Algorithmic Aspects of RNA Secondary Structures

Stéphane Vialette

CNRS & LIGM, Université Paris-Est Marne-la-Vallée, France

2016-1017

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 1 / 138

200

イロト イポト イヨト イヨト 一日

Plan



Pseudoknot prediction and alternate models

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

Э 2016-2017 2 / 138

5900

イロト イポト イヨト イヨト

Computational biology

- Computational biology involves the development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, behavioral, and social systems.
 The field is broadly defined and includes foundations in computer science, applied mathematics, animation, statistics, biochemistry, chemistry, biophysics, molecular biology, genetics, genomics, ecology, evolution, anatomy, neuroscience, and visualization.
- Computational biology is different from biological computation, which is a subfield of computer science and computer engineering using bioengineering and biology to build computers, but is similar to bioinformatics, which is an interdisciplinary science using computers to store and process biological data.

Sar

イロト イポト イヨト イヨト 一日

Computational biology - Subfields

- Computational biomodeling
- Computational genomics (Computational genetics)
- Computational neuroscience
- Computational pharmacology
- Computational evolutionary biology
- Cancer computational biology
- ...

200

イロト イポト イヨト イヨト 一日

Computational biology - Major conferences

- Workshop on Algorithms in Bioinformatics (WABI)
- Asia Pacific Bioinformatics Conference (APBC)
- Intelligent Systems for Molecular Biology (ISMB)
- European Conference on Computational Biology (ECCB)
- Research in Computational Molecular Biology (RECOMB)

San

イロト イポト イヨト イヨト 二日

Computational biology – RECOMB satellites

- RECOMB Satellite Workshop on Massively Parallel Sequencing (RECOMB-SEQ)
- RECOMB Satellite Conference on Bioinformatics Education (RECOMB-BE)
- RECOMB Satellite Workshop on Computational Cancer Biology (RECOMB-CCB)
- RECOMB Satellite Workshop on Computational Methods in Genetics (RECOMB-Genetics)
- RECOMB Satellite Workshop on Comparative Genomics (RECOMB-CG)
- RECOMB Satellite Conference on Open Problems in Algorithmic Biology

<ロ> <同> <目> <日> <日> <日> <日> <日> <日> <日> <日> <日</p>

Computational biology - Main journals

There are numerous journals dedicated to computational biology:

- Journal of Computational Biology
- PLOS Computational Biology
- Bionformatics
- IEEE/ACM Transactions on Computational Biology and Bioinformatics
- BMC Bioinformatics
- Journal of Bioinformatics and Computational Biology

• ...

Preamble

Société Française de Bioinformatique (SFBI)



RNA Secondary Structures

2016-2017 8 / 138

Sac

Association des Jeunes Bioinformaticiens de France (JeBiF)



Fête de la Science 2016



Dans le cadre de la Fête de la Science, des ateliers et des conférences autour de la Bioinformatique seront présentés dans différentes villes de France.

Découvrez ces ateliers et conférences ci-dessous.

RNA Secondary Structures

PROCHAINS JEBIF PUBS

Tables Ouvertes Bioinfo Paris 13/10/2016 Le Barilleur - Paris JeBIF Pub Lyon 13/10/2016 Les Fleurs du Malt - Lyon

À PROPOS

Suivez-nous sur : Twitter : @JeBiF Facebook : JeBiF RSG-France

200

9 / 138

2016-2017

Preamble

http://bioinfo-fr.net/



Isabelle, Nolwenn, Yoann et Yohan les admins de bioinfo-fr.net

S share 42 G+1 0

🏥 ven 23 Sep 2016 🛔 Yoann M. 💊 Editorial 🙊 0

TEctto



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 10 / 138

GdR de Bioinformatique Moléculaire (GdR BIM – INS2I CNRS)



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 11 / 138

Carrières académiques

- Maître de Conférences / Professeur des Universités
- CNRS : Section 06, Section 07 et CID 51 (Modélisation, et analyse des données et des systèmes biologiques : approches informatiques, mathématiques et physiques)
- INRIA
- INRA
- CEA
- Institut Pasteur
- Institut Curie

Preamble

AlgoB :: Algorithmics for Bioinformatics :: Home

	AlgoB
Home Research Events	
Home > Welcome	Some More
Home	Welcome
Welcome	Our welcome message to our dear visitors. Read it
Welcome to the webpage of the AlgoB group of the <u>Laboratoire d'Informatique Gaspard-</u> Monee, Itself part of the <u>Institut Gaspard-Monee</u> of <u>Université Paris-Est Marne-la-Vallée</u> , in France.	• Members The composition of the AlgoB group.
AlgoB stands for Algorithmics for Bioinformatics. The group is composed of approximately <u>10</u> researchers (including 6 permanent researchers) and is part of a bigger team called <u>Models</u>	Read it
and algorithms.	Contact us
AlgoB	Our coordinates : mail address, telephone number, fax number and e- mail address. More details
Hare you will find information about	● IGM
Our <u>members</u>	The Gaspard-Monge institute of electronics and computer science at
Our <u>publications</u> and <u>talks</u>	Marne-La-Vallée Learn more
 The research projects, events and software developments we're involved in How to contact us 	O VIDENA
The Laboratoire d'Informatique Gaspard-Monge (LIMP 8049)	• OPEM
The Institut Gaspard-Monge	The University of Paris-Est Marne-La- Vallée. Learn more

- The University of Paris-Est Marne-la-Vallée

- 2006-2015 ABligoo -

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

Э 2016-2017 13 / 138

200

Preamble

AlgoB :: Algorithmics for Bioinformatics :: Home



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

Э 2016-2017 14 / 138

5900

イロト イロト イヨト イヨト

Plan





3 RNA secondary structure prediction

Pseudoknot prediction and alternate models

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 15 / 138

Э

5900

イロト イポト イヨト イヨト

Timeline



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

Ξ 2016-2017 16 / 138

DQC

<ロト < 同ト < 三ト < 三ト

Central dogma of molecular biology



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 17 / 138

Э

DQC

イロト イロト イヨト イヨト

Central dogma of molecular biology

- 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 back from protein to either protein or nucleic acid.
- This has also been described as *DNA makes RNA makes protein*. However, this simplification does not make it clear that the central dogma as stated by Crick does not preclude the reverse flow of information

San

イロト イポト イヨト イヨト 一日

Transcription

- **Transcription** is the process by which the information contained in a section of DNA is transferred to a newly assembled piece of messenger RNA (mRNA).
- It is facilitated by RNA polymerase and transcription factors.
- In eukaryotic cells the primary transcript (pre-mRNA) must be processed further in order to ensure translation.
- This normally includes a 5' cap, a poly-A tail and splicing.
- Alternative splicing can also occur, which contributes to the diversity of proteins any single mRNA can produce.

200

イロト イポト イヨト イヨト 二日

Transcription



RNA Secondary Structures

2016-2017 20 / 138

Э

990

イロト イポト イヨト イヨト

Alternative splicing



RNA Secondary Structures

Ξ 2016-2017 21 / 138

500

<ロト < 同ト < 三ト < 三ト

Translation

- Eventually, the mature mRNA finds its way to a ribosome, where it is **translated**.
- In prokaryotic cells, which have no nuclear compartment, the process of transcription and translation may be linked together. In eukaryotic cells, the site of transcription (the cell nucleus) is usually separated from in the site of translation (the cytoplasm), so the mRNA must be transported out of the nucleus into the cytoplasm, where it can be bound by ribosomes
- The mRNA is read by the ribosome as triplet codons, usually beginning with an AUG (adenine-uracil-guanine), or initiator methionine codon downstream of the ribosome binding site.
- Translation ends with a *UAA*, *UGA*, or *UAG* stop codon.

S. Vialette (CNRS & LIGM)

San

イロト イポト イヨト イヨト 一日

Translation



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

1 2016-2017 23 / 138

Simultaneous translation and transcription

ribosomes



0.5 μm

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 24 / 138

Base pairing

- In molecular biology, two nucleotides on opposite complementary DNA or RNA strands that are connected via hydrogen bonds are called a base pair (often abbreviated bp).
- In the canonical Watson-Crick base pairing, adenine (*A*) forms a base pair with thymine (*T*), and guanine (*G*) forms one with cytosine (*C*) in DNA.
- In RNA, thymine is replaced by uracil (*U*).
- Alternate hydrogen bonding patterns, such as the wobble base pair and Hoogsteen base pair, also occur—particularly in RNA—giving rise to complex and functional tertiary structures.
- Importantly, pairing is the mechanism by which codons on messenger RNA molecules are recognized by anticodons on transfer RNA during protein translation

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 25 / 138

Base pairing



Left, an *AT* base pair demonstrating two intermolecular hydrogen bonds; Right, a *GC* base pair demonstrating three intermolecular hydrogen bonds.

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 26 / 138

Base pairing



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

Э 2016-2017 27 / 138

miRNA



Э 2016-2017 28 / 138

Non-coding RNA



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

3 2016-2017 29 / 138

200

イロト イポト イヨト イヨト

Structural conformations of biomolecules

- **Primary Structure**: sequence of monomeres (*ATCGAGATC*...)
- Secondary Structure: 2D-fold, defined by hydrogen bonds
- Tertiary Structure: 3D-fold
- **Quarternary Structure**: complex arrangement of multiple folded moleculesRNA tertiary structure



S. Vialette (CNRS & LIGM)

2016-2017 30 / 138

RNA seconday structure



The major role of tRNA is to translate mRNA sequence into amino acid sequence. A tRNA molecule consists of 70 - 80 nucleotides.

