2026

Prof. Dr. Sebastian Wild

Outline

8 RNA Structure Prediction

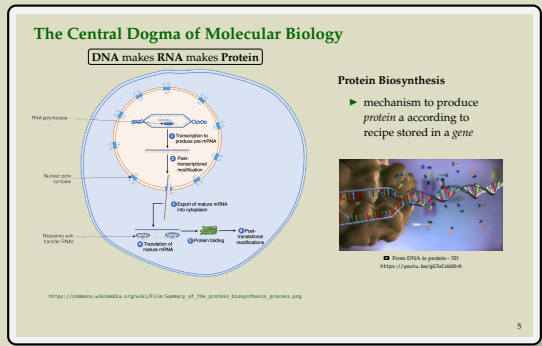
- 8.1 Noncoding RNA
- 8.2 RNA Secondary Structure
- 8.3 Pseudoknot-free secondary structures
- 8.5 Refined Models
- 8.6 Grammar-based Approaches
- 8.7 Probabilistic Parsing
- 8.8 The Grammar of RNA
- 8.9 RNA Prediction and Compression

8.1 Noncoding RNA

RNA

RNA (Ribonucleic acid)

- ▶ similar to DNA: polymer of *nucleotides*
- ↪ sequence of *nitrogenous bases*
Adenine, and *Cytosine*, *Guanine*, *Uracil*
- ▶ unlike DNA, typically *single-stranded*
- ▶ more “sticky” backbone
- ▶ mostly known as *messenger RNA (mRNA)*
 - ▶ including *mRNA vaccines*!
 - ▶ mRNA is a coding RNA since they encode a protein



Noncoding RNA

But RNA serves many other roles!



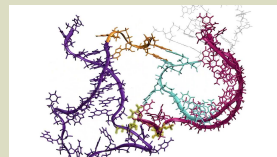
Introduction to Non-Coding RNA
<https://youtu.be/KIohfQsRRdQ>

- ▶ ironically, *ribosomes* (protein factories) themselves are mostly made of RNA
- ▶ for noncoding RNA, structure (3D folding form) crucial for function
- ▶ indeed, sequence often highly variable between species, but structure is similar!

RNA Secondary Structure Prediction

- ▶ Unfortunately, 3D shape hard and expensive to determine experimentally (*X-ray crystallography*)
 - ▶ Available (diverse) data much smaller than for proteins
 - ~> May **not** soon see successful machine-learning solutions similar to AlphaFold
- Rhiju Das, https://youtu.be/XqFq_zYx7Vo

- ▶ To make matters worse, often not a single static structure



▶ RNA folding in action
<https://youtu.be/2XTi9LG9NnU>

- ~> study *de-novo* approaches
- ~> and use simplified models of chemistry and shape to make progress

8.2 RNA Secondary Structure

Model of RNA Structure

- ▶ *RNA sequence / primary structure* $R[0..n) \in \Sigma^n$ $\Sigma = \{A, C, G, U\}$
- ▶ *RNA secondary structure*: matching of indices
 $S \subset [0..n)^2$ of pairs (i, j) that are
 - ▶ ordered $i \leq j$
 - ▶ disjoint: $(i, j), (k, l) \in S \wedge (i = k \vee j = l) \implies (i, j) = (k, l)$
 - ▶ not too close $(i, j) \in S \implies j - i \geq 4$ backbone can't bend more
↖ min. length of hairpin loop
- ▶ secondary structure S is valid for sequence R if
$$(i, j) \in S \implies (R[i], R[j]) \in \mathcal{C} = \{(A, U), (U, A), (C, G), (G, C), (G, U), (U, G)\}$$
- ▶ \mathcal{C} are the *canonical base pairs*: can form *hydrogen bonds* to stabilize RNA

Optimal RNA Structure – Attempt 1

- ▶ Since base pairs provide stability

Try to maximize $|S|$ (# pairs) among all valid secondary structures for $R[0..n]$.

↪ **maximum matching** in graph of all bases

- ▶ possible in polynomial time

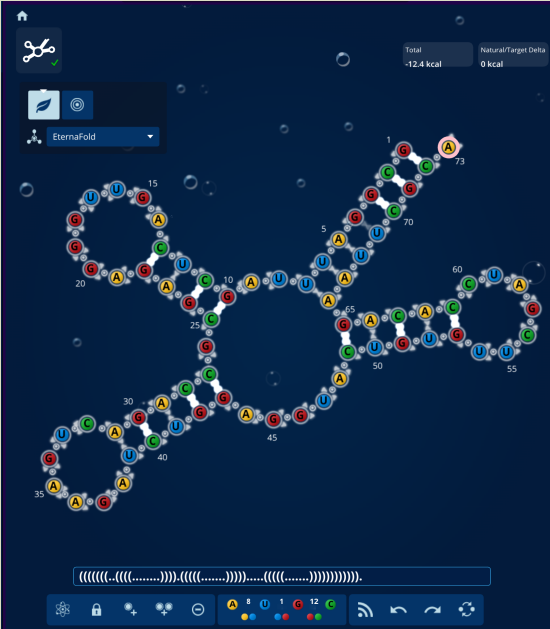
- ▶ actually, ignoring minimum hairpin length, trivial greedy approach is optimal:

1. form arbitrary C – G pairs (until we run out of Cs or Gs)
2. form arbitrary A – U pairs (until we run out)
3. form arbitrary G – U pairs (until we run out)

- ▶ unfortunately, useless predictions!

- ▶ number of pairs dictated by base counts
 - ▶ many equally good options exist
 - ▶ many “optimal” solutions force entire molecule crowd up in one place

Let's play a game!



phenylalanine transfer RNA from *Saccharomyces*
<https://rnacentral.org/rna/URS000011107D/4930>

EteRNA

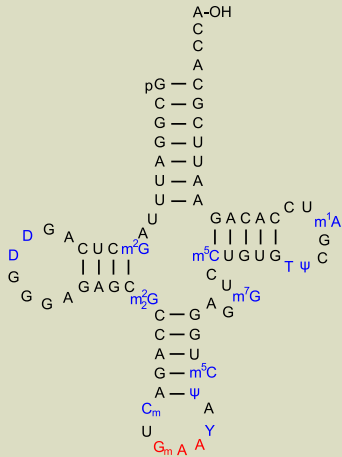
eternagame.org

- ▶ Eterna is a citizen scientist computer game running since 2010 lead by Rhiju Das (Stanford University School of Medicine)
 - ▶ You have to design an RNA sequence that folds into a given *target secondary structure*.
 - ▶ The game uses the best available simulation of RNA folding.
 - ▶ Simulation, prediction, and RNA design algorithms are co-evolving
 - ▶ RNA design crowdsourced to players
 - ▶ top designs synthesized and structure determined
- ~> growing dataset for RNA structures

2D Approximation

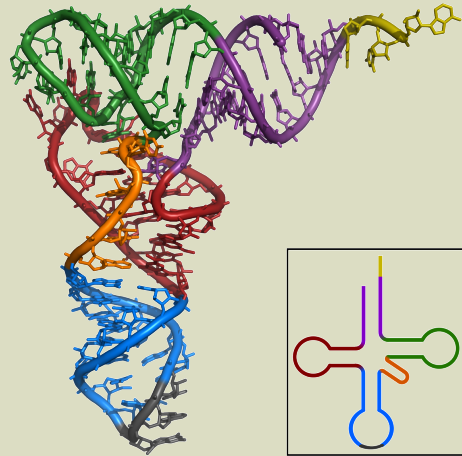
- As in Eterna, RNA secondary structure often drawn as “*roadkill diagrams*”

Roadkill diagram of yeast Phe tRNA

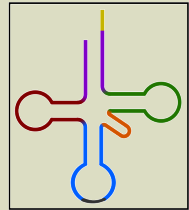


https://commons.wikimedia.org/wiki/File:TRNA-Phe_yeast_blanco.svg

3D Structure of yeast Phe tRNA



https://commons.wikimedia.org/wiki/File:TRNA-Phe_yeast_lehz.png

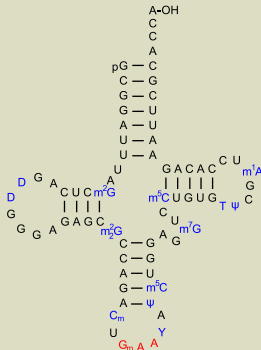


Stacks

Key Observation: Stable structures have many **adjacent pairs**

- ▶ “*stacked*” *pairs* forming a *stem* (the “ladder” regions)
- ▶ in 3D, stems form into a double helix (similar to DNA!)
- ▶ only reverse complement stems are stable

Roadkill diagram of yeast Phe tRNA



https://commons.wikimedia.org/wiki/File:TRNA-Phe_yeast_blanco.svg

3D Structure of yeast Phe tRNA

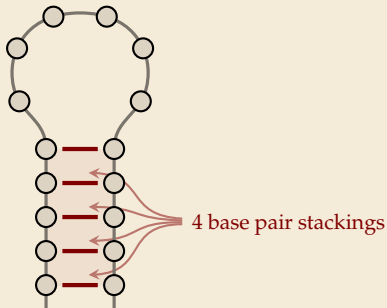


https://commons.wikimedia.org/wiki/File:TRNA-Phe_yeast_1ehz.png

Optimal RNA Structure – Attempt 2

- ▶ Recall: $S \subset [0..n)^2$ set of indices of paired bases
- ▶ instead of maximizing $|S|$ (# pairs), let's maximize number of **base pair stackings**!

$$\text{BPS}(S) = \left| \left\{ (i, j) \in S : (i + 1, j - 1) \in S \right\} \right|$$



General Secondary Structure Prediction

- ▶ **Given:** Sequence $R \in \{A, C, G, U\}^n$
- ▶ **Goal:** *Valid* secondary structure S with maximal $\text{BPS}(S)$

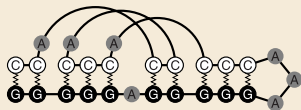
Hardness

Unfortunately, General Secondary Structure Prediction is **NP-hard**.

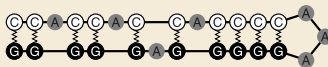
► reduction from BINPACKING



Lyngsø: Complexity of Pseudoknot Prediction in Simple Models, ICALP 2004



(a) An optimum structure for the RNA sequence constructed from an instance of BIN PACKING with four items of sizes 2, 2, 3, and 3, and two bins of capacity 5.



