

AMIB(io): Algorithms and Models for Integrative Biology

Team AMIBio Inria Saclay IdF CNRS UMR 7161 LIX



Algorithms and Models for Integrative Biology





Algorithms and Models for Integrative Biology





Algorithms and Models for Integrative Biology





Algorithms and Models for Integrative Biology





Algorithms and Models for Integrative Biology





Algorithms and Models for Integrative Biology





Algorithms and Models for Integrative Biology



Team Members

Mireille Régnier DR



Postdoc: Christelle Rovetta 2017

PhD Students:

Jorgelindo Da Veiga – 2016 Alice Héliou – 2014 T Amélie Héliou – 2014 T Wei Wang – 2014 T Juraj Michalik – 2016 U Afaf Saaidi – 2015 🛞 🛲 Antoine Soulé – 2014 M Ha Nguyen Ngoc – 2017 🕬 Pauline Pommeret – 2017 🕬



AMIBio's Methods:

- Stringology
- Dynamic Programming
- Formal Languages
- Random Generation
- Graph Theory
- Probabilistic models

- Analytic Combinatorics
- Clustering
- Complexity Theory
- Game Theory
- Kinematics
- Computational Geometry

- What are the active conformation(s) (2D&3D) of a molecule?
- What is the impact of genomic mutations at a structural level?
- Are there structured transcripts within an HTS dataset?
- Can we characterize the kinetics of biopolymers?
- How to predict the geometry of complexes?
- How to design functional molecules using rational principles?



AMIBio's Methods:

- Stringology
- Dynamic Programming
- Formal Languages
- Random Generation
- Graph Theory
- Probabilistic models

- Analytic Combinatorics
- Clustering
- Complexity Theory
- Game Theory
- Kinematics
- Computational Geometry

- What are the active conformation(s) (2D&3D) of a molecule?
- What is the impact of genomic mutations at a structural level?
- Are there structured transcripts within an HTS dataset?
- Can we characterize the kinetics of biopolymers?
- How to predict the geometry of complexes?
- How to design functional molecules using rational principles?



AMIBio's Methods:

- Stringology
- Dynamic Programming
- Formal Languages
- Random Generation
- Graph Theory
- Probabilistic models

- Analytic Combinatorics
- Clustering
- Complexity Theory
- Game Theory
- Kinematics
- Computational Geometry

- What are the active conformation(s) (2D&3D) of a molecule?
- What is the impact of genomic mutations at a structural level?
- Are there structured transcripts within an HTS dataset?
- Can we characterize the kinetics of biopolymers?
- How to predict the geometry of complexes?
- How to design functional molecules using rational principles?



AMIBio's Methods:

- Stringology
- Dynamic Programming
- Formal Languages
- Random Generation
- Graph Theory
- Probabilistic models

- Analytic Combinatorics
- Clustering
- Complexity Theory
- Game Theory
- Kinematics
- Computational Geometry

- What are the active conformation(s) (2D&3D) of a molecule?
- What is the impact of genomic mutations at a structural level?
- Are there structured transcripts within an HTS dataset?
- Can we characterize the kinetics of biopolymers?
- How to predict the geometry of complexes?
- How to design functional molecules using rational principles?



AMIBio's Methods:

- Stringology
- Dynamic Programming
- Formal Languages
- Random Generation
- Graph Theory
- Probabilistic models

- Analytic Combinatorics
- Clustering
- Complexity Theory
- Game Theory
- Kinematics
- Computational Geometry

- What are the active conformation(s) (2D&3D) of a molecule?
- What is the impact of genomic mutations at a structural level?
- Are there structured transcripts within an HTS dataset?
- Can we characterize the kinetics of biopolymers?
- How to predict the geometry of complexes?
- How to design functional molecules using rational principles?



Fundamental dogma of molecular biology





AMIB EPI/AMIBio EI – Inria Saclay – Computational Biology Evaluation Seminar

Fundamental dogma of molecular biology (v2.0)





AMIB EPI/AMIBio EI – Inria Saclay – Computational Biology Evaluation Seminar

Fundamental dogma of molecular biology





AMIB EPI/AMIBio EI – Inria Saclay – Computational Biology Evaluation Seminar

RNA world: Solving the chicken vs egg paradox at the origin of life...



