

RNA-RNA Interaction Prediction with Stochastic Grammars

Sebastian Wild Markus Nebel Anika Scheid

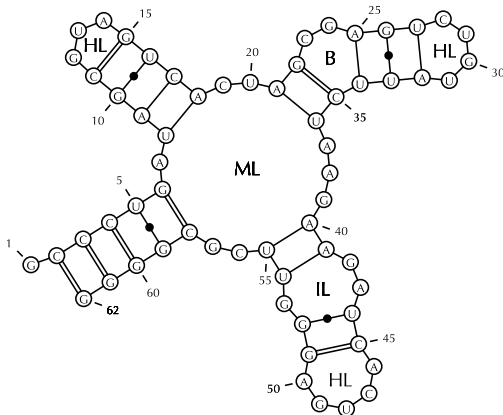
`{s_wild,nebel}@cs.uni-kl.de`



5. October 2012

10. Herbstseminar der Bioinformatik

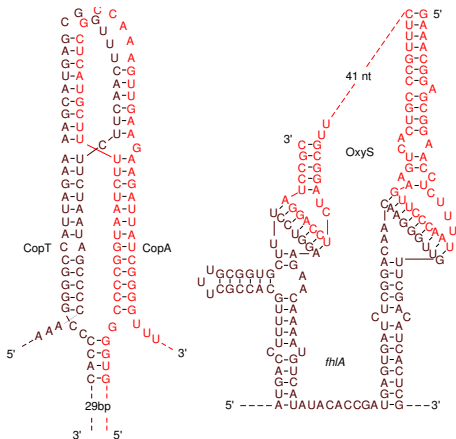
RNA secondary structure: model



- **primary structure:**
word over $\{a, c, g, u\}$
- **secondary structure:**
parentheses word $\{(, |,)\}$

gcccugauagcguagucacuagcgagucuguaauucuaagaagaucacugagggguucgcgggg
|(((((((|((((|))))))|((|(((|))))))|(((|(((|))))))|))))))

RNA also interacts!

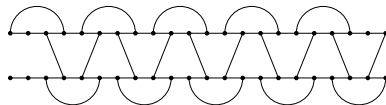


- bacterial antisense RNA
- interact non-trivially, “knotted” structure
- **goal:** predict whole interaction structure, not only interaction sites

RNA-RNA interaction problem (RIP)

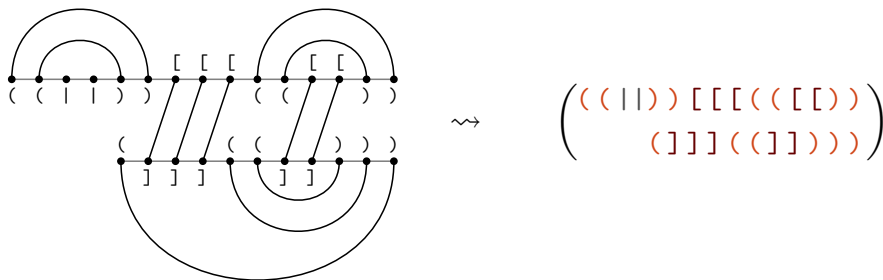
- **two** RNA molecules interact, i. e. form a **joint** secondary structure
- predict the “best” one possible
- RIP in general \mathcal{NP} -complete \rightsquigarrow restrict structure
 - 1 Exclude **pseudo knots** (internal & external)

- 2 Exclude **Zig-Zags**:



Example joint secondary structure

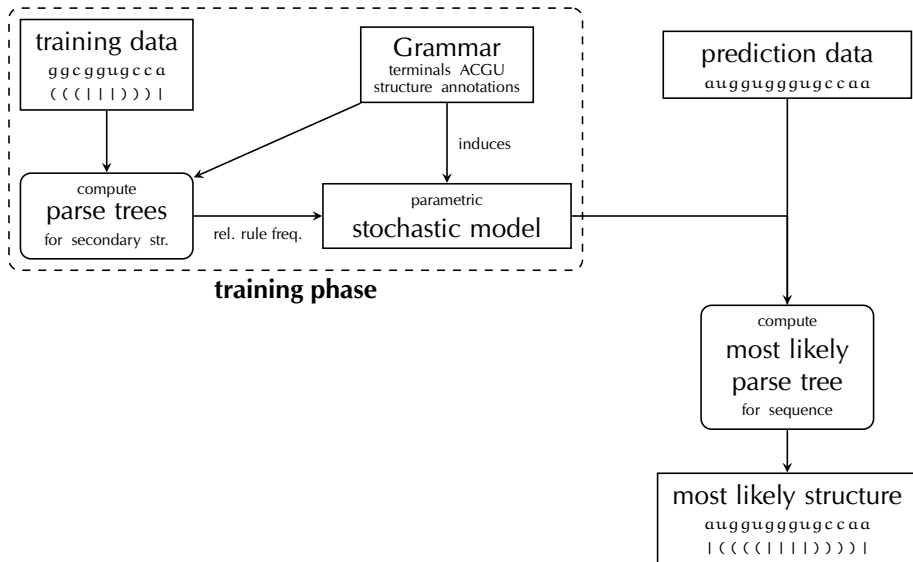
- \rightsquigarrow joint structure can be encoded as **“2D-word”**:
 - pair of upper and lower word
 - matching $()$ \rightsquigarrow internal bond
 - matching $[]$ \rightsquigarrow external bond (between two molecules)
 - $|$ \rightsquigarrow unpaired base



