# Counting Protein Structures
# by DFS with Dynamic Decomposition

Sebastian Will and Martin Mann

Albert-Ludwigs-University Freiburg
Bioinformatics at the Department of Computer Science

Workshop on
Constraint Based Methods for Bioinformatics 2006

'Counting Protein Structures
by DFS with Dynamic Decomposition'

### What we need to know:

- Protein Structures ?                                    ✘
- Prediction as CSP ?
- Counting by Decomposition ?

Counting
Protein
Structures by
DDFS

Overview

Lat. Proteins
HP-Model
Applications

The CSP
Idea
Modelling
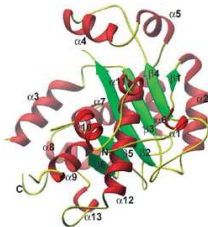Solving

Counting
DFS
Redundancy
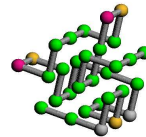Decomposition
DDFS

Results

# Protein structures

Real protein structure

Simple lattice proteins

$\Rightarrow$

## Lattice Proteins

- protein simplified to chain of monomers
- structure depends on underlying lattice
- energy depends on contact energy function

## Structure

- sequence monomers are placed on lattice positions
- structure = selfavoiding walk on the lattice



sequence = PHPHHPP

○ = H = hydrophic
● = P = polar

## Energy function

- focus on hydrophobic forces
- contact based → HH-contacts

$$\text{energy}(x,y) = \begin{cases} -1 & \text{if } (x,y == H) \text{ and } (x,y \text{ neighbors}) \\ 0 & \text{else} \end{cases}$$



energy = -2

⬭ = HH-contact

## Applications e.g.

- Neutral nets and protein evolution
- Exploring energy landscapes and protein kinetic
- Base for more complex protein models
- . . .

Therefore you need:

## Prediction of optimal structures

- NP-complete in 3D-lattice (Berger & Leighton, 1998) (even in 2D)
- can be solved by Constraint Programming !
  (Backofen & Will, 2006)

'Counting Protein Structures
by DFS with Dynamic Decomposition'

## What we need to know:

- Protein Structures ?                                    ✔
- Prediction as CSP ?                                     ✗
- Counting by Decomposition ?

## A CSP for optimal structure prediction in the HP-lattice-model



PHPHHP + H-core → CSP fomulation → X → solve CSP

Rolf Backofen and Sebastian Will
'A constraint-based approach to fast and exact structure prediction

in three-dimensional protein models' 2006

## H-Core of a given structure

- H-Core = set of H-monomer positions
- core energy $\leftrightarrow$ structure energy (only HH-contacts important)



- optimality implies optimal structure energy
- candidates can be precomputed based on H-number
- hard problem too $\rightarrow$ (solved via CP)

$\Rightarrow$ for now used as black box and given ... !

given

an HP-sequence

P-H-P-H-H-P-P

and

an optimal H-core



there is the question:

Exists a structure, so that all H-monomers are placed on H-core positions?

YES!

For a given HP-sequence and an optimal H-core:

## Variables

- one for each sequence monomer

## Domains = sets of lattice positions

- H-Monomers: H-core positions (ensures optimality)
- P-Monomers: remaining lattice

## Constraints

- binary Neighboring constraints along the chain (backbone)
- one global Alldifferent constraint (selfavoiding structure)

  ⇒ encodes the selfavoiding walk

Counting
Protein
Structures by
DDFS

Overview

Lat. Proteins
HP-Model
Applications

The CSP
Idea
Modelling
**Solving**

Counting
DFS
Redundancy
Decomposition
DDFS

Results

# Structure prediction as CSP

The CSP as Constraint Graph:



### From Solution to structure

- a CSP solution assigns a lattice position to each monomer
- solution = structure,     and optimal due to H-core !
- normal CSP-solving approaches can be applied
  e.g. DFS-branching combined with constraint propagation

'Counting Protein Structures
by DFS with Dynamic Decomposition'

## What we need to know:

- Protein Structures ?                                    ✔
- Prediction as CSP ?                                     ✔
- Counting by Decomposition ?                            ✘

## Counting all solutions

- is in complexity class #P-complete for counting problems
- = an important field of CP

- can be done by Constraint Propagation and DFS branching
- iterative process
- can be formulated as a recursion ...