A gene big enough to specify an enzyme would be too big to replicate accurately without the aid of an enzyme of the very kind that it is trying to specify. So the system apparently cannot get started.

[...] This is the RNA World. To see how plausible it is, we need to look at why proteins are good at being enzymes but bad at being replicators; at why DNA is good at replicating but bad at being an enzyme; and finally why *RNA might just be good enough at both roles to break out of the Catch-22.*

R. Dawkins. The Ancestor's Tale: A Pilgrimage to the Dawn of Evolution



RNA world: Solving the chicken vs egg paradox at the origin of life...



A gene big enough to specify an enzyme would be too big to replicate accurately without the aid of an enzyme of the very kind that it is trying to specify. So the system apparently cannot get started.

[...] This is the RNA World. To see how plausible it is, we need to look at why proteins are good at being enzymes but bad at being replicators; at why DNA is good at replicating but bad at being an enzyme; and finally why RNA might just be good enough at both roles to break out of the Catch-22.

R. Dawkins. The Ancestor's Tale: A Pilgrimage to the Dawn of Evolution



RNA sequence and structure(s)





AMIB EPI/AMIBio EI - Inria Saclay - Computational Biology Evaluation Seminar

Crossing interactions

Excluded from the secondary structure:

 Non-canonical base-pairs: Any base-pair other than {(A-U), (C-G), (G-U)} OR interacting in a non-standard way (WC/WC-Cis) [Leontis Westhof, RNA 2001].





Canonical CG base-pair (WC/WC-Cis)

Non-canonical base-pair (Sugar/WC-Trans)

(Pseudo?)knots: Crossing sets of nested stable base-pairs





Crossing interactions

Excluded from the secondary structure:

 Non-canonical base-pairs: Any base-pair other than {(A-U), (C-G), (G-U)} OR interacting in a non-standard way (WC/WC-Cis) [Leontis Westhof, RNA 2001].





Canonical CG base-pair (WC/WC-Cis)

Non-canonical base-pair (Sugar/WC-Trans)

(Pseudo?)knots: Crossing sets of nested stable base-pairs



Group I Intron (PDBID: 1Y0Q:A)



Crossing interactions

Excluded from the secondary structure:





Thermodynamics view

RNA folding can be adequately abstracted as a (continuous) Markov process, whose stationary distribution is the Boltzmann distribution.



Definition (Thermodynamic equilibrium)

Each structure S compatible with an RNA w observed with probability:

 $\mathbb{P}(S \mid w) = \frac{e^{\frac{-E_w(S)}{kT}}}{\mathcal{Z}_w} \quad \text{and} \quad \mathcal{Z}_w \equiv \sum_{S'} e^{\frac{-E_w(S')}{RT}} \quad \{\text{Partition function}\}\$ $E_w(S): \text{ free-energy of } S \text{ over } w; R: \text{Boltzmann constant; and } T: \text{ temperature.}$



Dynamic programming (DP) for RNA folding

Theorem (Exact solution – BP models [Nussinov1980]) Most stable structure computable in $O(n^3)/O(n^2)$ time/memory



 $E_{i,k}$: Free-energy contribution of base-pair (i, k).

 $N_{i,j}$: Max #base-pairs over interval [i, j]

$$N_{i,t} = 0, \quad \forall t \in [i, i + \theta]$$

$$N_{i,j} = \min \begin{cases} N_{i+1,j} & \{i \text{ unpaired}\}\\ \min_{k=i+\theta+1} E_{i,k} + N_{i+1,k-1} + N_{k+1,j} & \{i \text{ paired to } k\} \end{cases}$$



Dynamic programming (DP) for RNA folding

Theorem (Exact solution – BP models [Nussinov1980]) Most stable structure computable in $O(n^3)/O(n^2)$ time/memory



 $E_{i,k}$: Free-energy contribution of base-pair (i, k).