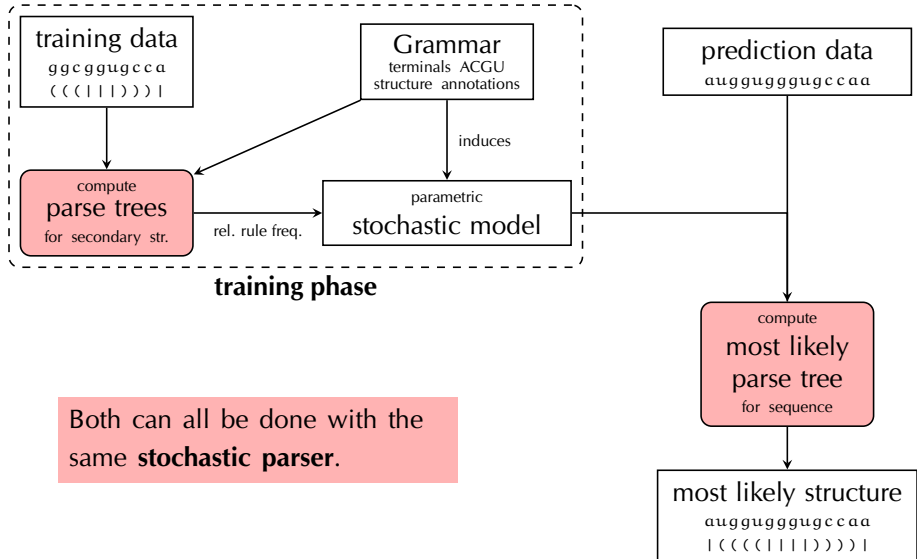
Stochastic Context-free Grammars

- $G = (N, \Sigma, R, S, P)$
 - nonterminals N
 - terminal alphabet Σ
 - rule set $R \subset N \times (N \cup \Sigma)^*$
 - start nonterminal $S \in N$
 - $P : R \rightarrow [0, 1]$: rule probabilities
- Probability of a **derivation tree**:
Product of used rules' probabilities
- For structure prediction:
 - $\Sigma = \{ |_{\text{b}}, (_{\text{b}},)_{\text{b}} : \text{b} \in \{\text{a}, \text{c}, \text{g}, \text{u}\} \}$
 - **unambiguous** w. r. t. structure

Structure Prediction with Formal Grammars



Structure Prediction with Formal Grammars



Both can all be done with the same **stochastic parser**.

Stochastic Context-free Grammars

- $G = (N, \Sigma, R, S, P)$
 - nonterminals N
 - terminal alphabet Σ
 - rule set $R \subset N \times (N \cup \Sigma)^*$
 - start nonterminal $S \in N$
 - $P : R \rightarrow [0, 1]$: rule probabilities

2-dimensional CFG

2-dimensional Context-free Grammars

- $G = (N, \Sigma, R, S, P)$
 - nonterminals N
 - terminal alphabet Σ
 - rule set $R \subset N \times (N \cup (\Sigma^*)^2)^*$
 - start nonterminal $S \in N$
 - $P : R \rightarrow [0, 1]$: rule probabilities

A Simplistic Grammar for RIP

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (|), \quad S \rightarrow (|) S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow ([]), \quad S \rightarrow ([]) S$$

Example:

S

A Simplistic Grammar for RIP

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (|), \quad S \rightarrow (|) S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow ([]), \quad S \rightarrow ([]) S$$

Example:

$(|)S$

A Simplistic Grammar for RIP

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (|), \quad S \rightarrow (|) S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow ([]), \quad S \rightarrow ([]) S$$

Example:

$$(|) (() S ()) S$$

A Simplistic Grammar for RIP

$$s \rightarrow (|), \quad s \rightarrow (|)s, \quad s \rightarrow (() s ()), \quad s \rightarrow (() s ())s$$

$$s \rightarrow (|), \quad s \rightarrow (|)s, \quad s \rightarrow (() s ()), \quad s \rightarrow (() s ())s$$

$$s \rightarrow ([]), \quad s \rightarrow ([])s$$

Example:

$$(|) (() (() s ()) ()) s$$

A Simplistic Grammar for RIP

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (|), \quad S \rightarrow (|) S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix}), \quad S \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix}) S$$