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 31 / 138

I D > I A

RNA tertiary structure



A hairpin loop from a pre-mRNA. Highlighted are the nucleobases (green) and the ribose-phosphate backbone (blue). Note that this is a single strand of RNA that folds back upon itself.

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 32 / 138

500

RNA tertiary structure



Three-dimensional representation of the 50S ribosomal subunit. RNA is in ochre, protein in blue. The active site is in the middle (red).

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 33 / 138

200

イロト イポト イヨト イヨト

Prediction of secondary structure: FASTA format

- **FASTA format** is a text-based format for representing either nucleotide sequences or peptide sequences, in which nucleotides or amino acids are represented using single-letter codes.
- The format also allows for sequence names and comments to precede the sequences.
- The format originates from the FASTA software package, but has now become a standard in the field of bioinformatics.
- The simplicity of FASTA format makes it easy to manipulate and parse sequences using text-processing tools and scripting languages like Python, Ruby, and Perl.

・ロト < 団ト < 三ト < 三ト < 三ト < 〇へ〇

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 34 / 138

Digression: BioXXX projects

- BioPython: http://biopython.org/wiki/Main_Page
- BioPerl: http://www.bioperl.org/wiki/Main_Page
- BioJava: http://biojava.org/wiki/Main_Page
- BioRuby: http://bioruby.org
- Bio (Haskell): http://biohaskell.org/Libraries/Bio
- BioCaml: http://biocaml.org

200

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Digression: BioPython



- The Biopython Project is an international association of developers of non-commercial Python tools for computational molecular biology, as well as bioinformatics.
- BioPython is one of a number of Bio* projects designed to reduce code duplication.
- http://biopython.org/wiki/Main_Page

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 36 / 138
Digression: BioPython

The main function is Bio.SeqIO.parse() which takes a file handle and format name, and returns a SeqRecord iterator.

```
from Bio import SeqIO
handle = open("example.fasta", "rU")
for record in SeqIO.parse(handle, "fasta") :
    print record.id
handle.close()
```

Digression: BioPython

```
from Bio import SeqIO
handle = open("example.fasta", "rU")
records = list(SeqIO.parse(handle, "fasta"))
handle.close()
print records[0].id #first record
print records[-1].id #last record
```

```
from Bio import SeqIO
handle = open("example.fasta", "rU")
record_dict = SeqIO.to_dict(SeqIO.parse(handle, "fasta"))
handle.close()
print record_dict["gi:12345678"] #use any record ID
```

<ロト < 同ト < 三ト < 三ト < 三ト 三 - のへ()

Introduction

Prediction of secondary structure: RNAfold





RNA Secondary Structures

2016-2017 39 / 138

Introduction

Prediction of secondary structure: RNAfold





RNA Secondary Structures

Prediction of secondary structure: RNAfold

RNAfold WebServer		2 View Result	
	(Home)	New job Hel	p]
Results for minimum free energy prediction			
The optimal secondary structure in dox-bracket notation with a minimum free energy of -28.80 kcal/mol is given below. [color by base-pairing probability [color by positional entropy] no coloring]			
1 00000EXXXX02CCXX0000ECCCCCCCCCCCCCCCCCCC			
1 IIII(((())))-III((1))))-III((1))))))))))			
You can download the minimum free energy (MFE) structure in [Vienna Format] Ct. Format]. You can get thermodynamic details on this structure by submitting to our RNAeval web server.			
Results for thermodynamic ensemble prediction			
The free energy of the thermodynamic ensemble is -30.14 kcal/mol.			
The frequency of the MFE structure in the ensemble is \$7.55 %.			
The ensemble diversity is 2.35.			
You may look at the dot plot containing the base pair processings [proj plater Conventent].			
The centroid secondary structure in dot-bracket notation with a minimum free energy of -29.80 kcal/mol is given below.			
[color by base-pairing probability color by positional entropy no coloring]			
1 GODDERNINGETENDUDGERNAMEGEOCOCCUTUEGENAARGUUMGEGOUMARGUUMGEGOUTENNEETEN			
1 (((((,,.)))),((((,,))))))))))))))			
You can download the minimum free energy (MFE) structure in {vienna Format]. You can get thermodynamic details on this structure by submitting to our RNAeval web server.			

http://rna.tbi.univie.ac.at/cgi-bin/RNAfold.cgi

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 40 / 138

トイヨトイヨト

5900

I D > I A

Prediction of secondary structure: RNAfold

Graphical output

You may look at the interactive drawing of the MET structure below. If you can see the interactive drawing and you are using interact Explorer, please install the Adobe SVC plugin. A note on base-pairing probabilities: The structure below is colored by base-pairing probabilities. The for unaived network the cord drawtest the the cord drawtest the the cord drawtest the the cord drawtest the the cord drawtest plane.



http://rna.tbi.univie.ac.at/cgi-bin/RNAfold.cgi

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 41 / 138

Plan





③ RNA secondary structure prediction

4 Pseudoknot prediction and alternate models

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 42 / 138

Э

5900

イロト イロト イヨト イヨト

RNA secondary structure prediction

- Many plausible secondary structures can be drawn from a sequence.
- The number increases exponentially with sequence length.
- An RNA only 200 bases long has over 10⁵⁰ possible base-paired structures.
- We must distinguish the biologically correct structure from all the incorrect.structures

200

イロト イポト イヨト イヨト 二日

Base pair maximisation: the Nussinov folding algorithm

- One (**naive**) approach is to find the structure with the most base pairs.
- Nussinov introduced an efficient dynamic programming algorithm for this problem.
- Although the criterion is too simplistic to give accurate structure predictions, the algorithm is instructive because the mechanics of the Nussinov folding algorithm are the same as those in the more sophisticated energy minimisation folding algorithms (and of probabilistic SCFG-based algorithms).

<ロト < 同ト < 三ト < 三ト < 三ト < 三 の < ○</p>

RNA secondary structure

Definition

Let $u \in \{A, C, G, U\}^*$ be a sequence. An **RNA-structure** over *u* is a set of pairs

 $P = \{(i, j) : i < j, u[i] \text{ and } u[j] \text{ form a a WC or non-standard pair} \}$

with the property that the associated graph has degree at most 1 (*i.e.*, every base can have at most one bond).

Remark

$$\begin{array}{ll} \forall (i,j), & (i,j) \in P \ \Rightarrow \forall i', (i',j) \notin P \\ \forall (i,j), & (i,j) \in P \ \Rightarrow \forall j', (i,j') \notin P \end{array}$$

500

<ロト < 回 > < 回 > < 回 > < 回 >

RNA secondary structure



Purine riboswitch (Rfam RF00167)

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 46 / 138

500

< D > < 🗗

The Nussinov folding algorithm

• **Idea (biological)**: Stacked base pairs of helical regions are considered to stabilize an RNA molecule.

Therefore, the goal is to maximize the number of base pairs.

- **Idea (algorithmic)**: The optimal structure *S*[*i*, *j*] on a subsequence u[i, j] can only be formed by two distinct ways from a shorter subsequence u[i + 1, j]:
 - Base *i* is unpaired, followed by an arbitrary shorter structure.
 - 2 Base *i* is paired with some partner base *k* requiring the computation of two independent substructures: the structure enclosed by the bp and the remaining structure behind the pair.



The Nussinov folding algorithm

Initialisation 0

$$\begin{aligned} \gamma(i, i-1) &= 0 \qquad 2 \leq i \leq n \\ \gamma(i, i) &= 0 \qquad 1 \leq i \leq n \end{aligned}$$

Recursion

$$\begin{split} \gamma(i,j) &= \max \begin{cases} \gamma(i+1,j) \\ \gamma(i,j-1) \\ \gamma(i+1,j-1) + \alpha(i,j) \\ \max_{i < k < j} \left\{ \gamma(i,k) + \gamma(k+1,j) \right\} \end{cases} \end{split}$$

• $O(n^3)$ time and $O(n^2)$ space.

Э 2016-2017 48 / 138

5900

イロト イロト イヨト イヨト

RNA secondary structure prediction

The Nussinov folding algorithm



2016-2017 49 / 138

=

500

< D > < 🗗



Э 2016-2017 50 / 138

200

ト イヨト イヨト

< □ ト < @



S. Vialette (CNRS & LIGM)

Э 2016-2017 50 / 138

200

ト イヨト イヨト

I D > I A



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 50 / 138

Э

ト イヨト イヨト

I D > I A

200



S. Vialette (CNRS & LIGM)

Э 2016-2017 50 / 138

ト イヨト イヨト

I D > I A

200

The Nussinov folding algorithm: Example A U С G G A A С G G G G А А Α U С С

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 50 / 138

Э

ト イヨト イヨト

< D > < P



S. Vialette (CNRS & LIGM)

Э 2016-2017 50 / 138

200

ト イヨト イヨト

I D > I A



S. Vialette (CNRS & LIGM)

Э 2016-2017 50 / 138

5900

3

- 4

< D > < P

▶ < Ξ >

The Nussinov folding algorithm

- The value *γ*(1, *n*) is the number of base pairs in the maximally base-paired structure.
- There are often a number of alternatives structures with the same number of base-pairs.
- To find one of these maximally base-paired structures, we trace back through the values we calculated in the dynamic programming matrix, beginning from *γ*(1, *n*).

San

イロト イポト イヨト イヨト 一日

The Nussinov folding algorithm: Traceback stage

```
1 S \leftarrow emptyStack
 2 push(S, (1, n))
 <sup>3</sup> while S \neq \emptyset do
           (i, j) \leftarrow \mathsf{pop}(S)
 4
          if i < j then
 5
                 if \gamma(i+1,j) = \gamma(i,j) then
 6
                       push(S, (i+1, j))
 7
                 else if \gamma(i, j-1) = \gamma(i, j) then
 8
                       push(S,(i,j-1))
 9
                 else if \gamma(i+1, j-1) + \alpha(i, j) = \gamma(i, j) then
10
                       record pair (i, j)
11
                       push(S, (i + 1, i - 1))
12
                 else
13
                       for j + 1 \le k \le j - 1 do
14
                             if \gamma(i,k) + \gamma(k+1,j) = \gamma(i,j) then
15
                              push(S, (k+1,j))
push(S, (i,k))
break
16
17
18
```

2016-2017 52 / 138

イロト イポト イヨト イヨト 二日

RNA secondary structure prediction

The Nussinov folding algorithm: Traceback stage





Э 2016-2017 53 / 138

500

3

- 4

▶ < Ξ >

< D > < A

The Nussinov folding algorithm: Traceback stage

- Linear time and space.
- In the canonical implementation of the traceback step, whenever there are multiple structures that are equivalent in terms of number of base-pairs the first structure that works is chosen because the algorithm does not care about anything besides the number of base-pairs, so any structure with the same number of base pairs as the optimal one will do.
- However, this ignores important information that can lead it to choose an unstable structure over a more stable one.

San

イロト イポト イヨト イヨト 一日

The Nussinov folding algorithm: Traceback stage

Possible improvements:

- Incorporates the following assumptions:
 - Longer stems (consecutive base pairs) are more stable than shorter stems,
 - A single loop or bulge is more stable than one split in two by a base pair in the middle
- Report all optimal solutions.

San

・ロト (四) (三) (三)

The Nussinov folding algorithm: Drawbacks

The Nussinov folding algorithm does not determine biological relevant structures since:

- There are many (all!?) possibilities to form base pairs.
- Stackings of base pairs are not considered.
- The size of internals loops are not considered.

200

イロト イポト イヨト イヨト 一日

An SCFG version of te Nussinov folding algorithm

- A single non-terminal *S*;
- 14 production rules with associated probability parameters.

$$S \rightarrow aS \mid cS \mid gS \mid uS \qquad (i \text{ unpaired})$$

$$S \rightarrow Sa \mid Sc \mid Sg \mid Su \qquad (j \text{ unpaired})$$

$$S \rightarrow aSu \mid cSg \mid gSc \mid uSa \qquad (i \text{ and } j \text{ paired})$$

$$S \rightarrow SS \qquad (bifurcation)$$

$$S \rightarrow \epsilon \qquad (termination)$$

<ロト < 同ト < 巨ト < 巨ト -

An SCFG version of te Nussinov folding algorithm

- Assume that the probabiliy parameters are known.
- The maximum probability parse of a sequence with this SCFG is an assignment of sequence positions to productions.
- Because the productions correspond to secondary structure elements (base pairs and single-stranded bases), the maximum probability parse is equivalent to the maximum probability secondary structure.
- If base pair productions have relatively high probability, the SCFG will favour parses which tend to maximise the number of base pairs in the structure.

San

イロト イポト イヨト イヨト 二日

RNA structure prediction: MFE-folding

- More realistic: thermodynamics and statistical mechanics.
- Stability of an RNA secondary structure coincides with thermodynamic stability.
- Quantified as the amount of free energy released/used by forming base pairs.

200

イロト イポト イヨト イヨト 二日

RNA structure prediction: MFE-folding

RNA molecules basically exist in a distribution of structures rather than a single ground-state conformation.

- "Most likely" conformation: structure exhibiting minimum of free energy (MFE).
- Energy contributions of different loop types have been measured.
- Since free energies are additive, a more sophisticated model, the standard energy model for RNA secondary structures, can be proposed.
- Based on loop decomposition, the total energy *E* of a structure *S* can be computed as the sum over the energy contributions of each constituent loop *ℓ*:

$$E(S) = \sum_{\ell \in S} E(\ell)$$

MFE folding: Example



2016-2017 61 / 138

Structural elements: Formal definition

Definition

Secondary structure elements Let u be a fixed sequence. Further, let S be an RNA secondary structure for u.

• A base pair $(i, j) \in S$ is a **hairpin loop** if

$$\forall i < i' \leq j' < i, \quad (i', j') \notin S$$

- A base pair $(i, j) \in S$ is called **stacking** if $(i + 1, j 1) \in P$.
- Two base pairs (i, j) ∈ S and (i', j') ∈ S form an internal loop (i, j, i', j') if
 i < i' < j' < j
 (i' i) + (j j') > 2 (no stack)
 there is no base pair (k, l) between (i, j) and (i', j')

500

Structural elements: Formal definition

Definition

Secondary structure elements (ctd)

- An internal loop (i, j, i', j') is called **left** (resp. **right**) **bulge**, if j = j' + 1 (resp. i' = i + 1).
- A *k*-multiloop consists of *k* base pairs $(i_1, j_1), (i_2, j_2), \ldots, (i_k, j_k) \in S$ and a closing base pair $(i, j) \in S$ with the property that

Structural elements: Formal definition



2016-2017 64 / 138

=

DQC

Structural elements: *k*-multiloop

Remarks

- Usually hairpin loops have minimal loop size of m = 3 (*i.e.*, for all $(i, j) \in S, i < j 3$).
- Each secondary structure element is defined uniquely by its closing basepair.
- For any basepair (*i*, *j*) we denote the corresponding secondary structure element with *sec*(*i*, *j*).


Energy of secondary structural elements

Definition (Energy contribution of loops) Energy contributions of the various structure elements:

hairpin loop (i,j):eH(i,j)stacking (i,j):eS(i,j)internal loop (i,j,i,j'):eL(i,j,i',j')multiloop $(i,j,i_1,j_2,\ldots,i_k,j_k)$: $eM(i,j,i_1,j_2,\ldots,i_k,j_k)$

Remarks

- General multi loop contribution will be too expensive in prediction: exponential explosion!
- Use a simplified contribution scheme.

200

イロト イロト イヨト イヨト

Energy of secondary structural elements

Definition (Simplified energy contribution of multiloops)

$$\mathbf{eM}(i,j,k,k') = a + bk + ck'$$

where

- *a*, *b* and *c* are weights (*a* is the energy contribution for closing of loop),
- *k* is the number of inner base pairs, and
- *k*′ is the number of unpaired bases within loop.

医子宫下子宫下 二

MFE folding

- The complexity of the dynamic programming algorithm is $O(n^4)$ time and $O(n^2)$ space.
- Using a trick, the time complexity can be reduced to $O(n^3)$.
- We assume traceback is done analogously to Nussinov-Traceback. Same reduced complexity. Only extension: trace through three matrices, (*i.e.*, keep track of matrix).

200

<ロト < 課 ト < 注 ト < 注 ト - 注

Plan



Pseudoknot prediction and alternate models

200

<ロト < 同ト < 三ト < 三ト

RNA pseudoknots

- A **pseudoknot** is a nucleic acid secondary structure containing at least two stem-loop structures in which half of one stem is intercalated between the two halves of another stem.
- The pseudoknot was first recognized in the turnip yellow mosaic virus in 1982.
- Pseudoknots fold into knot-shaped three-dimensional conformations but are not true topological knots.

San

イロト イポト イヨト イヨト

RNA pseudoknots



S. Vialette (CNRS & LIGM)

RNA Secondary Structures

Э 2016-2017 71 / 138

=

900

RNA pseudoknots



2016-2017 72 / 138

RNA pseudoknots: Prediction and identification

- The structural configuration of pseudoknots does not lend itself well to bio-computational detection due to its context-sensitivity or "overlapping" nature.
- The presence of pseudoknots in RNA sequences is more difficult to predict by the standard method of dynamic programming, which use a recursive scoring system to identify paired stems and consequently, most cannot detect non-nested base pairs.
- Popular secondary structure prediction methods do not predict pseudoknot structures present in a query sequence.
- It is possible to identify a limited class of pseudoknots using dynamic programming, but these methods are not exhaustive and scale worse as a function of sequence length than non-pseudoknotted algorithms.
- The general problem of predicting lowest free energy structures with pseudoknots has been shown to be **NP**-complete.

S. Vialette (CNRS & LIGM)

2016-2017 73 / 138

RNA pseudoknots: Biological significance

- Several important biological processes rely on RNA molecules that form pseudoknots, which are often RNAs with extensive tertiary structure.
- For example, the pseudoknot region of RNase P is one of the most conserved elements in all of evolution.
- The telomerase RNA component contains a pseudoknot that is critical for activity.
- Several viruses use a pseudoknot structure to form a tRNA-like motif to infiltrate the host cell.