(b) An optimum structure for the RNA sequence constructed from an instance of BIN PACKING with four items of sizes 2, 2, 2, and 4, and two bins of capacity 5.

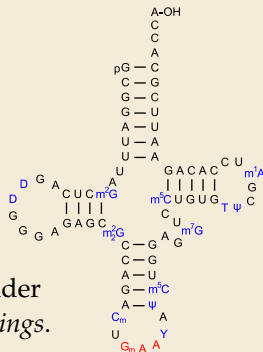
Fig. 3. Illustration of how the number of helices can be kept to one per item for an RNA sequence constructed from a ‘yes’ instance of BIN PACKING, while the base pairs of at least one substring corresponding to an item have to be split over at least two helices if the RNA sequence is constructed from a ‘no’ instance of BIN PACKING.

8.3 Pseudoknot-free secondary structures

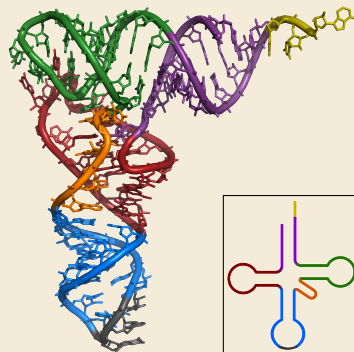
Flat Structures

Recall example tRNA structure

Roadkill diagram of yeast Phe tRNA



3D Structure of yeast Phe tRNA



⇒ Seems reasonable to only consider roadkill diagrams *without crossings*.

https://commons.wikimedia.org/wiki/File:TRNA-Phe_yeast_blanco.svg

https://commons.wikimedia.org/wiki/File:TRNA-Phe_yeast_1ehz.png

“Correct” formalization seems to be:

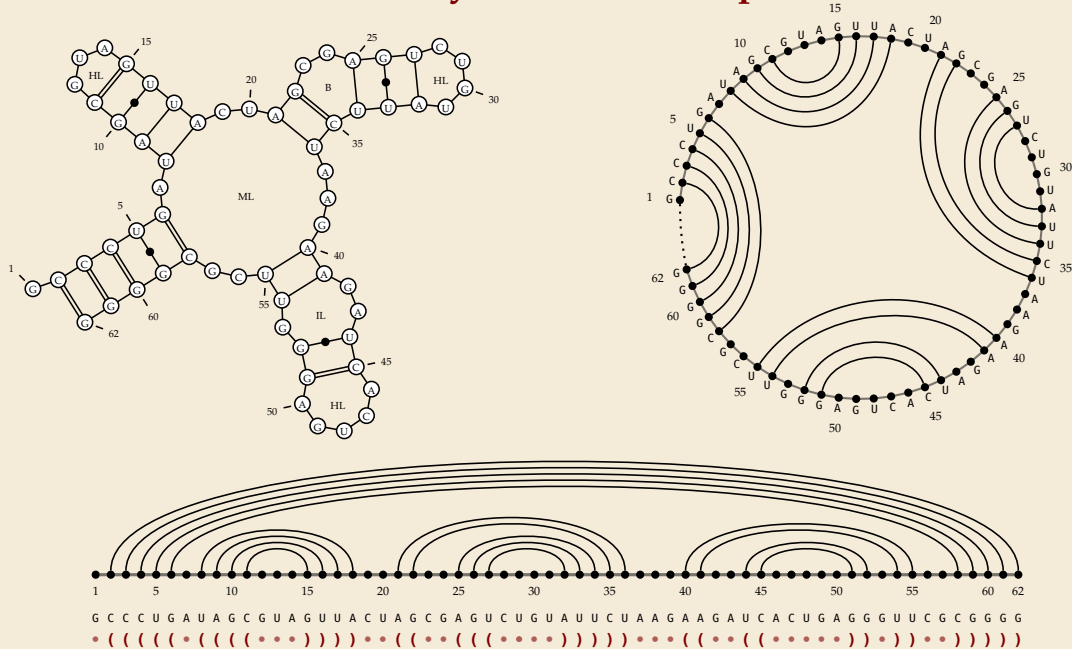
Require graph of pairs bases and backbone edges to be *outerplanar*.

Any other secondary structure is called a *pseudoknot*.

Pseudoknot-free secondary structures

- ▶ planar secondary structure (pairs) cover most of *free energy* of folding
- ▶ “coarse graining” of 3D structure biochemically useful
- ▶ natural intermediate step on folding pathway
- ▶ often well conserved between related species
- ▶ computationally tractable



Pseudoknot-free secondary structures – Representations



Nussinov's Algorithm

*Idea: Maximize total number of valid pairs among all **pseudoknot-free** structures.*

► back to maximum matching, but subject to outerplanar constraint . . .

► **outerplanar** iff $\forall (i, j), (k, \ell) \in S$ either
 $i < k < \ell < j$ (nested ) , or
 $i < j < k < \ell$ (disjoint ) .

Anything else () \rightsquigarrow *pseudoknot*.