Example:

$$(|) (() (() (\begin{bmatrix} | \\ | \end{bmatrix}) S ()) ()) S$$

A Simplistic Grammar for RIP

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (|), \quad S \rightarrow (|) S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix}), \quad S \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix}) S$$

Example:

$$(|) (() (() (\begin{bmatrix} | \\ | \end{bmatrix}) (\begin{bmatrix} | \\ | \end{bmatrix}) S () ()) S$$

A Simplistic Grammar for RIP

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ())S$$

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ())S$$

$$S \rightarrow \left(\begin{matrix} [\\] \end{matrix} \right), \quad S \rightarrow \left(\begin{matrix} [\\] \end{matrix} \right) S$$

Example:

$$(|)(()((\left(\begin{matrix} [\\] \end{matrix} \right) \left(\begin{matrix} [\\] \end{matrix} \right) \left(\begin{matrix} [\\] \end{matrix} \right) () ())S$$

A Simplistic Grammar for RIP

$$s \rightarrow (|), \quad s \rightarrow (|)s, \quad s \rightarrow (() s ()), \quad s \rightarrow (() s ())s$$

$$s \rightarrow (|), \quad s \rightarrow (|)s, \quad s \rightarrow (() s ()), \quad s \rightarrow (() s ())s$$

$$s \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix}), \quad s \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix})s$$

Example:

$$(|) (() (() \begin{bmatrix} | \\ | \end{bmatrix}) \begin{bmatrix} | \\ | \end{bmatrix}) \begin{bmatrix} | \\ | \end{bmatrix}) () () (|) s$$

A Simplistic Grammar for RIP

$$s \rightarrow (|), \quad s \rightarrow (|)s, \quad s \rightarrow (() s ()), \quad s \rightarrow (() s ())s$$

$$s \rightarrow (|), \quad s \rightarrow (|)s, \quad s \rightarrow (() s ()), \quad s \rightarrow (() s ())s$$

$$s \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix}), \quad s \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix})s$$

Example:

$$(|) (() (() \begin{bmatrix} | \\ | \end{bmatrix}) \begin{bmatrix} | \\ | \end{bmatrix}) \begin{bmatrix} | \\ | \end{bmatrix}) () () (|) (|)$$

A Simplistic Grammar for RIP

$$s \rightarrow (|), \quad s \rightarrow (|)s, \quad s \rightarrow (() s ()), \quad s \rightarrow (() s ())s$$

$$s \rightarrow (|), \quad s \rightarrow (|)s, \quad s \rightarrow (() s ()), \quad s \rightarrow (() s ())s$$

$$s \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix}), \quad s \rightarrow (\begin{bmatrix} | \\ | \end{bmatrix})s$$

Example:

$$(|) (() (() \begin{bmatrix} | \\ | \end{bmatrix}) \begin{bmatrix} | \\ | \end{bmatrix}) \begin{bmatrix} | \\ | \end{bmatrix}) () () (|) (|)$$

A Simplistic Grammar for RIP

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow ([]), \quad S \rightarrow ([]) S$$

Example:

| (([[[]]])) | |

A Simplistic Grammar for RIP

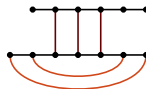
$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

$$S \rightarrow (|), \quad S \rightarrow (|)S, \quad S \rightarrow (() S ()), \quad S \rightarrow (() S ()) S$$

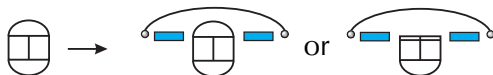
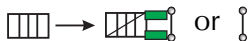
$$S \rightarrow ([]), \quad S \rightarrow ([]) S$$

Example:

$((([[[[]]]])))$



Used Grammar



Implementation

- stochastic parser, independent of grammar
- mixture of Earley-Parser and dynamic programming
- fast manual implementation in C++

RNA pair	n	m	runtime	memory
DIS DIS	35	35	2 min	300 MB
CopA CopT	56	57	1 h	2 GB
ompA MicA	137	72	2 d	18 GB
U2 and U6 snRNAs in yeast 21	144	95	1 week	34 GB

Summary

This Talk:

- 2D-CFGs give stochastic model for RNA joint structures (only slight extension of CFGs needed)
- Earley parsing can be used to train model and compute most likely structures
- efficient implementation available

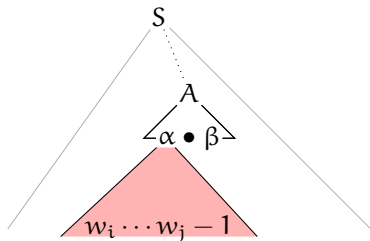
Open Problems:

- get good training data
- full empirical evaluation of prediction quality

Earley-Parsing

- does not need normal form
- here: as formal calculus
- defined in terms of **items**
 $(i \ j, A \rightarrow \alpha \bullet \beta)$ can be derived iff

$$S \Rightarrow^* w_{1,i-1} A \gamma \Rightarrow w_{1,i-1} \alpha \beta \gamma \Rightarrow^* w_{1,j-1} \beta \gamma$$



Earley-Parser for SCFGs

- Items: $(i \ j, A \rightarrow \alpha \bullet \beta)$

derivable **iff** $S' \Rightarrow^* w_{1,i-1} A \gamma \Rightarrow w_{1,i-1} \alpha \beta \gamma \Rightarrow^* w_{1,j-1} \beta \gamma$

- Start-Item: $(1 \ 1, S' \rightarrow \bullet S)$, Goal-Item $(1 \ n + 1, S' \rightarrow S \bullet)$
- Derivation Rules:

- Scanner

$$\frac{(i \ j - 1, A \rightarrow \alpha \bullet w_{j-1} \beta)}{(i \ j, A \rightarrow \alpha w_{j-1} \bullet \beta)}$$

- Predictor:

$$\frac{(i \ j, A \rightarrow \alpha \bullet B \gamma)}{(j \ j, B \rightarrow \bullet \beta)}$$

- Completer:

$$\frac{(i \ r, A \rightarrow \alpha \bullet B \gamma) \quad (r \ j, B \rightarrow \beta \bullet)}{(i \ j, A \rightarrow \alpha B \bullet \gamma)}$$

Earley-Parser for 2D-CFGs

- Items: $\left(\begin{smallmatrix} i & j \\ k & l \end{smallmatrix}, A \rightarrow \alpha_1 \bullet \beta_1 \atop \alpha_2 \bullet \beta_2 \right)$

derivable **iff** $S' \Rightarrow^* \begin{smallmatrix} u_{1,i-1} \\ v_{1,k-1} \end{smallmatrix} A \gamma \Rightarrow \begin{smallmatrix} u_{1,i-1} \\ v_{1,k-1} \end{smallmatrix} \alpha \beta \gamma \Rightarrow^* \begin{smallmatrix} u_{1,j-1} \\ v_{1,l-1} \end{smallmatrix} \beta \gamma$

- Start-Item: $\left(\begin{smallmatrix} 1 & 1 \\ 1 & 1 \end{smallmatrix}, S' \rightarrow \begin{smallmatrix} \bullet S_1 \\ \bullet S_2 \end{smallmatrix} \right)$, Goal-Item $\left(\begin{smallmatrix} 1 & n+1 \\ 1 & m+1 \end{smallmatrix}, S' \rightarrow \begin{smallmatrix} S_1 \bullet \\ S_2 \bullet \end{smallmatrix} \right)$

- Derivation Rules:

- Scanner (**upper**)
(lower similar):
$$\frac{\left(\begin{smallmatrix} i & j-1 \\ k & l \end{smallmatrix}, A \rightarrow \alpha_1 \bullet u_{j-1} \beta_1 \atop \alpha_2 \bullet \beta_2 \right)}{\left(\begin{smallmatrix} i & j \\ k & l \end{smallmatrix}, A \rightarrow \alpha_1 u_{j-1} \bullet \beta_1 \atop \alpha_2 \bullet \beta_2 \right)}$$

- Predictor:
$$\frac{\left(\begin{smallmatrix} i & j \\ k & l \end{smallmatrix}, A \rightarrow \alpha_1 \bullet B_1 \gamma_1 \atop \alpha_2 \bullet B_2 \gamma_2 \right)}{\left(\begin{smallmatrix} j & j \\ l & l \end{smallmatrix}, B \rightarrow \begin{smallmatrix} \bullet \beta_1 \\ \bullet \beta_2 \end{smallmatrix} \right)}$$

- Completer:
$$\frac{\left(\begin{smallmatrix} i & r \\ k & s \end{smallmatrix}, A \rightarrow \alpha_1 \bullet B_1 \gamma_1 \atop \alpha_2 \bullet B_2 \gamma_2 \right) \quad \left(\begin{smallmatrix} r & j \\ s & l \end{smallmatrix}, B \rightarrow \begin{smallmatrix} \beta_1 \bullet \\ \beta_2 \bullet \end{smallmatrix} \right)}{\left(\begin{smallmatrix} i & j \\ k & l \end{smallmatrix}, A \rightarrow \alpha_1 B_1 \bullet \gamma_1 \atop \alpha_2 B_2 \bullet \gamma_2 \right)}$$

Implementation Note: Dense Grammars

- Number of items: $\Theta(\frac{1}{4}n^2m^2)$
for n and m lengths of two RNA sequences
- structure prediction grammars are **dense**:
most items are derivable
- \rightsquigarrow compute values for **all** items à la dynamic programming