San

イロト イポト イヨト イヨト 二日

RNA pseudoknot type

• Simple, H-type



< D > < 🗗

Э 2016-2017 75 / 138

5900

1

< ∃ >

RNA pseudoknot type

Kissing hairpin



=

200

< □ ト < @

▶ < Ξ.</p>

RNA pseudoknot type

Three-knot



DQC

1

<ロト < 同ト < 三ト

RNA pseudoknot predition



Э 2016-2017 78 / 138

5900

ト イヨト イヨト

< □ ト < @

Nearest Neighbor Energy Model

For a secondary structure *S*

• the number of **base pairs stackings** is

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

• the number of **stacking base pairs** is

$$SBP(S) = |\{(i,j) \in S : (i+1,j-1) \in S \text{ or } (i-1,j+1) \in S\}|$$

200

I D > I A

Without pseudoknots

• Maximizing the number of base pairs is $O(n^3)$ time and $O(n^2)$ space:

$$\gamma(i,j) = \max \begin{cases} \gamma(i+1,j) \\ \gamma(i,j-1) \\ \gamma(i+1,j-1) + \alpha(i,j) \\ \max_{i < k < j} \{\gamma(i,k) + \gamma(k+1,j)\} \end{cases}$$

• To maximize BPS or SBP, dynamic programming can be extended.

5900

イロト イロト イヨト イヨト

With arbitrary pseudoknots

• Maximizing the number of base pairs is $O(n^{2/5})$ time.

The problem reduces to finding a maximum matching in a graph (solvable in $O(m\sqrt{n})$ time)

• To maximize BPS or SBP, matching becomes inadequate, and dynamic programming cannot be extended.

200

イロト イポト イヨト イヨト 二日

Theorem

Given a sequence u and a positive integer K, it is **NP**-complete to decide whether there exists a structure S that is legal under the canonical RNA folding model and with BPS $(S) \ge k$.

▶ < 프 ▶ < 프 ▶</p>

BIN PACKING

Definition (The BIN PACKING problem)

- **Input**: *k* items of sizes *a*₁, *a*₂, . . . , *a*_{*k*} and *B* bins each with capacity *C*.
- **Question**: Decide whether the items fit into the bins.

Remarks

- Strongly NP-complete.
- A straightforward greedy algorithm achieves an approximation factor of 2.
- Does not have a polynomial-time approximation scheme (PTAS) unless P = NP.
- for any $0 < \epsilon \le 1$, it is possible to find a solution using at most $(1 + \epsilon)$ **opt** + 1 bins in polynomial time (asymptotic PTAS).

5900

イロト イポト イヨト イヨト 二日

Strongly NP-complete

Definition

A problem is said to be **strongly NP-complete** (or **NP-complete in the strong sense**), if it remains so even when all of its numerical parameters are bounded by a polynomial in the length of the input.

- Normally numerical parameters to a problem are given in binary notation, so a problem of input size *n* might contain parameters whose size is exponential in *n*.
- If we redefine the problem to have the parameters given in unary notation, then the parameters must be bounded by the input size.
- Thus strong NP-completeness or NP-hardness may also be defined as the NP-completeness or NP-hardness of this unary version of the problem.

S. Vialette (CNRS & LIGM)

5900

Strongly NP-complete

Definition

A problem is said to be **strongly NP-complete** (or **NP-complete in the strong sense**), if it remains so even when all of its numerical parameters are bounded by a polynomial in the length of the input.

■ BIN PACKING is strongly **NP**-complete while the 0 − 1 KNAPSACK problem is only weakly **NP**-complete.

Thus the version of BIN PACKING where the object and bin sizes are integers bounded by a polynomial remains **NP**-complete, while the corresponding version of the KNAPSACK problem can be solved in polynomial time by dynamic programming.

• Any strongly **NP**-hard optimization problem with a polynomially bounded objective function cannot have an **FPTAS** unless P = NP.

S. Vialette (CNRS & LIGM)

Theorem

Given a sequence u and a positive integer K, it is **NP**-complete to decide whether there exists a structure S that is legal under the canonical RNA folding model and with BPS $(S) \ge k$.

Proof.

Construction:

$$u = C^{a_1} A C^{a_2} A \dots A C^{a_k} A A (A G^{\mathcal{C}})^B$$
$$K = \sum_{1 \le i \le k} a_i - k$$

As *A*'s can only form base pairs with *U*'s in the canonical folding model, all base pairs in a legal structure for u wil be $G \cdot C$.

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 85 / 138

< D > < B > < E > < E >

Theorem

Given a sequence u and a positive integer K, it is **NP**-complete to decide whether there exists a structure S that is legal under the canonical RNA folding model and with BPS $(S) \ge k$.

Proof.

•
$$|u| = \sum_{1 \le i \le k} a_i + B\mathcal{C} + B + k + 1.$$

• Since the BIN PACKING problem is **strongly NP**-complete we can assume that the $B, C, a_1, a_2, ..., a_k$ are all polynomially bounded by the size of the originally BIN PACKING instance.

イロト イロト イヨト イヨト

Theorem

Given a sequence $u \in \{0,1\}^*$ and a positive integer K, it is **NP**-complete to decide whether there exists a structure S that is legal under the general RNA folding model with $\mathbb{B} = \{(0,1), (1,0)\}$ with BPS $(S) \ge k$.

Proof.

Reduction from BIN PACKING:

▶
$$3 \le a_i \le C$$
 for $1 \le i \le k$, and
▶ $2 \le B \le k$.

• $u = 0^{a_1} 110^{a_2} 11 \dots 110^{a_k} (01^C)^B$.

•
$$K = \sum_{1 \le i \le k} a_i - k + B.$$

< D > < A

Theorem

Given an alphabet Σ , a set of legal base pairs $\mathbb{B} \subseteq \Sigma \times \Sigma$, a sequence $u \in \Sigma^*$ and a positive integer K, it is **NP**-complete to decide whether there exists a structure S that is legal under the general RNA folding model with SBP(S) $\geq k$.

Proof.

• Reduction from RESTRICTED SATISFIABILITY:



2-intervals

Definition (2-intervals)

- A 2-interval D = (I, J) consists of two disjoint (closed) intervals I and J such that I < J (*i.e.*, I is completely on the left of J).
- Two 2-intervals $D_1 = (I_1, J_1)$ and $D_2 = (I_2, J_2)$ are **disjoint** if the four intervals I_1, J_1, I_2 and J_2 are pairwise disjoint.

Definition (Relations between disjoint 2-intervals)

Let $D_1 = (I_1, J_1)$ and $D_2 = (I_2, J_2)$ be two disjoint 2-intervals. We have the following relations:

- $D_1 < D_2$ if $I_1 < J_1 < I_2 < J_2$.
- $D_1 \sqsubset D_2$ if $I_2 < I_1 < J_1 < J_2$.
- $D_1 \ \emptyset \ D_2$ if $I_1 < I_2 < J_1 < J_2$.

2-intervals: Models

Definition (Models)

- A **model** is a non-empty subset of $\{<, \sqsubset, \emptyset\}$.
- A set of 2-intervals \mathcal{D} is \mathcal{R} -structured if any two distinct 2-intervals in \mathcal{D} is R-comparable for some $R \in \mathcal{R}$.

Key idea: Model $\mathcal{R} = \{<, \sqsubset\}$ denotes pseudoknot-free structures.

<ロト < 課 > < 注 > < 注 > 二 注

2-intervals: Structured subsets

Definition (The 2-INTERVAL PATTERN problem)

- **Input**: A set of 2-intervals \mathcal{D} , a model \mathcal{R} and a positive integer *K*.
- **Question**: Decide whether there exists a \mathcal{R} -structured subset $\mathcal{D}' \subseteq \mathcal{D}$ of size *K*.

If each 2-interval $D \in \mathcal{D}$ is associated with a non-negative weight w(D), we are left with the WEIGHTED 2-INTERVAL PATTERN problem (*i.e.*, decide whether there exist a \mathcal{R} -structured subset $\mathcal{D}' \subseteq \mathcal{D}$ of total weight at least K).

2-intervals: Support and restriction

Definition (Support)

The **support** of a set of 2-intervals \mathcal{D} is the set of intervals $\{I, J : (I, J) \in \mathcal{D}\}.$

Definition (Restriction)

Support restrictions:

- Unlimited: no restriction.
- **Balanced**: every 2-intervals consists of two intervals of equal length.
- Unit: every 2-intervals consists of two intervals of unit length.
- **Point**: the intervals in the support are pairwise disjoint.

San

<ロト < 回 > < 回 > < 回 > < 回 >

2-intervals: State-of-the-art

Model	Unlimited	Balanced	Unit	Point		
$\overline{\{<,\sqsubset,\emptyset\}}$	APX-hard			$O(n\sqrt{n})$		
$\{\sqsubset, \emptyset\}$	APX-hard			$O(n\log(n) + \mathcal{L})$		
$\{<, \emptyset\}$	NP-complete					
{<,⊏}	$O(n\log(n) + dn)$					
{Ø}	$O(n\log(n) + \mathcal{L})$					
$\{<\}$	$O(n\log(n))$					
$\{\Box\}$	$O(n\log(n))$					

200

イロト イロト イヨト イヨト

2-intervals: Approximation ratios

Model	Unlimited	Balanced	Unit	Point		
$\{<,\sqsubset,\Diamond\}$	4	4	$2+\epsilon$	N/A		
$\{\sqsubset, \emptyset\}$	4	4	$2 + \epsilon$	N/A		
$\{<, \emptyset\}$	PTAS					

200

2-intervals: Model $\mathcal{R} = \{<, \sqsubset, \emptyset\}$

Theorem

The 2-INTERVAL PATTERN *problem for unlimited (resp. unit)* 2*-intervals* and model $\mathcal{R} = \{<, \sqsubset, \emptyset\}$ is approximable within ratio 4 (resp. 3.).