 $C_{i,j}$: Number of secondary structures restricted to [i, j]

$$\begin{aligned} \boldsymbol{C}_{i,t} &= \boldsymbol{1}, \quad \forall t \in [i, i + \theta] \\ \boldsymbol{C}_{i,j} &= \sum \left\{ \begin{array}{c} \boldsymbol{C}_{i+1,j} & \{i \text{ unpaired}\} \\ \sum_{k=i+\theta+1}^{j} \mathbb{1}_{\text{comp.}(i,k)} \times \boldsymbol{C}_{i+1,k-1} \times \boldsymbol{C}_{k+1,j} & \{i \text{ paired to } k\} \end{array} \right. \end{aligned}$$



Dynamic programming (DP) for RNA folding

Theorem (Exact solution – BP models [Nussinov1980]) Most stable structure computable in $O(n^3)/O(n^2)$ time/memory



 $E_{i,k}$: Free-energy contribution of base-pair (i, k).

 $\mathcal{Z}_{i,j} = \sum_{\substack{S \text{ comp.} \\ \text{with } w_{[i,j]}}} e^{-\frac{E_w(S)}{RT}} = \text{Partition function of structures over } [i,j]$

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i + \theta] \\ \mathcal{Z}_{i,j} &= \sum \left\{ \begin{array}{cc} \mathcal{Z}_{i+1,j} & \{i \text{ unpaired}\} \\ \sum_{k=i+\theta+1}^{j} e^{\frac{-\mathcal{E}_{i,k}}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} & \{i \text{ paired to } k\} \end{array} \right. \end{aligned}$$



Dynamic programming (DD) for DNA folding Many extensions:

- Nearest-neighbor/Turner energy model
- Comparative folding
- Equilibrium base-pairing probabilities
- Moments of additive features
- ► Δ kcal.mol⁻¹ suboptimal structures of MFE
- Basic crossing structures
- Exact sampling in Boltzmann distr.
- Moments of additive features
- Maximum expected accuracy structure structures over [/ [Do2006]
- Distance-classified partitioning of Boltzmann ens. [E.Freyhult2007a]

Made possible by:

- Completeness/Unambiguity of decomposition
- Objective function additive with respect to DP scheme
 - ⇒ Combinatorial Dynamic Programming



[Ding2003,Ponty2008] [Miklos2005.Ponty2011]

[Miklos2005,Ponty2011]

[Wuchty1999]

[Zuker1981]

[Sankoff1985]

[McCaskill1990]

[Rivas1999] . . .

SHAPE probing

SHAPE = Selective 2'-Hydroxyl Acylation by Primer Extension

Produces accessibility profiles, i.e. *projections* of RNA structure.

in silico analysis of SHAPE data remains complex and misleading.



Our goal: To develop *hybrid* modeling approaches with experimentalists (Paris V) and bioinformaticians (McGill)

- Massively parallel derivation of profiles (PCR/NGS/EM);
- Clustering to model structure of viruses (HIV, Ebola);
- ► Joint analysis of multiple SHAPE profiles.

Funded by Fondation pour la Recherche Médicale (2015-2018)

[Desforges, ..., Saaidi, Ponty, Ohlmann, Sargueil, NAR 2017]



SHAPE probing

SHAPE = Selective 2'-Hydroxyl Acylation by Primer Extension

Produces accessibility profiles, i.e. *projections* of RNA structure.

in silico analysis of SHAPE data remains complex and misleading.

Our goal: To develop *hybrid* modeling approaches with experimentalists (Paris V) and bioinformaticians (McGill)

- Massively parallel derivation of profiles (PCR/NGS/EM);
- Clustering to model structure of viruses (HIV, Ebola);
- ► Joint analysis of multiple SHAPE profiles.

Funded by Fondation pour la Recherche Médicale (2015-2018)

[Desforges, ..., Saaidi, Ponty, Ohlmann, Sargueil, NAR 2017]







Interface: Dyn. Prog./Rand. Gen. (RNA Design) 1/3



5s rRNA (PDBID: 1K73:B)



Interface: Dyn. Prog./Rand. Gen. (RNA Design) 1/3





Dyn. Prog./Rand. Gen. (RNA Design) 2/3

Input: Set of constraints

- Secondary structure
- Pattern avoidance/occurrence
- Energy/robustness
- Putative interactions

▶ ...