► key insight: *decomposability!* see arc diagram / dot-bracket representation



\rightsquigarrow Apply dynamic programming
on subproblems $R[i..j]$

$D(i, j) = \max$ valid pairs in pseudoknot-free structure for $R[i..j]$

Nussinov's Algorithm – DP

$D(i, j) = \max$ valid pairs in any pseudoknot-free structure for $R[i..j]$

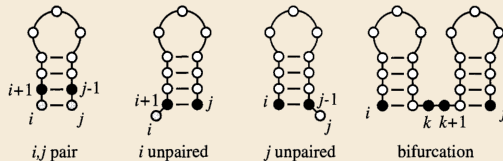


Figure 10.7 from Durbin et al. 1998

$$\rightsquigarrow D(i, j) = \begin{cases} 0, & \text{if } j - i \leq 4; \\ \max \begin{cases} D(i + 1, j - 1) + [(R[i], R[j-1]) \in \mathcal{C}], \\ D(i + 1, j), \\ D(i, j - 1), \\ \max_{k \in [i..j)} D(i, k) + D(k + 1, j) \end{cases} & \text{else.} \end{cases}$$

$\rightsquigarrow O(n^3)$ time, $O(n^2)$ space

8.5 Refined Models

Back to Base Pair Stackings

- ▶ While maximum outerplanar matching is well-defined and tractable, it doesn't usually yield natural structures.
- ▶ already know that we should count base pair stackings!
- ▶ We can extend the DP solution to count those instead!

Graphical notation for DP recursions

$$\text{gray arc} = \min \left\{ \text{empty arc}, \text{arc with 1 sub-arc}, \text{arc with 2 sub-arcs} \right\}$$

Key

- ▶ dots bases; if touching, neighbors on backbone
- ▶ horizontal line RNA backbone
- ▶ wiggly arcs base pair
- ▶ dashed arcs boundary; could be paired or not
- ▶ white area no arcs here
- ▶ gray area potentially further arcs

Counting Base Pair Stackings

Idea: Need to remember whether outermost bases paired.

$$\text{Diagram}(i, j) = \min \left\{ \begin{array}{l} \text{Diagram}(i, j) \\ \text{Diagram}(i, i_1) + \text{Diagram}(j_1, j) \\ \text{Diagram}(i, i+1) + \text{Diagram}(p, p+1) + \text{Diagram}(j-1, j) \end{array} \right\}$$

► In the middle case, if $(i_1, j_1) = (i, j)$, count stacked base pair for (i, j)

$$\text{Diagram}(i, j) = \min \left\{ \begin{array}{l} \text{Diagram}(i, j) \\ \text{Diagram}(i, i+1) + \text{Diagram}(j-1, j) \\ \text{Diagram}(i, p) + \text{Diagram}(p+1, j) \end{array} \right\}$$

↪ Same $O(n^3)$ time, $O(n^2)$ space complexity

Turner Energy Model

- ▶ Simply counting base pair stackings is still a **very crude approximation**
- ▶ Which bases are paired influences bonding strength
- ▶ Which bases are adjacent in stems influences stabilization contribution of stem
- ▶ Which bases form first unpaired base in hairpin loop influences stability
- ▶ ... (play Eterna a bit for more 😊)

↪ More refined models to compute free energy (\approx instability) of structure

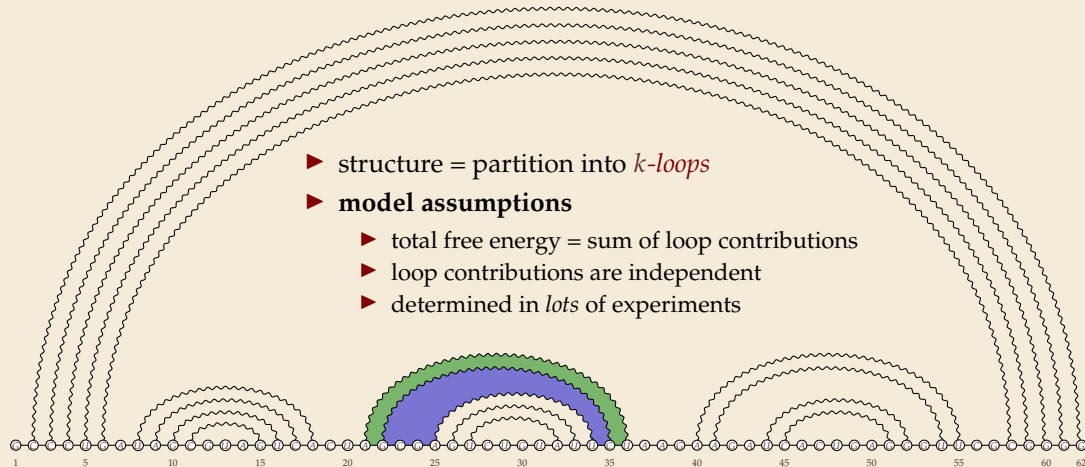


Mathews, Sabina, Zuker, Turner: *Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure*, J Molecul. Biolog. 1999



Mathews, Disney, Childs, Schroeder, Zuker, Turner: *Incorporating chemical modification constraints into a dynamic programming algorithm for prediction of RNA secondary structure*, PNAS 2004

Turner Energy Model [2]



Conceptually unbounded sum

$$\text{shaded loop} = \min \left\{ \text{empty loop}, \text{shaded loop}, \text{shaded loop}, \text{shaded loop}, \dots \right\} \quad \text{⚡ too many variables!}$$

Zuker's Algorithm

- ▶ Only compute exactly up to 2-loops (2 enclosed pairs)
- ▶ additive approximation for bigger multiloops

↪ *same* mutually recursive cost as for pair stackings

8.6 Grammar-based Approaches

Can't machine learning help?