Remarks

- The approximation algorithm for unit 2-intervals is $O(n \log(n))$ time, where *n* is the number of input 2-intervals.
- The approximation algorithm for unlimited 2-interval uses linear programming techniques, which in practice are very often too time costly.

200

イロト イポト イヨト イヨト 一日

2-intervals: Model $\mathcal{R} = \{<, \sqsubset, \emptyset\}$

Theorem

The 2-INTERVAL PATTERN *problem for balanced* 2*-intervals and model* $\mathcal{R} = \{<, \sqsubset, \emptyset\}$ *is approximable within ratio* 4 ($O(n^2)$ *time algorithm*).

Proof.

 $\begin{array}{c|c} \textbf{Data: A set of balanced 2-intervals } \mathcal{D} \\ \textbf{Result: A } \{<, \sqsubset, \check{\mathbb{Q}}\}\text{-structured subset of } \mathcal{D} \\ 1 \quad \mathcal{D}_{sol} \leftarrow \varnothing \\ 2 \quad \textbf{while } \mathcal{D} \neq \oslash \textbf{do} \\ 3 \quad & \text{Let } D_{min} \text{ be the smallest 2-interval left in } \mathcal{D} \\ 4 \quad & \mathcal{D}_{sol} \leftarrow \mathcal{D}_{sol} \cup \{D_{min}\} \\ 5 \quad & \mathcal{D} \leftarrow \mathcal{D} \setminus \{D \in \mathcal{D} : D \cap D_{min} \neq \varnothing\} \\ 6 \quad \textbf{return } \mathcal{D}_{sol} \end{array}$

5900

イロト イポト イヨト イヨト 二日

2-intervals: Covering intervals

Definition

Let D = (I, J) be a 2-interval. The **covering interval** of D, denoted c(D), is the smallest interval that covers D (*i.e.*, c(D) = [l(I), r(J)], where l(i) (resp. r(J)) is the left (resp. right) endpoint of I (resp. J).

Observation

Let \mathcal{D} be a set of 2-intervals. For any $\{\sqsubset, \emptyset\}$ -structured subset $\mathcal{D}' \subseteq \mathcal{D}$, the associated covering intervals $c(\mathcal{D}')$ are pairwise intersecting.

San

イロト イポト イヨト イヨト 二日

Interval graphs

Definition (Interval graph)

An **interval graph** is the intersection graph of a multiset of intervals on the real line. It has one vertex for each interval in the set, and an edge between every pair of vertices corresponding to intervals that intersect.



Interval graphs

Remarks

- Determining whether a given graph G = (V, E) is an interval graph can be done in O(|V| + |E|) time by seeking an ordering of the maximal cliques of *G* that is consecutive with respect to vertex inclusion.
- A graph is an interval graph if and only if it is chordal and its complement is a comparability graph.
- The number of maximal cliques in a chordal graph is linear in the size of the graph.

5900

イロト イポト イヨト イヨト 二日

Chordal graphs

Definition (Chordal graph)

A graph is **chordal** if each of its cycles of four or more vertices has a chord, which is an edge that is not part of the cycle but connects two vertices of the cycle.


Comarability graphs

Definition (Comparability graph)

A **comparability graph** is a graph that has a transitive orientation, an assignment of directions to the edges of the graph (*i.e.*, an orientation of the graph) such that the adjacency relation of the resulting directed graph is transitive: whenever there exist directed edges (x, y) and (y, z), there must exist an edge (x, z).



2-intervals: $\mathcal{R} = \{\Box, \emptyset\}$

Theorem

The 2-INTERVAL PATTERN problem for unlimited (resp. unit) 2-intervals and model $\mathcal{R} = \{\Box, \emptyset\}$ is approximable within ratio 4 (resp. 3). (The algorithm is $O(n^2 \log(n))$ time for unit 2-intervals.)

Proof.

Data: A set of 2-intervals \mathcal{D} **Result**: A $\{\Box, \emptyset\}$ -structured subset of \mathcal{D} 1 $c(\mathcal{D}) \leftarrow \{c(D) : D \in \mathcal{D}\}$ 2 $\mathcal{K} \leftarrow$ all maximal cliques of $\Omega(c(\mathcal{D}))$ **foreach** *maximal clique* $K \in \mathcal{K}$ **do** 3 $| \mathcal{D}_K \leftarrow \{ D \in \mathcal{D} : c(D) \in K \}$ 4 $\mathcal{D}'_{K} \leftarrow$ (Approximate) pairwise disjoint subset of \mathcal{D}_{K} 5 6 return the largest \mathcal{D}'_{κ} found

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

Э 2016-2017 101 / 138

200

<ロト < 同ト < 巨ト < 巨ト

Trapezoid graphs

Definition (Trapezoid)

Consider two intervals *I* and *J* defined over two distinct horizontal lines. The **trapezoid** T = (I, J) is the convex set of points bounded by *I* and *J*, and the two arcs connecting the right and left endpoints of *I* and *J*. The interval *I* and *J* are the **top interval** and the **bottom interval** of *T*.

A **family of trapezoids** is a finite set of trapezoids which are all defined over the same two horizontal lines.

Definition (Trapezoid graph) A **trapezoid graph** is the intersection graph of a family of trapezoids.

200

イロト イポト イヨト イヨト 二日

Trapezoid graphs



10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 12 13 14 15 16 17 18 19 20 21 22 23 24

Э

<ロト < 同ト < 三ト < 三ト

Definition (Corresponding trapezoid family)

Let \mathcal{D} be a set of 2-intervals and let α and β be two distinct horizontal lines such that α is below β . The **corresponding trapezoid family** of \mathcal{D} , denoted $\mathcal{T}(\mathcal{D})$, is defined as the family containing a single trapezoid T = (I', J') for each 2-interval $D = (I, J) \in \mathcal{D}$, where I' is defined over α , J' is defined over β , and I' = I and J' = J.



Observation

Any two **disjoint** 2-intervals in \mathcal{D} are $\{<, \emptyset\}$ -comparable if and only their corresponding trapezoids in $\mathcal{T}(\mathcal{D})$ are disjoint.

Remarks

- Felsner et al. gave a $O(n \log(n))$ time algorithm for finding a maximum disjoint subset in a family of trapezoids.
- But there may be disjoint trapezoids in T(D) that correspond to non-disjoint 2-intervals in D.
- {□}-comparable 2-intervals in D correspond to intersecting trapezoids in T(D).

San

< ロト < 回 > < 三 > < 三 >





3 107 / 138 2016-2017

990

Definition (Clashing intervals)

Let $I_1 = [l(I_1), r(I_1)]$ and $I_2 = [l(I_2), r(I_2)]$ be two distinct intervals defined over two distinct horizontal lines such that $l(I_1) < l(I_2)$. The two intervals I_1 and I_2 **clash** if either

•
$$l(I_1) \le l(I_2) \le r(I_2) \le r(I_1)$$
, or
• $l(I_1) \le l(I_2) \le r(I_1) \le r(I_2)$.



Definition (Clashing trapezoids)

Let $T_1 = (I_1, J_1)$ and $T_2 = (I_2, J_2)$ be two distinct trapezoids in a family of trapezoids. The two trapezoids T_1 and T_2 **clash** if either

• I_1 and J_2 clash, or

• I_2 and J_1 clash.



Observation

Any two distinct 2-intervals in \mathcal{D} are $\{<, \emptyset\}$ -comparable if and only their corresponding trapezoids in $\mathcal{T}(\mathcal{D})$ are disjoint and do not clash.

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 110 / 138

Theorem

The 2-INTERVAL PATTERN *problem for unit (resp. point)* 2*-intervals and model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 3 (*resp.* 2)*. (The algorithm is* $O(n^2)$ *time for unit* 2*-intervals.)*

Proof.

```
Data: A set of 2-intervals \mathcal{D}

Result: A {<, §}-structured subset of \mathcal{D}

1 Construct the corresponding trapezoid family \mathcal{T}(\mathcal{D})

2 Compute \mathcal{T}' \subseteq \mathcal{T}(\mathcal{D}): the maximum pairwise disjoint subset o

3 \mathcal{T}_{sol} \leftarrow \emptyset

4 while \mathcal{T}' \neq \emptyset do

5 Let T_0 be the leftmost trapezoid in \mathcal{T}'

6 Let \mathcal{T}_{sol} \leftarrow \mathcal{T}_{sol} \cup \{T_0\}

7 Omit T_0 and all trapezoids clashing with T_0 in \mathcal{T}'

8 return the 2-intervals corresponding to trapezoids in \mathcal{T}_{sol}
```

Theorem

The 2-INTERVAL PATTERN *problem for unit (resp. point)* 2*-intervals and model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 3 (*resp.* 2)*. (The algorithm is* $O(n^2)$ *time for unit* 2*-intervals.)*

Proof.