Goal:

Design of active RNAs

Method:

Generation + Selection



Dyn. Prog./Rand. Gen. (RNA Design) 3/3

Fact #1: Selection is expensive

⇒ Capture constraints during generation stage

Fact #2: Goals of synthetic biology are evolving ⇒ Need for modular approaches Our approach: Non-uniform *Boltzmann* random generation



AMIB EPI/AMIBio EI – Inria Saclay – Computational Biology Evaluation Seminar

Dyn. Prog./Rand. Gen. (RNA Design) 3/3

Fact #1: Selection is expensive

 \Rightarrow Capture constraints during generation stage

Fact #2: Goals of synthetic biology are evolving

 \Rightarrow Need for modular approaches

Our approach: Non-uniform *Boltzmann* random generation



Complex features

Robustness Predicted folding (2D/MFold, 3D/MCFold...) Stability (Molecular dynamics) Interactions (RNACofold, Docking) Non-redundant generation [5]

Posterior selection



AMIB EPI/AMIBio EI - Inria Saclay - Computational Biology Evaluation Seminar



Image : Lorenz et al, GCB'09

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!





Landscape : Lorenz et al, GCB'09, Structures : Varna - Darty et al, Bioinf. (2009)

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!





Landscape : Lorenz et al, GCB'09, Structures : Varna - Darty et al, Bioinf. (2009)

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!





Landscape : Lorenz et al, GCB'09, Structures : Varna - Darty et al, Bioinf. (2009)

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!





Landscape : Lorenz et al, GCB'09, Structures : Varna - Darty et al, Bioinf. (2009)

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!





Landscape : Lorenz et al, GCB'09, Structures : Varna - Darty et al, Bioinf. (2009)

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!





Landscape : Lorenz et al, GCB'09, Structures : Varna - Darty et al, Bioinf. (2009)

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!





Landscape : Lorenz et al, GCB'09, Structures : Varna - Darty et al, Bioinf. (2009)

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!





Landscape : Lorenz et al, GCB'09, Structures : Varna - Darty et al, Bioinf. (2009)

Assuming a thermodynamic equilibrium sometimes misrepresents the reality of RNA folding in finite time \rightarrow RNA kinetics!



Sampling for RNA 2D kinetics



- RNA 2D studies reveal dynamic behaviors (Markov process);
- Bottleneck: Build kinetic landscape (combinatorial explosion);

 Our project: Advanced sampling methods to approximate landscapes, enabling kinetics studies beyond 1k NTs.

 \Rightarrow ANR/FWF-funded RNALands project (2015 – 2019)

Partners: TBI Vienna, EPI Bonsai (Inria Lille), Paris-Sud

[Michalik, Touzet, Ponty, ECCB-ISMB 2017/Bioinformatics]



Sampling for RNA 2D kinetics



- RNA 2D studies reveal dynamic behaviors (Markov process);
- Bottleneck: Build kinetic landscape (combinatorial explosion);
- Our project: Advanced sampling methods to approximate landscapes, enabling kinetics studies beyond 1k NTs.

 \Rightarrow ANR/FWF-funded RNALands project (2015 – 2019)

Partners: TBI Vienna, EPI Bonsai (Inria Lille), Paris-Sud

[Michalik, Touzet, Ponty, ECCB-ISMB 2017/Bioinformatics]



RNA is (still) a nice testing ground for efficient sampling methods

- Discrete abstractions yield reasonable predictions
- ▶ ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ...fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions ...
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions ...
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



- RNA is (still) a nice testing ground for efficient sampling methods
- Discrete abstractions yield reasonable predictions
- ... from first principle, using robust energy models
- Despite exponential #structures, Poly-time algorithms...
- ... fast enough to be applied at a genomic/transcriptomic scales
- Provide handle into evolution of structure/(function?)
- Kinetics is the next frontier!

- Alain Denise & Wei Wang
- Bruno Sargueil & Delphine Allouche
- Jerome Waldispühl & Vlad Reinharz
- Andrea Tanzer
- Afaf Saaidi
- Juraj Michalik...



References I



R. Nussinov and A.B. Jacobson.

Fast algorithm for predicting the secondary structure of single-stranded RNA. *Proc Natl Acad Sci U S A*, 77:6903–13, 1980.

(nría_



Inría