- ▶ free-energy models are great *ab initio* methods
 - ▶ however, they remain limited in accuracy
 - ▶ with growing datasets, tempting to improve structure prediction using machine learning
 - ▶ but: available data much too few for blackbox learning
- ⇒ statistical learning with curated probabilistic model

Probabilistic Context-Free Grammars

Recap from your formal languages intro course ...

Context-free grammars (CFG)

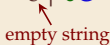
$$G = (N, T, R, S)$$

- ▶ nonterminals N
- ▶ terminals T
- ▶ rules $R \subseteq N \times (N \cup T)^*$
- ▶ start symbol $S \in N$

Applying rules to replace nonterminals

$$S \Rightarrow^* w \rightsquigarrow w \in \mathcal{L}(G)$$

Example

- ▶ $N = \{E, I, V, C, C'\}$
- ▶ $T = \{x, y, 0, \dots, 9, +, \cdot\}$
- ▶ $E \rightarrow (E + E) \mid (E \cdot E) \mid I$
 $I \rightarrow C \mid V$
 $V \rightarrow x \mid y$
 $C \rightarrow 0 \mid 1C' \mid \dots \mid 9C'$
 $C' \rightarrow \varepsilon \mid 0C' \mid \dots \mid 9C'$


Probabilistic Context-Free Grammars (PCFG)

- ▶ For each nonterminal, assign *probabilities* to right-hand sides.

\rightsquigarrow prob of a derivation in G = product of rule probabilities.

RNA Grammars

Example Grammar to produce RNA sequence and structure pairs

(with some fictitious rule probabilities)

rule	prob.	rule	prob.	rule	prob.
$S \rightarrow LS$	0.65	$L \rightarrow \begin{bmatrix} C \\ \cdot \end{bmatrix} S \begin{bmatrix} G \\ \cdot \end{bmatrix}$	0.10	$L \rightarrow \begin{bmatrix} A \\ \cdot \end{bmatrix}$	0.10
$S \rightarrow \varepsilon$	0.35	$L \rightarrow \begin{bmatrix} G \\ \cdot \end{bmatrix} S \begin{bmatrix} C \\ \cdot \end{bmatrix}$	0.05	$L \rightarrow \begin{bmatrix} U \\ \cdot \end{bmatrix}$	0.15
$L \rightarrow \begin{bmatrix} A \\ \cdot \end{bmatrix} S \begin{bmatrix} U \\ \cdot \end{bmatrix}$	0.05	$L \rightarrow \begin{bmatrix} U \\ \cdot \end{bmatrix} S \begin{bmatrix} G \\ \cdot \end{bmatrix}$	0.05	$L \rightarrow \begin{bmatrix} C \\ \cdot \end{bmatrix}$	0.10
$L \rightarrow \begin{bmatrix} U \\ \cdot \end{bmatrix} S \begin{bmatrix} A \\ \cdot \end{bmatrix}$	0.15	$L \rightarrow \begin{bmatrix} G \\ \cdot \end{bmatrix} S \begin{bmatrix} U \\ \cdot \end{bmatrix}$	0.10	$L \rightarrow \begin{bmatrix} G \\ \cdot \end{bmatrix}$	0.15

Can prove: Grammar has unique parse tree *iff* valid RNA sequence structure pair.

Key idea for structure prediction:

- Given only sequence of bases, cannot distinguish $\begin{bmatrix} A \\ \cdot \end{bmatrix}$, $\begin{bmatrix} A \\ \cdot \end{bmatrix}$, and $\begin{bmatrix} A \\ \cdot \end{bmatrix}$
- ↪ Many possible parse trees compatible with sequence
- ↪ Predict structure corresponding to *most likely parse tree*

Structure Prediction with PCFG

rule	prob.	rule	prob.	rule	prob.
$S \rightarrow LS$	0.65	$L \rightarrow \begin{bmatrix} G \\ \cdot \end{bmatrix} S \begin{bmatrix} G \\ \cdot \end{bmatrix}$	0.10	$L \rightarrow \begin{bmatrix} A \\ \cdot \end{bmatrix}$	0.10
$S \rightarrow \varepsilon$	0.35	$L \rightarrow \begin{bmatrix} G \\ \cdot \end{bmatrix} S \begin{bmatrix} C \\ \cdot \end{bmatrix}$	0.05	$L \rightarrow \begin{bmatrix} U \\ \cdot \end{bmatrix}$	0.15
$L \rightarrow \begin{bmatrix} A \\ \cdot \end{bmatrix} S \begin{bmatrix} U \\ \cdot \end{bmatrix}$	0.05	$L \rightarrow \begin{bmatrix} U \\ \cdot \end{bmatrix} S \begin{bmatrix} G \\ \cdot \end{bmatrix}$	0.05	$L \rightarrow \begin{bmatrix} C \\ \cdot \end{bmatrix}$	0.10
$L \rightarrow \begin{bmatrix} U \\ \cdot \end{bmatrix} S \begin{bmatrix} A \\ \cdot \end{bmatrix}$	0.15	$L \rightarrow \begin{bmatrix} G \\ \cdot \end{bmatrix} S \begin{bmatrix} U \\ \cdot \end{bmatrix}$	0.10	$L \rightarrow \begin{bmatrix} G \\ \cdot \end{bmatrix}$	0.15