Data: A set of 2-intervals \mathcal{D} **Result**: A $\{<, \emptyset\}$ -structured subset of \mathcal{D} Construct the corresponding trapezoid family $\mathcal{T}(\mathcal{D})$ Compute $\mathcal{T}' \subseteq \mathcal{T}(\mathcal{D})$: the maximum pairwise disjoint subset of $\mathcal{T}(\mathcal{D})$ $\mathcal{T}_{sol} \leftarrow \emptyset$ 3 while $\mathcal{T}' \neq \emptyset$ do Let T_0 be the leftmost trapezoid in \mathcal{T}' 5 $\mathcal{T}_{sol} \leftarrow \mathcal{T}_{sol} \cup \{T_0\}$ 6 Omit T_0 and all trapezoids clashing with T_0 in \mathcal{T}' 7 8 return the 2-intervals corresponding to trapezoids in \mathcal{T}_{sol} S. Vialette (CNRS & LIGM) RNA Secondary Structures 2016-2017 111 / 138

2-intervals: Model $\mathcal{R} = \{\langle, 0 \rangle\}$

Theorem

The 2-INTERVAL PATTERN problem for unit (resp. point) 2-intervals and model $\mathcal{R} = \{\langle, 0\rangle\}$ is approximable within ratio 3 (resp. 2). (The algorithm is $O(n^2)$ time for unit 2-intervals.)

Proof.

- Let T_0 be the leftmost trapezoid \mathcal{T}' and let D_0 be its corresponding 2-interval in \mathcal{D} .
- By our definition of a 2-interval and of $\mathcal{T}(\mathcal{D})$, any trapezoid in $\mathcal{T}(\mathcal{D})$ has a bottom interval which is completely to the left of its top interval.
- Thus, T_0 can only clash with trapezoids on its right in \mathcal{T}' .

200

Theorem

The 2-INTERVAL PATTERN *problem for unit (resp. point)* 2*-intervals and model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 3 (*resp.* 2)*. (The algorithm is* $O(n^2)$ *time for unit* 2*-intervals.)*

Proof.

if D is a point 2-interval set, then all 2-intervals with left intervals intersecting the right interval of D₀ have the same left interval, and as T' is pairwise disjoint, at most one of these has a corresponding trapezoid in T'.

2-intervals: Model $\mathcal{R} = \{\langle, 0 \rangle\}$

Theorem

The 2-INTERVAL PATTERN problem for unit (resp. point) 2-intervals and model $\mathcal{R} = \{\langle, 0\rangle\}$ is approximable within ratio 3 (resp. 2). (The algorithm is $O(n^2)$ time for unit 2-intervals.)

Proof.

- if \mathcal{D} is a unit 2-interval set, intersecting intervals involved in \mathcal{D} must overlap.
- Thus, any trapezoid in \mathcal{T}' clashing with T_0 corresponds to a 2-interval with a left interval which contains either endpoints, but not both, of the right interval of D_0 .
- Since \mathcal{T}' is pairwise disjoint, there can be at most two such trapezoids in \mathcal{T}' .

200

Theorem

The 2-INTERVAL PATTERN problem for balanced 2-intervals and model $\mathcal{R} = \{<, \emptyset\}$ is approximable within ratio 5. (The algorithm is $O(n^2)$ time.)

Proof.

 $\begin{array}{l} \textbf{Data: A set of balanced 2-intervals } \mathcal{D} \\ \textbf{Result: A } \{<, \check{0}\}\text{-structured subset of } \mathcal{D} \\ \text{1 Construct the corresponding trapezoid family } \mathcal{T}(\mathcal{D}) \\ \text{2 Compute } \mathcal{T}' \subseteq \mathcal{T}(\mathcal{D})\text{: the maximum pairwise disjoint subset of } \mathcal{T}(\mathcal{D}) \\ \text{3 } \mathcal{T}_{sol} \leftarrow \varnothing \\ \text{4 while } \mathcal{T}' \neq \varnothing \ \textbf{do} \\ \text{5 } \\ \text{6 } \\ \text{7 } \\ \text{Compute } \mathcal{T}_{sol} \leftarrow \mathcal{T}_{sol} \cup \{T_0\} \\ \text{7 } \\ \text{7 } \\ \text{7 } \\ \text{7 } \\ \text{8 } \mathbf{return } \textit{the 2-intervals corresponding to trapezoids in } \mathcal{T}_{sol} \\ \end{array}$

200

Definition (Proper trapezoid family)

A family of trapezoids \mathcal{T} is **proper** if for any two distinct trapezoids $T_1 = (I_1, J_1)$ and $T_2 = (I_2, J_2)$ in $\mathcal{T}, I_1 \cap I_2 = \emptyset$ and $J_1 \cap J_2 = \emptyset$.



2016-2017 114 / 138

Definition (Strongly proper trapezoid family)

A proper family of trapezoids \mathcal{T} is **strongly proper** if for any two distinct trapezoids $T_1 = (I_1, J_1)$ and $T_2 = (I_2, J_2)$ in \mathcal{T} , if J_1 and I_2 clash then $l(I_2) \leq l(J_1) < r(J_1) \leq r(I_2)$, where $l(J_1), r(J_1)$ and $l(I_2), r(I_2)$ are the left and right endpoints of J_1 and I_2 , respectively.



< ロト < 同ト < 臣ト < 臣ト -

Remarks

- Any pairwise disjoint family of trapezoids is proper (but the converse is not true).
- Thus Step 2 of the preceding algorithm yields a proper family of trapezoids \mathcal{T}' .

2-intervals: Model $\mathcal{R} = \{\langle, 0 \rangle\}$

Remarks

• Computing a strongly proper subset $\mathcal{T}' \subseteq \mathcal{T}''$ can be done easily by adjusting the loop step: Instead of omitting all trapezoids clashing with the leftmost trapezoid in this iteration, we need only to omit a small subset of these trapezoids.

More specifically, let $T_0 = (I_0, J_0)$ be the leftmost trapezoid in \mathcal{T}' . We only omit trapezoids T = (I, J) with

▶
$$l(I) \le l(I_0) \le r(I)$$
, or
▶ $l(I) < r(I_0) < r(I)$.

• We obtain a strongly proper trapezoid family $\mathcal{T}'' \subseteq \mathcal{T}'$ if we proceed in this fashion such that $3 |\mathcal{T}''| > |\mathcal{T}'|$.

San

Definition (Clashing trapezoid graph)

Let \mathcal{T} be a family of trapezoids. The **clashing trapezoid graph** of \mathcal{T} , denoted $G_{\mathcal{T}}$, is the graph such that each vertex in $G_{\mathcal{T}}$ correspond to a distinct trapezoid in \mathcal{T} , and two vertices are connected by an edge if and only if their corresponding trapezoid clash.



Theorem

Let T *be a family of trapezoids. If* T *is strongly proper then* G_T *is a forest.*

Proof.

Theorem

Let T *be a family of trapezoids. If* T *is strongly proper then* G_T *is a forest.*

Proof.

$$V(\vec{G}_{\mathcal{T}}) = V(G_{\mathcal{T}})$$

 $E(\vec{G}_{\mathcal{T}}) = \{(T_1, T_2) : \{T_1, T_2\} \in E(G_{\mathcal{T}}) \text{ and } T_1 < T_2\}$

• Since \mathcal{T} is strongly proper, every trapezoid in \mathcal{T} clashes with at most one trapezoid on its left, and hence the in-degree of every vertex $T \in V(\vec{G}_{\mathcal{T}})$ is at most 1.

< ロト < 回 > < 三 > < 三 >

2-intervals: Model $\mathcal{R} = \{\langle, \rangle\}$

Theorem

Let \mathcal{T} be a family of trapezoids. If \mathcal{T} is strongly proper then $G_{\mathcal{T}}$ is a forest.

Proof.

- Hence any cycle (T_0, \ldots, T_k, T_0) in $G_{\mathcal{T}}$ is a directed cycle in $G_{\mathcal{T}}$.
- Then we must have $T_0 < T_k < T_0$ by definition of \vec{G}_{τ} . This is a contradiction.
- Then it follows that $G_{\mathcal{T}}$ is acyclic.

200

<ロト < 回 > < 回 > < 回 > < 回 >

Remarks

- A maximum independent set in any forest of size *n* is of size at least ⁿ/₂. (This set can be found in linear time with respect to *n*.)
- if *T* is a pairwise disjoint family of trapezoids, then any independent set of *G*_{*T*} corresponds to a pairwise disjoint non-clashing set of trapezoids, and hence corresponds to a {<, ≬}-comparable set of 2-intervals.

San

< ロ ト < 同 ト < 三 ト < 三 ト -

Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 6. *(The algorithm is O*(n^2) *time.)*

Proof.

Data: A set of 2-intervals \mathcal{D} **Result**: A {<, \emptyset }-structured subset of \mathcal{D}

- 1 Construct $\mathcal{T}(\mathcal{D})$, the corresponding trapezoid set of \mathcal{D}
- ² Compute \mathcal{T}' , the maximum pairwise disjoint subset of $\mathcal{T}(\mathcal{D})$
- ³ Compute \mathcal{T}'' , a strongly proper subset of \mathcal{T}' , such that $3|\mathcal{T}''| \ge |\mathcal{T}'|$
- 4 Compute $G_{\mathcal{T}''}$ and the maximum independent set of $G_{\mathcal{T}''}$
- 5 return the 2-intervals corresponding to the maximum independent set of $G_{T''}$

5900

< D > < B > < E > < E > <</p>

Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 2. *(The algorithm is* $O(n^3 \log(n))$ *) time.)*

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

2016-2017 122 / 138

Definition (Precedence or crossing)

For two 2-intervals $D_1 = (I_1, J_1)$ and $D_2 = (I_2, J_2)$, the **precedence or crossing** relation \notin is defined by:

 $D_1 \notin D_2$ if and only if $D_1 < D_2$ or $D_1 \notin D_2$

Remarks

• If $D_1 \notin D_2$ then $I_1 < I_2$ and $J_1 < J_2$.

The & relation specifies a total order for any {<, ≬}-structured
 2-interval set.

S. Vialette (CNRS & LIGM)

2016-2017 123 / 138

<ロト < 同ト < 三ト < 三ト < 三ト < 三 の < ○</p>

Let \mathcal{D} be a {<, §}-structured 2-interval set and let \mathcal{D}_{\notin} denotes \mathcal{D} ordered by the \notin relation (viewed as an ordered sequence).

- $\mathcal{D}_{\mathfrak{F}}[i]$ denotes the 2-interval with rank *i* in \mathcal{D} .
- $\mathcal{D}_{\check{\&}}[i,j]$ denotes the subsequence $\mathcal{D}_{\check{\&}}[i], \mathcal{D}_{\check{\&}}[i+1], \dots, \mathcal{D}_{\check{\&}}[j]$.

- For each index $1 \le i \le |\mathcal{D}_{\mathfrak{F}}|$, define $\operatorname{next}(i)$ as the smallest index j, $1 \le j \le |\mathcal{D}_{\mathfrak{F}}|$ such that $\mathcal{D}_{\mathfrak{F}}[i] < \mathcal{D}_{\mathfrak{F}}[j]$. If such an index j does not exist define $\operatorname{next}(i) = |\mathcal{D}_{\mathfrak{F}}| + 1$.
- Define the **backbone indices** of $\mathcal{D}_{\&}$ as the sequence of indices i_1, i_2, \ldots, i_k such that $i_1 = 1, i_j = \text{next}(i_{j-1})$ and $\text{next}(i_k) = |\mathcal{D}_{\&}| + 1$.

(For convenience, we define next(i_k) = $|\mathcal{D}_{\check{\mathfrak{A}}}| + 1$ and imagine a 2-interval $\mathcal{D}_{\check{\mathfrak{A}}}[i_{k+1}]$ such that $\mathcal{D}_{\check{\mathfrak{A}}}[i] < \mathcal{D}_{\check{\mathfrak{A}}}[i_{k+1}]$ for all $1 \leq i \leq |\mathcal{D}_{\check{\mathfrak{A}}}|$.)

<ロト < 同ト < 三ト < 三ト < 三ト < 三 の < ○</p>



3 2016-2017 126 / 138

200

- For each backbone index $1 \le i_s \le k$, define a **stripe** $\mathcal{T}(i_s) = \mathcal{D}_{\&}[i_s + 1, i_{s+1} 1].$
- The stripe is **odd** if *s* is odd; it is **even** if *s* is even.
- For each 2-interval $D \in \mathcal{T}(i_s)$, we observe that $\mathcal{D}_{\&}[i_s] \& D \& \mathcal{D}_{\&}[i_s+1]$.
- Every stripe of $\mathcal{D}_{\&}$ is $\{\emptyset\}$ -structured.



2016-2017 128 / 138

3

200

イロト 不良 とうほう 不良 とう

- A {<, $\hfill\)}-structured sequence <math display="inline">\mathcal{D}_{\hfill\)}$ is striped if either
 - ► its odd stripes are all empty, or
 - its even stripes are all empty.
- Although D_𝔅 is not always striped, it contains two striped subsequences:

$$\mathcal{D}_{\check{\mathfrak{A}}}[i_1] \ \mathcal{T}(i_1) \ \mathcal{D}_{\check{\mathfrak{A}}}[i_2] \ \mathcal{D}_{\check{\mathfrak{A}}}[i_3] \ \mathcal{T}(i_3) \ \mathcal{D}_{\check{\mathfrak{A}}}[i_4] \ \dots \\ \mathcal{D}_{\check{\mathfrak{A}}}[i_1] \ \mathcal{D}_{\check{\mathfrak{A}}}[i_2] \ \mathcal{T}(i_2) \ \mathcal{D}_{\check{\mathfrak{A}}}[i_3] \ \mathcal{D}_{\check{\mathfrak{A}}}[i_4] \ \mathcal{T}(i_4) \ \dots$$

These two subsequences together cover the sequence $\mathcal{D}_{\hat{k}}$: the 2-intervals at the backbone indices are covered twice, each remaining 2-interval is covered once.

One of the two subsequences has a length of at least $|\mathcal{D}_{\&}|/2$.

2016-2017 129 / 138

5900

<ロト < 同ト < 巨ト < 巨ト = 巨 =</p>

- A {<, $\hfill\)}-structured sequence <math display="inline">\mathcal{D}_{\hfill\)}$ is striped if either
 - ► its odd stripes are all empty, or
 - its even stripes are all empty.
- Although D_𝔅 is not always striped, it contains two striped subsequences:

$$\mathcal{D}_{\mathfrak{F}}[i_1] \ \mathcal{T}(i_1) \ \mathcal{D}_{\mathfrak{F}}[i_2] \ \mathcal{D}_{\mathfrak{F}}[i_3] \ \mathcal{T}(i_3) \ \mathcal{D}_{\mathfrak{F}}[i_4] \ \dots$$
$$\mathcal{D}_{\mathfrak{F}}[i_1] \ \mathcal{D}_{\mathfrak{F}}[i_2] \ \mathcal{T}(i_2) \ \mathcal{D}_{\mathfrak{F}}[i_3] \ \mathcal{D}_{\mathfrak{F}}[i_4] \ \mathcal{T}(i_4) \ \dots$$

These two subsequences together cover the sequence $\mathcal{D}_{\hat{\mathbf{x}}}$: the 2-intervals at the backbone indices are covered twice, each remaining 2-interval is covered once.

One of the two subsequences has a length of at least $|\mathcal{D}_{\check{\&}}|/2$.

5900

イロト (四) (注) (注) (注) [- [- [-
Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 2. *(The algorithm is* $O(n^3 \log(n))$ *) time.)*

Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 2. *(The algorithm is* $O(n^3 \log(n))$ *) time.)*

Proof.

• **Step 1.** Make a dummy 2-interval D_{ω} such that $D_{\gamma} < D_{\omega}$ for all $D_{\gamma} \in \mathcal{D}$.

Set $\mathcal{D}^+ = \mathcal{D} \cup \{D_\omega\}.$

Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 2. *(The algorithm is* $O(n^3 \log(n))$ *) time.)*

Proof.

• **Step 2.** For each pair of 2-intervals D_{α} and D_{β} in \mathcal{D}^+ , $D_{\alpha} < D_{\beta}$, find the subset of 2-intervals

$$\mathcal{D}^+_{\alpha,\beta} = \{D_{\gamma} : D_{\gamma} \in \mathcal{D}^+ \text{ and } D_{\alpha} \And D_{\gamma} \And D_{\beta}\}$$

Then compute $C_{\alpha,\beta}$, a maximum size $\{\emptyset\}$ -structured subset of $\mathcal{D}^+_{\alpha,\beta}$.

Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 2. *(The algorithm is* $O(n^3 \log(n))$ *) time.)*

Proof.

Step 3.1. Process the 2-intervals in D⁺_{α,β} in an arbitrary order that conforms to the partial order specified by the < relation. For each 2-interval D_β in D⁺, find the subset of 2-intervals

$$\mathcal{D}^+_{eta} = \{D_{lpha}: D_{lpha} \in \mathcal{D}^+ ext{ and } D_{lpha} < D_{eta}\}$$

Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 2. *(The algorithm is* $O(n^3 \log(n))$ *) time.)*

Proof.

• Step 3.2. If
$$\mathcal{D}_{\beta}^{+} = \emptyset$$
, $\mathcal{A}_{\beta} \leftarrow \{\beta\}$ and $\mathcal{B}_{\beta} \leftarrow \{\beta\}$. Otherwise

Find D_α ∈ D⁺_β such that |B_α| is maximum and set A_β ← B_β ∪ {D_β},
Find D_α ∈ D⁺_β such that |A_α| + |C_{α,β}| is maximum and set B_β ← A_α ∪ C_{α,β} ∪ {D_β}.

Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 2. *(The algorithm is* $O(n^3 \log(n))$ *) time.)*

Proof.

Step 4. Let D_{sol} be either A_ω or B_ω such that |D_{sol}| is maximum.
 Return D_{sol} \ {ω}.

Theorem

The 2-INTERVAL PATTERN *problem for model* $\mathcal{R} = \{<, \emptyset\}$ *is approximable within ratio* 2. *(The algorithm is* $O(n^3 \log(n))$ *) time.)*

Proof.

Notes

- In the algorithm we use A_{β} and B_{β} to represent the two different alternating patterns, with β as both the last element backbone element.
- The 2-interval D_α in steps 3.1 and 3.2 represents the second-to-last backbone element in A_β and B_β.
- The subset C_{α,β} represents the maximum size stripe between the two backbone elements D_α and D_β.

d-claw free graphs

Definition (d-claw, d-claw-free)

For an undirected graph *G*, a *d*-claw *C* is an induced subgraph $K_{1,d}$ that consists of an independent set T_c of *d* vertices (called **talons**) and a **center** vertex z_C that is connected to all the talons.

A graph is *d*-claw-free if it has no *d*-claws.

(A 3-claw is commonly called a claw so that a graph is claw-free if and only if it does not contain the complete bipartite graph $K_{1,3}$ (known as the "**claw graph**") as an induced subgraph.)



