Exploring Life through Logic Programming

Alessandro Dal Palú Andrea Formisano Agostino Dovier Enrico Pontelli

Dept. Computer Science, New Mexico State University, USA

University of Udine, Italy December 2015

Exploring Life through LP

4 B 6 4 B 6

< 🗇 🕨

Protein Structure Prediction

38 N

Proteins and Central Dogma



- The translation phase starts from a mRNA sequence and associates a protein sequence
- Proteins are made of amino acids (20 common different types)
- Amino acids are defined by letters $\{A, \ldots, Z\} \setminus \{B, J, O, U, X, Z\}$

・ 同 ト ・ ヨ ト ・ ヨ

Universal code



- The translation selects 3 RNA basis (codon) and associates 1 amino acid.
- The translation rules are encoded in the universal code.
- The code contains *stop* symbol and some redundant RNA triplets.

E. Pontelli (NMSU)

Exploring Life through LP

Proteins Amino acids

- Proteins are molecules made of a linear sequence of amino acids.
- Amino acids are combined through *peptide bond*.



→ E → < E</p>

A 10

Proteins Amino acids

- Proteins are molecules made of a linear sequence of amino acids.
- Amino acids are combined through *peptide bond*.



- The purple dots represent the *side chains*, that depend on the amino acid type
- Side chains have different shape, size, charge, polarity, etc.
- A side chain contains from 1 (Glycine) up to 18 (Tryptophan) atoms.

E. Pontelli (NMSU)

Exploring Life through LP

Proteins Amino acids



- There are 2 degrees of freedom (black arrows) for each amino acid
- A protein with *n* amino acids has 2*n* degrees of freedom (ignoring the side chains)
- Typical size range from 50 to 500 amino acids

Image: A math

Proteins

Primary sequence

Amino acid sequence: GPEILCGAELVDALQFVCGDRGFYFNKPTGYGSSS RRAPQTGIVDECCFRSCDLRRLEMYCAPLKPAKSA



Primary sequence...



...embedded in the 3D space

E. Pontelli (NMSU)

Exploring Life through LP

Udine, Dec. 21-23 2015 7 / 23

Proteins Tertiary structure





Primary sequence...

... tertiary structure

• Proteins FOLD spontaneously (when in their natural environment)

Exploring Life through LP

글 🕨 🖌 글

Proteins Facts

- Folding is consistent ⇒ same protein folds in the same way [Anfinsen74]
- Folding is fast $\Rightarrow 1\mu S 1mS$
- Driven by non covalent forces: electrostatic interactions, volume constraints, Hydrogen Bonding, van der Waals, Salt/disulfide Bridges
- Backbone is rigid

< ロ > < 同 > < 回 > < 回 > < 回 > <

Proteins Simplified structure

- The *backbone* (green) links consecutive amino acids
- Each amino acid can be identified by its carbon alpha (C_α)
- The distance between consecutive Cα points is constant (3.81Å).

・ 同 ト ・ ヨ ト ・ ヨ ト

Proteins

Backbone

- The distance between $C\alpha$ points is constant (3.81Å).
- Let us consider a backbone made of 4 amino acids.

There are two bend angles

And one torsional angle.

Udine, Dec. 21-23 2015

э

11/23

E. Pontelli (NMSU)

Exploring Life through LP

The structure prediction problem

- Given the primary structure of a protein (its amino acid sequence)
- For each amino acid, output its position in the space (tertiary structure of a protein)



B b 4

The structure prediction problem

... and this is the hard part:

- In nature a protein has a unique/stable 3D conformation
- A cost function (that mimics physics laws) can be used to score each conformation
- Searching for the optimal score produces the best candidate is difficult (NP-complete even in extremely simplified modelings)

The protein structure prediction problem

- In this presentation we have two simplifications:
- Protein model: only one atom per amino acid, and only 2 classes of amino acids (Hydrophobic and Polar)





.

Modeling

The protein structure prediction problem

- In this presentation we have two simplifications:
- Protein model: only one atom per amino acid, and only 2 classes of amino acids (Hydrophobic and Polar)
- Spatial model: 2D square lattice to represent amino acid positions





4 3 5 4

The protein structure prediction problem Model

- The input is a list *S* of amino acids $S = s_1, \ldots, s_n$, where $s_i \in \{h, p\}$
- Each *s_i* is placed on a 2D grid with integer coordinates
- Any pair of two amino acids cannot occupy the same position
- If two amino acids are at distance 1, they are in contact





・ 同 ト ・ ヨ ト ・ ヨ

The protein structure prediction problem Model

- Let: $next(\langle X_1, Y_1 \rangle, \langle X_2, Y_2 \rangle) \iff |X_1 X_2| + |Y_1 Y_2| = 1.$
- A folding is a function $\omega : \{1, \ldots, n\} \longrightarrow \mathbb{N}^2$ where:
 - For each $1 \le i < n$: $next(\omega(i), \omega(i+1))$
 - For each $1 \le i, j \le n$: $(i \ne j \rightarrow \omega(i) \ne \omega(j))$