1: **function** $\mathrm{CDFS}(\mathcal{X}, \mathcal{D}, \mathcal{C})$

▷ reduce domains

2:     $(\mathcal{D}', \mathcal{C}') \leftarrow \mathrm{PROPAGATE}(\mathcal{X}, \mathcal{D}, \mathcal{C})$

▷ check for recursion stop

3:     **if** $\mathrm{ISFAILED}(\mathcal{X}, \mathcal{D}', \mathcal{C}')$ **then return** 0
4:     **else if** $\mathrm{ISSOLVED}(\mathcal{X}, \mathcal{D}')$ **then return** 1
5:     **else**

▷ branch search

6:         $c \leftarrow \mathrm{SELECT}(\mathcal{X}, \mathcal{D}')$
7:         **return** $\mathrm{CDFS}(\mathcal{X}, \mathcal{D}', \mathcal{C}' \cup \{c\}) + \mathrm{CDFS}(\mathcal{X}, \mathcal{D}', \mathcal{C}' \cup \{\neg c\})$
8:     **end if**
9: **end function**

**redundant work**

## The problem

✗ Redundant work due to independent partial problems!

## How to avoid?

- detect decomposition into independent partial problems
- do it dynamically during search (CSP changes)
- solve each partial problem independently
- generate the overall solution number (via product)

## Counting using dynamic decomposing DFS

1: **function** $\mathrm{DDFS}(\mathcal{X}, \mathcal{D}, \mathcal{C})$

                                            ▷ reduce domains

2:     $(\mathcal{D}', \mathcal{C}') \leftarrow \mathrm{PROPAGATE}(\mathcal{X}, \mathcal{D}, \mathcal{C})$

                                 ▷ check for recursion stop

3:     **if** $\mathrm{ISFAILED}(\mathcal{X}, \mathcal{D}', \mathcal{C}')$ **then return** $0$

4:     **else if** $\mathrm{ISSOLVED}(\mathcal{X}, \mathcal{D}')$ **then return** $1$

                                ▷ decomposing branching

5:     **else**  $s \leftarrow 1$                   ▷ initialize counter

6:         $\mathfrak{D} \leftarrow \mathrm{DECOMPOSE}(\mathcal{X}, \mathcal{D}', \mathcal{C}')$

7:         **for all** $(\hat{\mathcal{X}}, \hat{\mathcal{D}}, \hat{\mathcal{C}}) \in \mathfrak{D}$ **do**

8:            $c \leftarrow \mathrm{SELECT}(\hat{\mathcal{X}}, \hat{\mathcal{D}})$

9:            $s = s \cdot \left( \mathrm{DDFS}(\hat{\mathcal{X}}, \hat{\mathcal{D}}, \hat{\mathcal{C}} \cup \{c\}) + \mathrm{DDFS}(\hat{\mathcal{X}}, \hat{\mathcal{D}}, \hat{\mathcal{C}} \cup \{\neg c\}) \right)$

10:         **end for**

11:         **return** $s$

12:     **end if**

13: **end function**

## DDFS

- DFS $\rightarrow$ redundant solving of independent partial problems
- can be avoided by DDFS via dynamic decomposition
- overall solutions are generated during backtracking
- $=$ general approach for exhaustive solution enumeration

## But . . .

- early and strong decompositions neccessary
- $\Rightarrow$ new requirements to variable and value selection
- $\Rightarrow$ problem specific heuristics important
- applicable for global constraints too! (using binarisation)
- first recursion draft can be improved (see paper ) ☺

'Counting Protein Structures
by DFS with Dynamic Decomposition'

### What we need to know:

- Protein Structures ?                                    ✔
- Prediction as CSP ?                                     ✔
- Counting by Decomposition ?                             ✔

OK, but does it help? Results . . .

Using DDFS for structure counting in HP-model

## Decomposing versus 'normal' DFS

- in a first implementation:
  - 10x less branchings
  - 2x faster

- possible speedup higher !!!
  - further algorithmic improvements
  - optimize implementation (e.g. constraint graph handling)
  - better problem specific branching heuristics

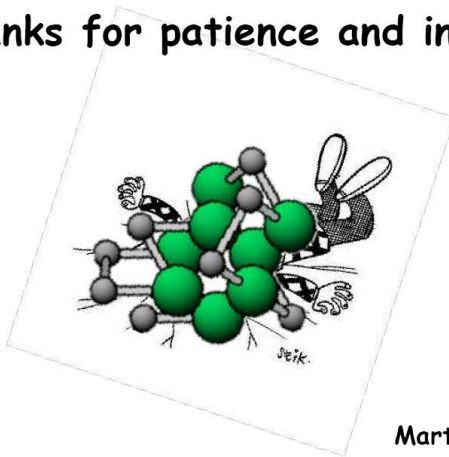'Counting Protein Structures
by DFS with Dynamic Decomposition'

Lets summarise:

1. optimal structure prediction for lattice proteins can be formulated as CSP

2. counting all solutions via DFS yields redundancy

3. can be avoided by dynamic decompositions $\Rightarrow$ DDFS

4. leads to big speedups and less branchings

5. DDFS is a general approach $\Rightarrow$ other CSPs

# Thanks for patience and interest



**Martin Mann**