4 Ξ ≥ < Ξ ≥ ...</p>

d-claw free graphs



The regular icosahedron, a polyhedron whose vertices and edges form a claw-free graph.

S. Vialette (CNRS & LIGM)

RNA Secondary Structures

< □ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 ♪ < ○ ○</p>
2016-2017 132 / 138

d-claw free graphs

Definition (The MAXIMUM WEIGHT INDEPENDENT SET problem)

- **Input**: A graph G = (V, E) and a weight function $w : V \to \mathbb{N}$.
- Solution: A set of independent vertices $V' \subseteq V$ that maximises $\sum_{v \in V'} w(v)$.

Remarks

In *d*-claw-free graphs

- Arbitrary weight: $(d/2 + \epsilon)$ -approximation
- Small weight: (*d*/2)-approximation

• Unit weight: $((d-1)/2 + \epsilon)$ -approximation

5900

イロト イポト イヨト イヨト 二日

2-intervals and *d*-claw free graphs

Theorem

For a set of 2-intervals D with interval length ℓ , $a \leq \ell \leq b$, the 2-interval graph G(D) is d-claw-free for

$$d = 5 + \frac{2(b-2)}{a}.$$

Proof.

- Let *I* be an interval and let *I* be a set of disjoint intervals that intersect *I*.
- All intervals in *I* are completely contained in *I* except possibly the leftmost one and the rightmost one.

San

イロト 不得下 イヨト イヨト

2-intervals and *d*-claw free graphs

Theorem

For a set of 2-intervals D with interval length ℓ , $a \leq \ell \leq b$, the 2-interval graph G(D) is d-claw-free for

$$d = 5 + \frac{2(b-2)}{a}.$$

Proof.

- Let *I* be an interval and let *I* be a set of disjoint intervals that intersect *I*.
- All intervals in *I* are completely contained in *I* except possibly the leftmost one and the rightmost one.

San

< ロト < 回 > < 回 > < 回 >

2-intervals and *d*-claw free graphs

Corollary

Let D *be a set of* 2*-intervals.*

- *if all intervals have the same length (unit support), the associated* 2-*interval graph* $G(\mathcal{D})$ *is* 5-*claw-free;*
- if all intervals have length 2 or 3, the associated 2-interval graph $G(\mathcal{D})$ is 5-claw-free

Corollary

The (WEIGHTED) 2-INTERVAL PATTERN *problem is approximable within ratio* 2.5 + ϵ *for arbitrary weights and* 2 + ϵ *for unit weights.*

Nearest Neighbor Energy Model

For a secondary structure *S*

• the number of **base pairs stackings** is

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

• the number of **stacking base pairs** is

$$SBP(S) = |\{(i,j) \in S : (i+1,j-1) \in S \text{ or } (i-1,j+1) \in S\}|$$

500

- E - E

I D > I A

Nearest Neighbor Energy Model

Definition (The MAXIMUM BASE PAIRS STACKING (BPS) problem)

- Input: A sequence *u*.
- **Solution**: A secondary structure *S* for *u* that maximises BPS(*S*).

Definition (The MAXIMUM STACKINGBASE PAIRS (SBP) problem)

Input: A sequence *u*. •

• **Solution**: A secondary structure *S* for *u* that maximises SBP(*S*).

San

イロト イポト イヨト イヨト

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

- **Step 1**. Repeatedly find the leftmost 5 consecutive stacking loops (*i.e.*, find the 2-interval ([x, x + 5], [y 5, y]) where x is as small as possible). Add these stacking loops to *S*.
- **Step 2**. Repeatedly find any 4 consecutive stacking loops. Add these stacking loops to *S*.
- **Step 3**. Repeatedly find any 3 consecutive stacking loops. Add these stacking loops to *S*.

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

- **Step 4.1**. Construct a 2-interval set \mathcal{D} by associating a 2-interval to each 2 consecutive stacking loop.
- **Step 4.2**. Construct the 2-interval graph *G*(*D*) and assign each vertex a weight: 1 for a single stacking loop and 2 for two consecutive stacking loops.

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

- **Step 4.3**. Find a maximum weight independent set \mathcal{D}' in $G(\mathcal{D})$ (5/2-approximation algorithm for 5-claw-free graphs).
- **Step 4.4**. For each 2-interval in \mathcal{D}' , add the corresponding stacking loop in *S*.

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

Proof.

• Let *s*₁, *s*₂, *s*₃ and *s*₄, respectively, be the number of stacking loops found by the first, second, third and fourth of our algorithm.

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

- Let *S*^{*} be the set of stacking loops in an optimal secondary structure.
- Let s_1^*, s_2^* and s_3^* , respectively, be the number of stacking loops in S^* that intersect the stacking loops found by the first, second and third step of our algorithm.
- Let s_4^* be the number of remaining stacking loops in S^* which are represented by 2-intervals in \mathcal{D} .

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

$$|S| = s_1 + s_2 + s_3 + s_4$$
$$|S^*| = s_1^* + s_2^* + s_3^* + s_4^*$$

Theorem

The MAXIMUM BASE PAIRS STACKING *problem is approximable within ratio* 8/3.

Proof.

For each *k* consecutive stacking loops *D* found by the first three steps of our algorithm, the number of stacking loops in S^* that intersect them is at most 2(k + 2) (*i.e.*, k + 2 for each interval of the 2-interval *D*).

Theorem

The MAXIMUM BASE PAIRS STACKING *problem is approximable within ratio* 8/3.

Proof.

Step 1.

 By always choosing the leftmost 5-consecutive stacking loop D₅, we can guarantee that the left interval of the 2-interval D₅ intersects at most 5 + 1 stacking loops in S*.

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

Proof.

Step 1.

- Suppose the contrary that the left interval of *D*₅ intersects 7 stacking loops in *S*^{*}.
- Then these 7 stacking loops must be consecutive, and the leftmost 5 of these stacking loops should have been choosen instead of *D*₅.

Theorem

The MAXIMUM BASE PAIRS STACKING *problem is approximable within ratio* 8/3.

Proof.

Step 1.

$$\frac{s_1^*}{s_1} \le \frac{5+1+5+2}{5} = \frac{13}{5} = 2.6$$

Theorem

The MAXIMUM BASE PAIRS STACKING *problem is approximable within ratio* 8/3.

Proof.

Step 2.

With all 5 consecutive stacking loops found by the first step, we can guarantee that each interval of a 2-interval D₄ (consisting of 4 consecutive stacking loops) found by the second step of our algorithm intersects at most 4 + 1 stacking loops in S*.

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

Proof.

Step 2.

- Suppose the contrary that an interval of *D*₅ intersects 6 stacking loops in *S**.
- Then these 6 stacking loops must be consecutive, and hence must contain 5 consecutive stacking loops.
- This is a contradiction.

Theorem

The MAXIMUM BASE PAIRS STACKING *problem is approximable within ratio* 8/3.

Proof.

Step 2.

$$\frac{s_2^*}{s_2} \le \frac{4+1+4+1}{4} = \frac{10}{4} = 2.5$$

Theorem

The MAXIMUM BASE PAIRS STACKING *problem is approximable within ratio* 8/3.

Proof.

Step 3.

$$\frac{s_3^*}{s_3} \le \frac{3+1+3+1}{4} = \frac{8}{3} \simeq 2.67$$

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

Proof.

Step 4.

- Each 2-interval in \mathcal{D} is balanced and corresponds to either a single stacking loop (with interval length 2) of two consecutive stacking loops (with interval length 3).
- Therefore the 2-interval graph $G(\mathcal{D})$ is 5-claw-free.

Theorem

The MAXIMUM BASE PAIRS STACKING *problem is approximable within ratio* 8/3.

Proof.

Step 4.

• The MAXIMUM WEIGHT INDEPENDENT SET problem in 5-claw-free graphs is approximable within ratio 5/2.

$$\frac{s_4^*}{s_4} \le \frac{5}{2} = 2.5$$

Theorem

The **MAXIMUM BASE PAIRS STACKING** *problem is approximable within ratio* 8/3.

$$\frac{S^*|}{|S|} = \frac{\sum_{i=1}^4 s_i^*}{\sum_{i=1}^4 s_i} \\ = \sum_{i=1}^4 \frac{s_i^*}{\sum_{j=1}^4 s_j}$$

Theorem

The MAXIMUM BASE PAIRS STACKING *problem is approximable within ratio* 8/3.

$$\begin{aligned} \frac{|S^*|}{|S|} &\leq \frac{13}{5} \left(\frac{s_1}{\sum_{j=1}^4 s_j} \right) + \frac{10}{4} \left(\frac{s_2}{\sum_{j=1}^4 s_j} \right) + \frac{8}{3} \left(\frac{s_3}{\sum_{j=1}^4 s_j} \right) + \frac{5}{2} \left(\frac{s_4}{\sum_{j=1}^4 s_j} \right) \\ &\leq \frac{8}{3} \left(\frac{s_1}{\sum_{j=1}^4 s_j} \right) + \frac{8}{3} \left(\frac{s_2}{\sum_{j=1}^4 s_j} \right) + \frac{8}{3} \left(\frac{s_3}{\sum_{j=1}^4 s_j} \right) + \frac{8}{3} \left(\frac{s_4}{\sum_{j=1}^4 s_j} \right) \\ &= \frac{8}{3} \end{aligned}$$