Leftmost derivations for $R = \text{GAC}$ using $A \equiv \begin{bmatrix} A \\ * \end{bmatrix} \equiv \begin{bmatrix} A \\ \cdot \end{bmatrix} \equiv \begin{bmatrix} A \\ \cdot \end{bmatrix} \equiv \begin{bmatrix} A \\ \cdot \end{bmatrix}$

$$1. \dot{S} \xRightarrow{0.65} \dot{L}S \xRightarrow{0.15} \begin{bmatrix} G \\ \cdot \end{bmatrix} \dot{S} \xRightarrow{0.65} \begin{bmatrix} G \\ \cdot \end{bmatrix} \dot{L}S \xRightarrow{0.10} \begin{bmatrix} G \\ \cdot \end{bmatrix} \begin{bmatrix} A \\ \cdot \end{bmatrix} \dot{S} \xRightarrow{0.65} \begin{bmatrix} G \\ \cdot \end{bmatrix} \begin{bmatrix} A \\ \cdot \end{bmatrix} \dot{L}S \xRightarrow{0.10} \begin{bmatrix} G \\ \cdot \end{bmatrix} \begin{bmatrix} A \\ \cdot \end{bmatrix} \begin{bmatrix} C \\ \cdot \end{bmatrix} \dot{S} \xRightarrow{0.35} \begin{bmatrix} G \\ \cdot \end{bmatrix} \begin{bmatrix} A \\ \cdot \end{bmatrix} \begin{bmatrix} C \\ \cdot \end{bmatrix}$$

Total probability 0.000144

$$2. \dot{S} \xRightarrow{0.65} \dot{L}S \xRightarrow{0.05} \begin{bmatrix} G \\ \cdot \end{bmatrix} \dot{S} \begin{bmatrix} C \\ \cdot \end{bmatrix} S \xRightarrow{0.65} \begin{bmatrix} G \\ \cdot \end{bmatrix} \dot{L}S \begin{bmatrix} C \\ \cdot \end{bmatrix} S \xRightarrow{0.10} \begin{bmatrix} G \\ \cdot \end{bmatrix} \begin{bmatrix} A \\ \cdot \end{bmatrix} \dot{S} \begin{bmatrix} C \\ \cdot \end{bmatrix} S \xRightarrow{0.35} \begin{bmatrix} G \\ \cdot \end{bmatrix} \begin{bmatrix} A \\ \cdot \end{bmatrix} \begin{bmatrix} C \\ \cdot \end{bmatrix} \dot{S} \xRightarrow{0.35} \begin{bmatrix} G \\ \cdot \end{bmatrix} \begin{bmatrix} A \\ \cdot \end{bmatrix} \begin{bmatrix} C \\ \cdot \end{bmatrix}$$

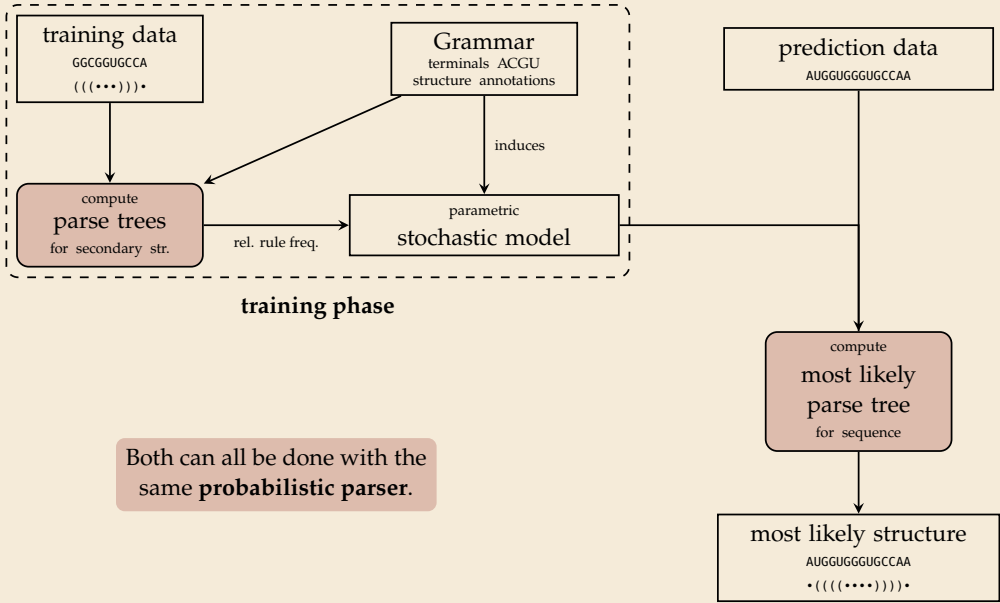
Total probability 0.000259

\rightsquigarrow Predict structure “ (\bullet) ” (Simple grammar ignores min length of hairpin loops!)

Here trivial, since only 2 different derivations . . . how to solve in general?

And where do the probabilities come from?

Stochastic Approach to Secondary Structure Prediction



8.7 Probabilistic Parsing

Refresher: Chomsky Normal Form

Describing parsing easier with specific form of grammars. (reduces case distinctions)

Theorem 8.1 (Chomsky-Normalform (CNF))

For every context-free grammar (CFG) $G = (N, T, R, S)$ with $\varepsilon \notin \mathcal{L}(G)$, we can effectively construct CFG $G' = (N', T, R', S')$ with $\mathcal{L}(G) = \mathcal{L}(G')$ and $R' \subset N' \times N' N' \cup N' \times T$. ◀

\rightsquigarrow Only rules of form $A \rightarrow a$ or $A \rightarrow BC$

Proof idea:

Successively eliminate

1. ε -rules $A \rightarrow \varepsilon$
2. chain rules $A \rightarrow B$
3. “unproductive nonterminals”
i. e., A where there is *no* derivation $S \Rightarrow^* \alpha A \beta \Rightarrow^* w$
4. terminals in RHS $A \rightarrow \alpha a \beta \rightsquigarrow A \rightarrow \alpha A_a \beta$ und $A_a \rightarrow a$
5. overlong RHS $A \rightarrow BCD \rightsquigarrow A \rightarrow BA'$ and $A' \rightarrow CD$

\rightsquigarrow need to verify:

Later steps don't re-introduce
problems eliminated earlier
(easy)

Refresher: Cocke-Younger-Kasami

CNF \rightsquigarrow Only rules of form $A \rightarrow a$ or $A \rightarrow BC$

Given: CFG $G = (N, T, R, S)$ in Chomsky normal form, word $w \in T^n$, $n \geq 1$

Decide: $w \in \mathcal{L}(G)$?



Solve **more general** problem: $A \Rightarrow^* w[i..j] ?$ for **all** $A \in N$, $0 \leq i < j \leq n$.

$\rightsquigarrow w \in \mathcal{L}(G)$ iff $A \Rightarrow^* w[i..j]$ for $A = S$, $i = 0$, $j = n$

► Guess **first used rule** \rightsquigarrow two cases for the two rule types

► $A \Rightarrow^* w[i..i+1)$ iff $A \rightarrow w[i] \in R$

► $A \Rightarrow^* w[i..j)$ for $j \geq i+2$ iff $\exists k \in [i+1..j)$, $B, C \in N$ with $A \rightarrow BC \in R$ and

instances of same problem $\begin{cases} \rightarrow B \Rightarrow^* w[i..k) \text{ and} \\ \rightarrow C \Rightarrow^* w[k..j) \end{cases}$

\rightsquigarrow **dynamic programming** $\rightsquigarrow O(n^3 g)$ time for g the number of rules in G

\rightsquigarrow via backtrace also find a **leftmost derivation** for w (if $w \in \mathcal{L}(G)$)

Probabilistic Parsing

Notation

- ▶ derivation r_1, \dots, r_t = a sequence of rules from R
- ▶ probability of a derivation $\mathbb{P}[r_1, \dots, r_t] = \prod_{i=1}^t P(r_i)$
- ▶ write $\text{lmd}_{A \rightarrow \gamma}(x A \beta) = x \gamma \beta$ for one leftmost derivation step
- ▶ extend $\text{lmd}_r(\circ)$ to several rules $\text{lmd}_{r_1, \dots, r_t}(\circ)$ by successively applying rules

Standard CYK solves the **word/parsing problem**

Given: CFG $G = (N, T, R, S)$ in Chomsky normal form, word $w \in T^n, n \geq 1$

Output: any r_1, \dots, r_t such that $\text{lmd}_{r_1, \dots, r_t}(S) = w$ or NOT_IN_LANGUAGE

For RNA structure prediction, need the most likely among many possible derivations!

Viterbi parse problem

Given: PCFG $G = (N, T, R, S, P)$, word $w \in T^n, n \geq 1$

Output: $r_1, \dots, r_t = \arg \max_{r_1, \dots, r_t} \mathbb{P}[r_1, \dots, r_t]$
 $\text{lmd}_{r_1, \dots, r_t}(S) = w$

Viterbi value vs. probability

Note that in probabilistic parsing, there two “probabilities” for a string w

1. the *Viterbi value*

$$V(w) = \max_{\substack{r_1, \dots, r_t \\ \text{lmd}_{r_1, \dots, r_t}(S) = w}} \mathbb{P}[r_1, \dots, r_t]$$

probability of a single, most likely derivation $S \Rightarrow^* w$

2. the *production probability*

$$\mathbb{P}[w] = \mathbb{P}[S \Rightarrow^* w] = \sum_{\substack{r_1, \dots, r_t \\ \text{lmd}_{r_1, \dots, r_t}(S) = w}} \mathbb{P}[r_1, \dots, r_t]$$

total probability that we obtain w

when starting with S and applying randomly chosen rules to expand nonterminals.

Both can make sense; for us only Viterbi value needed.

Probabilistic CYK

In Chomsky normal form, easy to make CYK probabilistic

General Viterbi values: $V_A(w) = \max_{r_1, \dots, r_t} \mathbb{P}[r_1, \dots, r_t]$
 $\text{Imd}_{r_1, \dots, r_t}(A) = w$

$$V_A(w[i..i+1]) = \begin{cases} P(A \rightarrow w[i]) & \text{if } A \rightarrow w[i] \in R \\ 0 & \text{otherwise} \end{cases}$$

$$V_A(w[i..j]) = \max_{\substack{k \in [i+1..j] \\ A \rightarrow BC \in R}} P(A \rightarrow BC) V_B(w[i..k]) V_C(w[k..j])$$

\leadsto Same DP works!

$O(n^3 g)$ time, backtrace yields most likely derivation

8.8 The Grammar of RNA

RNA Grammar Normal Form

Chomsky Normal Form rather inconvenient for RNA.

- ▶ Really want a single rule to introduce a pair bond.
- ▶ Nonterminals often represent a logical abstraction of several special cases would want to have chain rules for that

Stochastic RNA Form (SRF)

Normal form from



Onokpasa, Wild, Wong: *Towards Optimal Grammars for RNA Structures*, DCC 2024

- ▶ Assume $N = \{A_1, \dots, A_k\}$
- ▶ *start symbol* must be A_k
- ▶ all rules must of the form
 - $A_i \rightarrow A_j A_\ell$
 - $A_i \rightarrow \bullet$
 - $A_i \rightarrow (A_j)$
 - $A_i \rightarrow A_j$ and $j < i$(assumes a *total order* on the nonterminals)

Viterbi values

Generalization of CYK computes Viterbi values

$$V_A(w[i..i+1]) = \begin{cases} P(A \rightarrow w[i]) & \text{if } A \rightarrow w[i] \in R \quad \text{type (ii) rules} \\ 0 & \text{otherwise} \end{cases}$$

$$V_A(w[i..j]) = \max \begin{cases} \max_{\substack{k \in [i+1..j] \\ A \rightarrow BC \in R}} P(A \rightarrow BC) V_B(w[i..k]) V_C(w[k..j]) & \text{type (i) rules} \\ \max_{A \rightarrow w[i] B w[j-1] \in R} P(A \rightarrow w[i] B w[j-1]) V_B(w[i+1..j-1]) & \text{type (iii) rules} \\ \max_{A \rightarrow B \in R} P(A \rightarrow B) V_B(w[i..j]) & \text{type (iv) rules} \end{cases}$$

Towards Subcubic RNA Folding

- ▶ both general context-free parsing as well as minimum-free energy models led to $\Theta(n^3)$ worst case time algorithms
- ▶ at core, both problems resemble matrix multiplication
 - ▶ dominant running time comes from bifurcation rules

▶ e. g., Recall Nussinov $D(i, j) = \begin{cases} 0, & \text{if } j - i \leq 4; \\ \max \begin{cases} D(i + 1, j - 1) + [(R[i], R[j-1]) \in \mathcal{C}], \\ D(i + 1, j), \\ D(i, j - 1), \\ \boxed{\max_{k \in [i..j)} D(i, k) + D(k + 1, j)} \end{cases} & \text{else.} \end{cases}$

Towards Subcubic RNA Folding [2]

Approaches

1. Exhaustive Tabulation can be made to work ... shaves off \log -factors
2. for (probabilistic) context-free parsing *Valiant's algorithm* reduces the problem to (Boolean) matrix multiplication $\rightsquigarrow O(n^\omega)$ possible



Valiant: *General context-free recognition in less than cubic time*, J Comput. System Sci. 1975

- ▶ constant factors are huge; even asymptotically slower Strassen
- ▶ beyond plain word recognition problem, more complications to be overcome



Zakov, Tsur, Ziv-Ukelson: *Reducing the worst case running times of a family of RNA and CFG problems, using Valiant's approach*, Alg. Molecul. Biol. 2011

3. Closer to $(\min, +)$ -matrix multiplication than regular matrix $(+, \cdot)$ -multiplication
 - ▶ in general, likely no $O(n^{3-\epsilon})$ algorithms to be expected
 - ▶ but for bounded difference matrices, speedups known



Bringmann, Grandoni, Saha, Williams: *Truly Subcubic Algorithms for Language Edit Distance and RNA Folding via Fast Bounded-Difference Min-Plus Product*, SIAM J. Comput. 2019

What Grammars to Choose?

- ▶ #parameters = #rules -1 per nonterminal
- ▶ more parameters means more expressive
but also more prone to overfitting
- ▶ finding a grammar that captures structure well at this point *more art than science*

~> Good results, e. g., for this grammar

$$U \rightarrow \bullet$$
$$B \rightarrow (M)$$
$$T \rightarrow B \mid U$$
$$M \rightarrow B \mid TS \mid T$$
$$S \rightarrow TS \mid T$$

Grammar G_6 from



Dowell, Eddy: *Evaluation of several lightweight stochastic context-free grammars for RNA secondary structure prediction*, BMC Bioinform. 2004

8.9 RNA Prediction and Compression


Prediction = Compression + Generalization

- ▶ RNA grammar turns RNA sequence-structure pair into single string:
The sequence of used rules!
 - ▶ Can compress these rules using arithmetic coding
- ↪ efficient compression of RNA

Can we use the compression ratio as a proxy to compare RNA grammars?

 Onokpasa, Wild, Wong: RNA secondary structures: from ab initio prediction to better compression, and back, DCC 2023

Can we find better grammars than Dowell & Eddy?

 Onokpasa, Wild, Wong: Towards Optimal Grammars for RNA Structures, DCC 2024