

# Designing RNA Structures is Hard

- A report on the paper by Édouard Bonnet, Paweł Rzążewski, and Florian Sikora

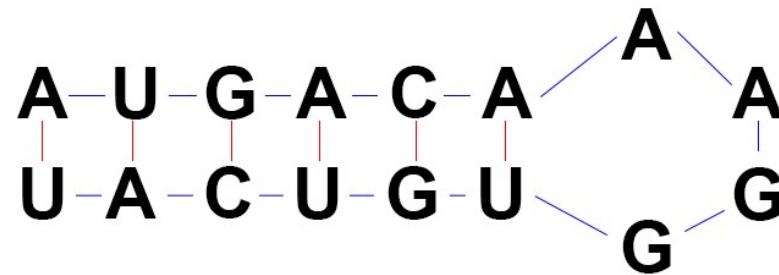
# RNA Secondary Structure Design

- Secondary structure is important in the biological function of the molecule
- Solving the problem has real world applications:
  - Pharmaceutical research
  - Biochemistry
  - Synthetic biology
  - RNA nanostructures

# RNA Structures

- RNA is made up of 4 nucleotides, labeled A, U, C, and G
- The **primary** structure of an RNA molecule is just a string over the alphabet {A, U, C, G} – each element is a **base**
- Nucleotides can bind together: A with U, C with G in the Watson-Crick energy model
- Each binding reduces the free energy of the molecule
- The **secondary** structure of a molecule is the set of positional pair bindings

# RNA Structures



- This structure is called a **stem loop**
- In a **pseudo-knot-free** structure, no stem contains part of a stem from another stem loop
- A **minimum free energy** (MFE) structure is one with the maximal number of bases paired (each pair is worth -1 energy)
- Sequences want to fold into an MFE configuration



# Robustness of Proof

- Uses Watson-Crick energy model (simplest)
  - Other models are not uniform, more complex
  - Prove hardness independently of energy model
- Ignore pseudo-knots
  - Again, proving hardness in the easier case
- Reductions structure maps well to stem loops
  - This is realistic, stem loops are a basic RNA building block

# RNA-DESIGN-EXTENSION is in NP

- The inverse problem is RNA-FOLDING
  - Given an RNA sequence, compute its MFE folding
  - Can be solved with DP in polynomial time
- We can use RNA-FOLDING as an oracle to verify solutions to RNA-DESIGN-EXTENSION
- Having such an oracle means that RNA-DESIGN-EXTENSION is in NP

# Proving RDE is NP-Complete

- Reduce from E3-SAT
  - Each clause contains 3 distinct variables
  - Known NP-Hard if each variable is used up to 4 times
- Map SAT instance onto a string representation of RDE
- Label the string representation with bases
- Show that the base sequence labels are a solution to the RDE instance iff the SAT instance is satisfiable



# Representing Secondary Structures

- A **structure** a well-parenthesized expression with dots (((..)(..)))
- ()'s represent a base pair, . represents an unpaired base
- A **sequence** is a string of the same length from {1,2,3,4}
  - 1=A, 2=C, 3=G, 4=U so proper pairs sum to 5
  - Sequence w corresponds to structure S if proper pairs match
- Sequence is a **design** if it can't fold into anything with more pairs
- A **partial sequence** is a sequence with some ?'s (unassigned)
- A partial sequence w' is a **design extension** if the ?'s can be filled in to turn it into a design.



## Building the RDE instance: clauses

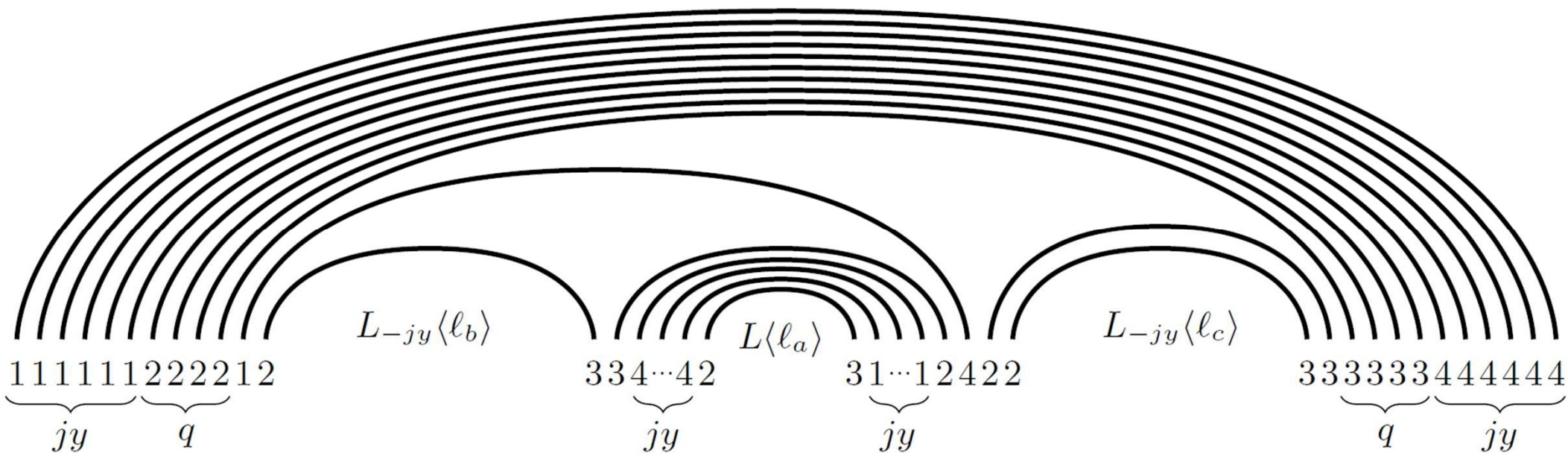
- Let a clause  $C_j$  contain three literals  $l_a, l_b, l_c$ ;  $a < b < c$ 
  - $l_a$  is the same gadget as  $V\langle X_a \rangle$
  - $l_b$  and  $l_c$  are  $V\langle X_b \rangle, V\langle X_c \rangle$  with  $jy$  parentheses pairs remove
  - The clause is