Exploring Life through LP

Udine, Dec. 21-23 2015 16 / 23

-

The protein structure prediction problem Model

• Let: $next(\langle X_1, Y_1 \rangle, \langle X_2, Y_2 \rangle) \iff |X_1 - X_2| + |Y_1 - Y_2| = 1.$

• A folding is a function $\omega : \{1, ..., n\} \longrightarrow \mathbb{N}^2$ where: • For each $1 \le i < n$: • For each $1 \le i, j \le n$: $(i \ne j \rightarrow \omega(i) \ne \omega(j))$

• Find a folding that minimizes the (simplified) energy function:

$$E(S,\omega) = \sum_{\substack{1 \le i \le n-2\\i+2 \le j \le n}} \operatorname{Pot}(s_i, s_j) \cdot \operatorname{next}(\omega(i), \omega(j))$$

where Pot(p, p) = Pot(h, p) = Pot(p, h) = 0 and Pot(h, h) = -1.



Modeling

Model

- H: Yellow; P: Grey
- All foldings have energy -6







The protein structure prediction problem Complexity

 Establishing whether there is a folding with energy < k in the HP model and N² space is NP-complete

> (Crescenzi, Goldman, Papadimitriou, Piccolboni, Yannakakis. On the Complexity of Protein Folding. Journal of Computational Biology 5(3): 423-466 (1998))

The protein structure prediction problem ASP Model

- The positions can be restricted to $[1 \dots 2n] \times [1 \dots 2n]$
- The first two amino acids are set to $\omega(1) = \langle n, n \rangle$ and $\omega(2) = \langle n, n+1 \rangle$
- Sequence is represented by *n* facts:

```
prot(1,p). prot(2,h). prot(3,p).
prot(4,h). prot(5,p). prot(6,p).
prot(7,h). prot(8,p).
```

• Each placement is stored as sol(i, X, Y).

・ロット (四)・ (日)・ (日)・

- Count the protein length and defines the board size
- range is loaded with all combinations from 1 to 2N

・ロット (四)・ (日)・ (日)・

20/23

Udine, Dec. 21-23 2015

- (4) sol(1,N,N) :- size(N).
- (5) sol(2,N,N+1) :- size(N).

- (4) sol(1,N,N) :- size(N).
- (5) sol(2,N,N+1) :- size(N).
- (6) 1 { sol(I,X,Y) : range(X,Y) } 1 :- prot(I,Amino).
 - Each amino acid / is associated to a unique position

- (4) sol(1,N,N) :- size(N).
- (5) sol(2,N,N+1) :- size(N).
- (6) 1 { sol(I,X,Y) : range(X,Y) } 1 :- prot(I,Amino).
- - Constraint: two distinct amino acids cannot overlap

• Constraint: two consecutive amino acids in sequence are next in space

< 同 > < 三 > < 三

Modeling

- If two amino acids of type h are next, count one contact (energy_pair)
- Optimization: at least one amino acid in sequence between the pair
- Optimization: only odd distances in space can contribute to a contact

Modeling

Extensions

- Generalization to 3d cubic space is trivial.
- organization of the space have been explored
- The Face Centered Cube lattice models the discrete space in which the protein can fold.

- Two points are *connected* (next) iff $|x_i - x_j|^2 + |y_i - y_j|^2 +$
 - $|z_i-z_j|^2=2,$
- Each point has 12 neighbors

22/23

Some References

- R. Backofen and S. Will, A constraint-based approach to fast and exact structure prediction in 3-dimensional protein models, Constraints 11(1):5-30, 2006.
- A. Dal Palù, A. Dovier and F. Fogolari. Constraint logic programming approach to protein structure prediction, BMC Bioinformatics 5(186), 2004.
- A. Dal Palù, A. Dovier and E. Pontelli, A constraint solver for discrete lattices, its parallelization, and application to protein structure prediction, Software Practice and Experience 37(13):1405-1449, 2007. (COLA)
- A. Dal Palù, A. Dovier and E. Pontelli. Computing approximate solutions of the protein structure determination problem using global constraints on discrete crystal lattices, Int'l Journal of Data Mining and Bioinformatics 4(1):1–20, 2010. Also in WCB 06 and WCB 07
- P. Barahona and L. Krippahl, Constraint programming in structural bioinformatics, Constraints 13(1-2):3-20, 2008.
- A. Dovier. Recent constraint/logic programming based advances in the solution of the protein folding problem. Intelligenza Artificiale 5(1):113-117, 2011.
- Approximated results with local search and/or LNS by Hoos et al. and by Van Hentenryck et al.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >