



Optimal Structure Prediction and Application

Freiburg Bioinformatics Group

Martin Mann
Freiburg

Martin Mann, Sebastian Will and Rolf Backofen

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

EMBio Meeting Leipzig 2007



Overview

CPSP

Degeneracy

HPdesign

Summary



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The goal

- Prediction of the native structure given an AA-sequence

Assumptions

- Sequence determines structure
- Native structure has lowest energy

Problems

- too complex energy function
- too many degrees of freedom

LGGYMLGSA...RESQAYYQR



Human prion 1HJM



► Zoom in



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The goal

- Prediction of the native structure given an AA-sequence

Assumptions

- Sequence determines structure
- Native structure has lowest energy

Problems

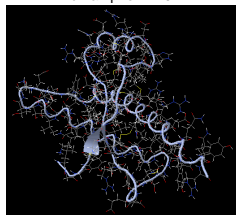
- too complex energy function
- too many degrees of freedom

Computationally not capable!

LGGYMLGSA...RESQAYYQR



Human prion 1HJM



▶ Zoom in



Simplified Off-Lattice Protein Models

One possible abstraction



Martin Mann
Freiburg

Overview

CPSP

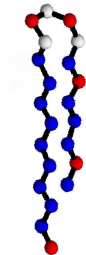
Degeneracy

HPdesign

Summary

Backbone Structure

C_α sequence only



Full 3D Space

all angles etc. allowed

Reduced Alphabet

e.g. HP, HPNX, ...

Contact Energy Function

e.g.

| | H | P |
|---|-----|------|
| H | -1 | 0.5 |
| P | 0.5 | -0.5 |



Simplified Off-Lattice Protein Models

One possible abstraction



Martin Mann
Freiburg

Overview

CPSP

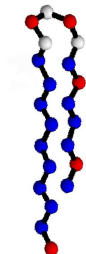
Degeneracy

HPdesign

Summary

Backbone Structure

C_{α} sequence only



Full 3D Space

all angles etc. allowed

Reduced Alphabet

e.g. HP, HPNX, ...

Contact Energy Function

e.g.

| | H | P |
|---|-----|------|
| H | -1 | 0.5 |
| P | 0.5 | -0.5 |

⇒ **Still too many degrees of freedom ...**



Simplified Lattice Protein Models

An other possible discrete abstraction



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

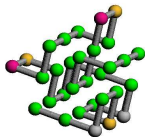
Summary

Backbone Structure

C_α sequence only

Reduced Alphabet

e.g. HP, HPNX, ...



Lattice Space

e.g. cubic, fcc, ...

Contact Energy Function

e.g.

| | H | P |
|---|----|---|
| H | -1 | 0 |
| P | 0 | 0 |



Simplified Lattice Protein Models

An other possible discrete abstraction



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

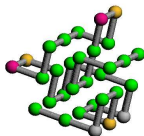
Summary

Backbone Structure

C_α sequence only

Reduced Alphabet

e.g. HP, HPNX, ...



Lattice Space

e.g. cubic, fcc, ...

Contact Energy Function

e.g.

| | H | P |
|---|----|---|
| H | -1 | 0 |
| P | 0 | 0 |

Discrete, Enumerable, Computationally Capable



Simple and nice... But what for?

Applications



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

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)





Simple and nice... But what for?

Applications



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

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)





Martin Mann
Freiburg

Overview

CPSP

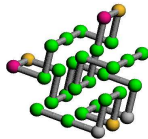
Degeneracy

HPdesign

Summary

Constraint-based Protein Structure Prediction¹

- **Short introduction**
 - **Optimal structure prediction (CPSP)**
- **Applications in**
 - **Protein stability**
 - **Inverse folding problem**



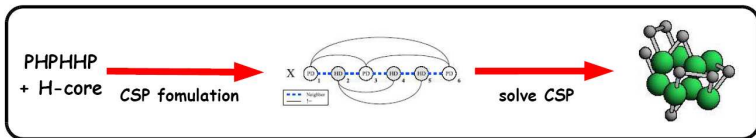
¹R. Backofen and S. Will



Martin Mann
Freiburg

Constraint-based Protein Structure Prediction (CPSP)

An Approach for optimal structure prediction
in the HP-lattice-model



Rolf Backofen and Sebastian Will

'A constraint-based approach to fast and exact structure prediction

in three-dimensional protein models' 2006



The CPSP Approach

The Main Idea



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The HP-Model

- simplest energy function available
- focus on hydrophobicity (hydrophobic cores)
- maximizing HH-contacts \leftrightarrow minimizing surface

CPSP - The central idea

- optimal H-monomer distribution is *sequence independent*
- precalculate such optimal and suboptimal so called *H-cores*
- try to find a mapping of a given sequence to H-cores
 \Rightarrow *Sequence-Threading*



Martin Mann
Freiburg

Overview

CPSP

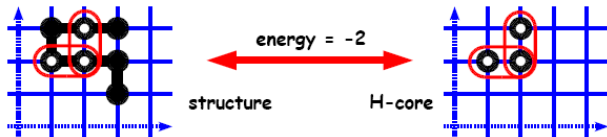
Degeneracy

HPdesign

Summary

H-Core of a given structure

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



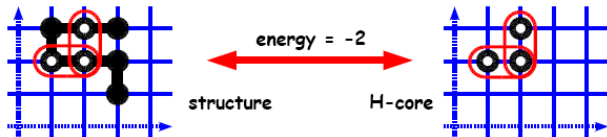
- optimality implies optimal structure energy
- candidates can be precomputed based on H-number
- hard problem too \rightarrow (solved via CP)

\Rightarrow for now used as black box and given in a DB ... !



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 in a DB ... !



The CPSP Approach

Sequence Threading : The Question to solve ...



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

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?



The CPSP Approach

Sequence Threading : The Question to solve ...



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

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?



The CPSP Approach

Sequence Threading : The Question to solve ...



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

Given

an HP-sequence

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

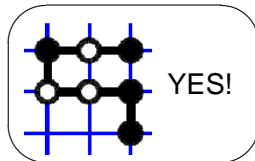
and

an optimal H-core



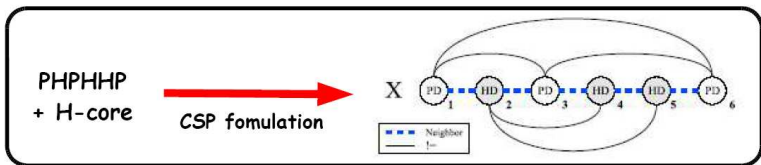
there is the question:

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





Modelling of the question as CSP



From Solution to structure

- fast standard CP-solvers can be applied for solving
- a CSP solution assigns a lattice position to each monomer
- solution = structure, and optimal due to H-core !
- usually a huge number of solutions / optimal structures



The CPSP Approach

Fast and very flexible ...



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The CPSP-Approach

- proven optimal structures via precalculated H-cores
- fast (first hit within seconds)
- deterministic (no stochastic structure space exploration)
- CP yields a very flexible, extensible modelling

Extensibility

- lattices (cubic, face centered cubic, ...)
- energy functions (HP, HPNX, ...)
- exclusion of symmetric solutions during enumeration
- solution space sampling via distance constraints (D-10)
- advanced CP-techniques for solution counting (D-9)
- ... future: side chain models, structure shapes, ...



The CPSP Approach

Fast and very flexible ...



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The CPSP-Approach

- proven optimal structures via precalculated H-cores
- fast (first hit within seconds)
- deterministic (no stochastic structure space exploration)
- CP yields a very flexible, extensibel modelling

Extensibility

- lattices (cubic, face centered cubic, ...)
- energy functions (HP, HPNX, ...)
- exclusion of symmetric solutions during enumeration
- solution space sampling via distance constraints (D-10)
- advanced CP-techniques for solution counting (D-9)
- ... future: side chain models, structure shapes, ...



The CPSP Approach

Summarizing



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The CPSP-Approach

- Prediction of optimal structures **without** folding simulation
- Exhaustive enumeration possible

Use

- Verifying folding simulation results
- Creation of interesting test sets for analysis
- Clustering of sequences into proteinlike or not
- Enumeration of the low energy part of the landscape
→ base for kinetic studies
- Base to solve other problems e.g. sequence evolution, inverse folding / designability



The CPSP Approach

Summarizing



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The CPSP-Approach

- Prediction of optimal structures **without** folding simulation
- Exhaustive enumeration possible

Use

- Verifying folding simulation results
- Creation of interesting test sets for analysis
- **Clustering of sequences into proteinlike or not**
- Enumeration of the low energy part of the landscape
→ base for kinetic studies
- Base to solve other problems e.g. sequence evolution,
inverse folding / designability



Martin Mann
Freiburg

Overview

CPSP

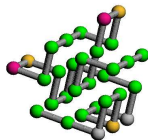
Degeneracy

HPdesign

Summary

Constraint-based Protein Structure Prediction

- **Short introduction**
 - **Optimal structure prediction (CPSP)**
- **Applications in**
 - **Protein stability**
 - **Inverse folding problem**





Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

Observations

- Not all possible AA-sequences used
- Fast folding process
(folding funnel hypothesis)
- Usually **one** stable, native structure
- ...

Human prion 1HJM





Degeneracy

A measure of protein stability



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

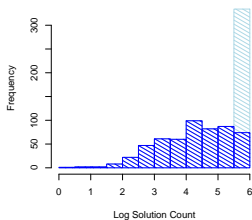
Summary

Degeneracy

- = the number of optimal structures
- important for protein stability
- deg. 1 = indicator for being stable

Degeneracy in Lattice-Models

- is high for most sequences
- due to simple energy function
- assumption: deg. 1 = stable



$\log_{10}(\text{ degeneracy })$ histogram
of 809 HP-sequences (H:P=1:1)
(32% with deg. $> 10^6$)

⇒ exact determinable via CPSP approach



Degeneracy

Summarizing



Martin Mann
Freiburg

Overview

CPSP

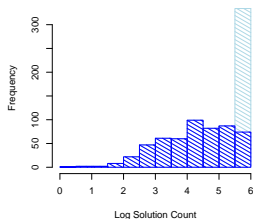
Degeneracy

HPdesign

Summary

Degeneracy

- = # of optimal structures
→ measure for protein stability
- exactly determinable using CPSP
- usually very high in HP-Model
- e.g. to distinguish proteinlike or random sequences
- counting can be improved ('06)



$\log_{10}(\text{ degeneracy })$ histogram
of 809 HP-sequences (H:P=1:1)
(32% with deg. > 10^6)



Martin Mann
Freiburg

Overview

CPSP

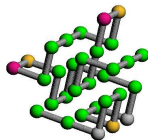
Degeneracy

HPdesign

Summary

Constraint-based Protein Structure Prediction

- Short introduction
 - Optimal structure prediction (CPSP)
- Applications in
 - Protein stability
 - Inverse folding problem





Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The goal

- design of sequences for a given structure

Sequence constraints

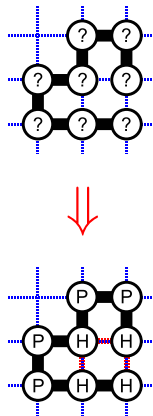
- proteinlike (low degeneracy)
- forms structure as optimal one

Problem

- # of sequences is exponentially in length

Addressed Questions (Designability)

- Is a structure X codeable?
- How many sequences code X ?





Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

The goal

- design of sequences for a given structure

Sequence constraints

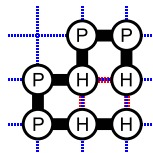
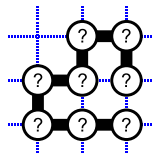
- proteinlike (low degeneracy)
- forms structure as optimal one

Problem

- # of sequences is exponentially in length

Addressed Questions (Designability)

- Is a structure X codeable?
- How many sequences code X ?



⇒ Solved using H-cores and CPSP ...



HPdesign

A 'Generate and Test' Workflow



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

Input

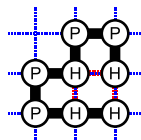
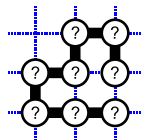
- a structure X

HPdesign workflow

- 1 Generate good candidates
- 2 Validate the sequences

Output

- a set of sequences that:
 - form X as their optimum
 - are stable (degeneracy 1)



HPPHPPHH
HPHHPHHH

...





Martin Mann
Freiburg

Overview

CSP

Degeneracy

HPdesign

Summary

Observation

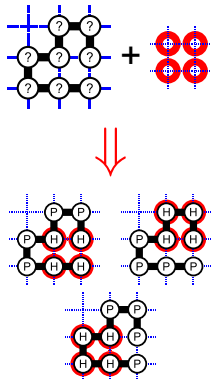
- optimal structure \Leftrightarrow optimal H-core

Generation

- take an arbitrary optimal H-core C
- shift C through structure X
- store resulting sequence for each hit

Result

- sequences with high chance to form X as their optimal structure





Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

Task

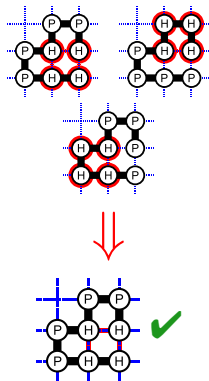
- Check for each sequences if it
 - is stable (degeneracy 1)
 - forms X as optimal structure

Workflow using CPSP

- 1: **for all** sequences S **do**
- 2: $\mathcal{X} \leftarrow \text{CPSP}(S, \text{max} = 2)$
- 3: **if** $(|\mathcal{X}| = 1 \wedge X \in \mathcal{X})$ **then**
- 4: STORE(S)
- 5: **end if**
- 6: **end for**

Result of Filtering

- stable S that form X as optimum





Martin Mann
Freiburg

Overview

CPSP

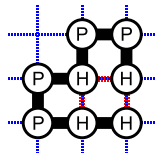
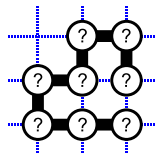
Degeneracy

HPdesign

Summary

The Inverse Folding Problem

- Design of stable sequences that form a given structure as optimum is a hard task
- Can be solved using 'Generate and Test'
- Candidate set can be shrinked (H-cores)
- Validation via CPSP approach possible





Outlook

Questions to answer ...



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

Open ...

- What distinguishes proteinlike and random sequences?
- Relations to real protein sequence properties?
- Are there common patterns in stable structures?
- How big are the bassins of attraction of stable optima?
- Leads one optimal structure to a folding funnel?
- What makes a structure designable?
- ...



Martin Mann
Freiburg

Overview

CPSP

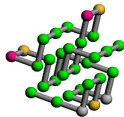
Degeneracy

HPdesign

Summary

Constraint-based Protein Structure Prediction

- **CPSP Approach**
 - Enumeration of optimal structures
 - Fast and extensible
- **Degeneracy**
 - Important measure of protein stability
- **Inverse folding problem**
 - Find stable sequences that form a structure X as optimum

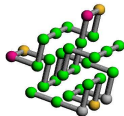




Constraint-based Protein Structure Prediction

CPSP-tools

- **HPstruct:** CPSP Approach
 - Enumeration of optimal structures
 - Fast and extensible
- **HPdeg:** Degeneracy
 - Important measure of protein stability
- **HPdesign:** Inverse folding problem
 - Find stable sequences that form a structure X as optimum
- ...





CPSP-tools

An implementation of CPSP and related methods



Martin Mann
Freiburg

Overview

CPSP

Degeneracy

HPdesign

Summary

CPSP-tools

- Implements the CPSP approach etc.
- Object oriented C++
- Library of core classes and functionality
- Completely documented / API
- Freely available

<http://www.bioinf.uni-freiburg.de/sw/cpcp/>



That's all folks!



Martin Mann
Freiburg

Overview

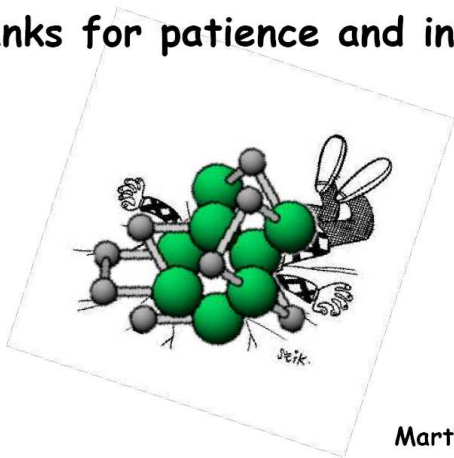
CPSP

Degeneracy

HPdesign

Summary

Thanks for patience and interest



Martin Mann



Martin Mann
Freiburg

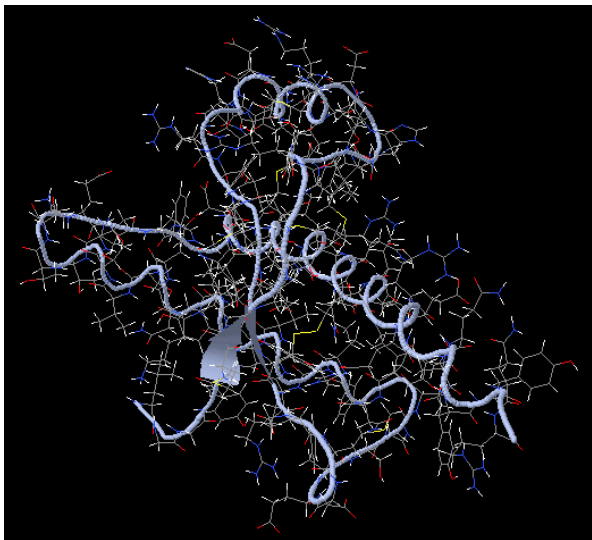
Appendix section



Protein folding and native structure prediction



Martin Mann
Freiburg



◀ Zoom out



Simplified Lattice Protein Models

Plus and Minus



Martin Mann
Freiburg

Advantages

- discrete \rightarrow full enumeration for result validation
- computationally capable
- folding dynamics similar to real proteins (time scale)
- unique folders

Possible Critics

- Energy function (HP) \Rightarrow HPNX, ...
- Lattice type (3D-cubic) \Rightarrow FCC, ...
- Lattice vs. angles \Rightarrow discrete angle model, ...

Return



The CSP

A very very simple Formulation



Martin Mann
Freiburg

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



The CSP

A very very simple Formulation

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



Martin Mann
Freiburg



The CSP

A very very simple Formulation



Martin Mann
Freiburg

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



The CSP

A very very simple Formulation



Martin Mann
Freiburg

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