$$S\langle C_j \rangle := \underbrace{(\dots)}_{jy} \underbrace{((\dots))}_{q} (((L_{-jy}\langle l_b \rangle) (\underbrace{(\dots)}_{jy} ((L\langle l_a \rangle) \dots))) \underbrace{((L_{-jy}\langle l_c \rangle))}_{jy} \dots) \dots) \underbrace{)}_{q} \underbrace{)}_{jy}$$

- The  $jy+q$  parentheses are the **arch** of the clause
  - The outer  $jy$  parentheses are the first arch layer
  - The inner  $q$  parentheses are the second arch layer

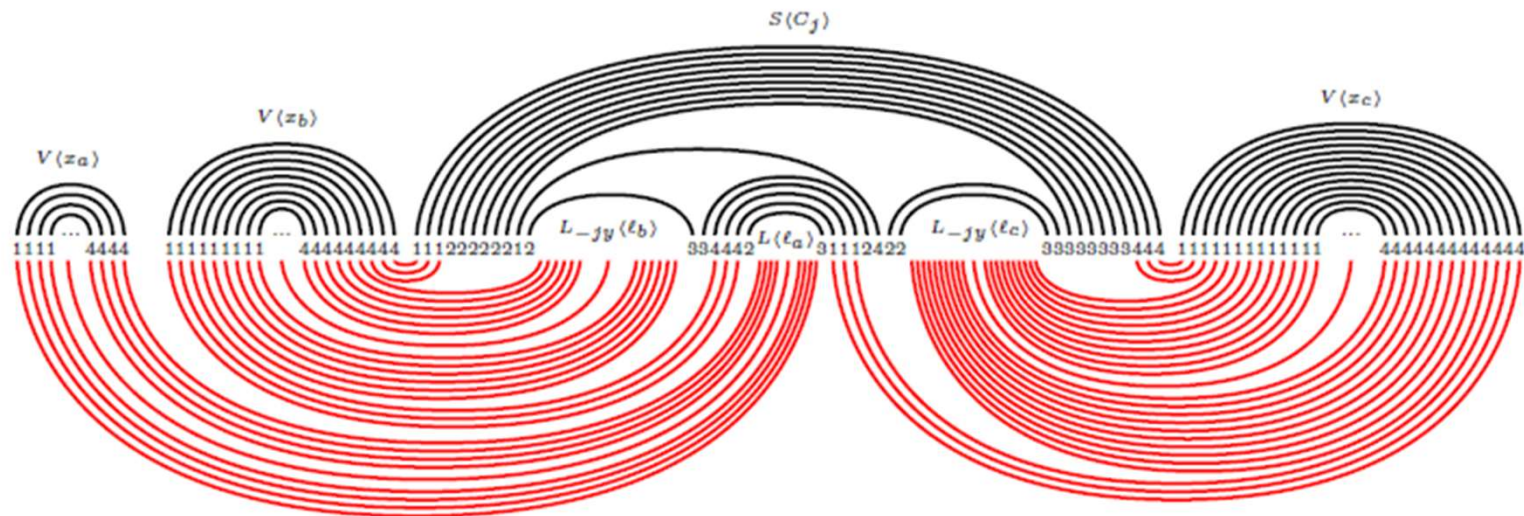
# Building the RDE instance

- The entire instance is clauses interleaved with variables such that every clause is between the corresponding  $V\langle Xb \rangle$  and  $V\langle Xc \rangle$  gadgets



# SAT unsatisfiable -> no design extension

- Black lines are by given construction, red lines are a re-matching by removing some parentheses
- If clause is unsatisfiable, then it is possible to rematch to a better structure
- This implies the originally constructed structure is NOT a design extension



## SAT satisfiable $\rightarrow$ design extension

- This direction is too gory for full details
- Given a structure  $S$  with a satisfiable SAT instance, assume
  - There is a better structure  $S'$  for the corresponding sequence  $w$
  - Assume wlog that  $S'$  actually is the maximal matching
- By an argument of counting matching parts of  $S'$ , can prove that there exists  $S''$  which matches even MORE bases in  $w$
- But this is a contradiction, so  $S'$  cannot exist
- Therefore  $S$  itself must be the maximal matching and is a design extension
- We have proven SAT satisfiable iff the corresponding sequence is a design extension, thus RDE is NP-Hard
- We showed earlier that RDE is in NP, thus RDE is NP-Complete

# Algorithmic Consequences

- Taking advantage of the structures in the proof leads to a faster algorithm
- Naïve:  $O^*(4^n)$
- Using insights from the proof, can prune search space
  - $\sqrt{3}^n n^{O(1)}$ , where  $n$  is the length of the input structure
  - $2^s n^{O(1)}$ , where  $s$  is the number of unlabeled elements in the input structure
- RNA-DESIGN is known to be in P for saturated structures
  - Using DP based on ideas from the proof, RDE is also tractable on saturated structures
- There are many other tree-structured problems in computer science
  - This gadget mapping may have uses elsewhere

**Questions?**