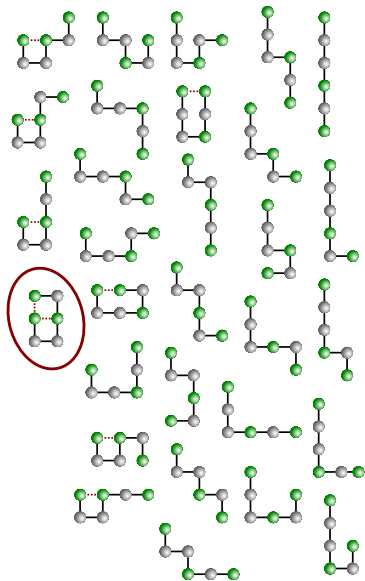
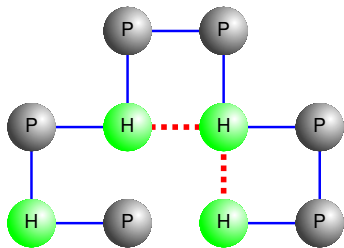


Application: Protein Structure Prediction



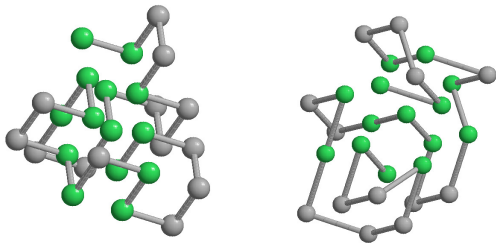
Exact Prediction in 3D cubic & FCC

The problem

IN: sequence s in $\{H, P\}^n$

HHPPPHHPHHPHHPHHPHPPHHPHPPHPPHH

OUT: self avoiding walk ω on cubic/fcc lattice with
minimal HP-energy $E_{HP}(s, \omega)$



A First Constraint Model

- Variables $X_1, \dots, X_n, Y_1, \dots, Y_n, Z_1, \dots, Z_n$ and $HHContacts$

$\begin{pmatrix} X_i \\ Y_i \\ Z_i \end{pmatrix}$ is the position of the i th monomer $\omega(i)$

- Domains

$$D(X_i) = D(Y_i) = D(Z_i) = \{-n, \dots, n\}$$

- Constraints

1. positions i and $i + 1$ are neighbored (**chain**)
2. all positions differ (**self-avoidance**)
3. relate $HHContacts$ to X_i, Y_i, Z_i

4. $\begin{pmatrix} X_1 \\ Y_1 \\ Z_1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$

A First Constraint Model

- Variables $X_1, \dots, X_n, Y_1, \dots, Y_n, Z_1, \dots, Z_n$ and $HHContacts$

$\begin{pmatrix} X_i \\ Y_i \\ Z_i \end{pmatrix}$ is the position of the i th monomer $\omega(i)$

- Domains

$$D(X_i) = D(Y_i) = D(Z_i) = \{-n, \dots, n\}$$

- Constraints

- positions i and $i + 1$ are neighbored (**chain**)
- all positions differ (**self-avoidance**)
- relate $HHContacts$ to X_i, Y_i, Z_i

4. $\begin{pmatrix} X_1 \\ Y_1 \\ Z_1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$

A First Constraint Model

- Variables $X_1, \dots, X_n, Y_1, \dots, Y_n, Z_1, \dots, Z_n$ and $HHContacts$

$\begin{pmatrix} X_i \\ Y_i \\ Z_i \end{pmatrix}$ is the position of the i th monomer $\omega(i)$

- Domains

$$D(X_i) = D(Y_i) = D(Z_i) = \{-n, \dots, n\}$$

- Constraints

- positions i and $i + 1$ are neighbored (**chain**)
- all positions differ (**self-avoidance**)
- relate $HHContacts$ to X_i, Y_i, Z_i

- $\begin{pmatrix} X_1 \\ Y_1 \\ Z_1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$

The First Model in More Detail (Cubic Lattice)

The Constraints cannot be expressed directly, i.e. we need auxiliary variables

$$Xdiff_{ij} = |X_i - X_j| \quad Ydiff_{ij} = |Y_i - Y_j| \quad Zdiff_{ij} = |Z_i - Z_j|$$

1. Positions i and $i + 1$ neighbored (**chain**)

$$Xdiff_{i(i+1)} + Ydiff_{i(i+1)} + Zdiff_{i(i+1)} = 1$$

2. All positions differ (**self-avoidance**)

$$Xdiff_{ij} + Ydiff_{ij} + Zdiff_{ij} \neq 0 \quad (\text{for } i \neq j).$$

3. Relate *HHContacts* to X_i, Y_i, Z_i

Detect HH-contact, if $Xdiff_{ij} + Ydiff_{ij} + Zdiff_{ij} = 1$ for $s_i = s_j = H$. Then add 1 to *HHContacts*.

(Technically, use **reified constraints**)

The First Model in More Detail (Cubic Lattice)

The Constraints cannot be expressed directly, i.e. we need auxiliary variables

$$Xdiff_{ij} = |X_i - X_j| \quad Ydiff_{ij} = |Y_i - Y_j| \quad Zdiff_{ij} = |Z_i - Z_j|$$

1. Positions i and $i + 1$ neighbored (**chain**)

$$Xdiff_{i(i+1)} + Ydiff_{i(i+1)} + Zdiff_{i(i+1)} = 1$$

2. All positions differ (**self-avoidance**)

$$Xdiff_{ij} + Ydiff_{ij} + Zdiff_{ij} \neq 0 \quad (\text{for } i \neq j).$$

3. Relate *HHContacts* to X_i, Y_i, Z_i

Detect HH-contact, if $Xdiff_{ij} + Ydiff_{ij} + Zdiff_{ij} = 1$ for $s_i = s_j = H$. Then add 1 to *HHContacts*.

(Technically, use **reified constraints**)

The First Model in More Detail (Cubic Lattice)

The Constraints cannot be expressed directly, i.e. we need auxiliary variables

$$Xdiff_{ij} = |X_i - X_j| \quad Ydiff_{ij} = |Y_i - Y_j| \quad Zdiff_{ij} = |Z_i - Z_j|$$

1. Positions i and $i + 1$ neighbored (**chain**)

$$Xdiff_{i(i+1)} + Ydiff_{i(i+1)} + Zdiff_{i(i+1)} = 1$$

2. All positions differ (**self-avoidance**)

$$Xdiff_{ij} + Ydiff_{ij} + Zdiff_{ij} \neq 0 \quad (\text{for } i \neq j).$$

3. Relate *HHContacts* to X_i, Y_i, Z_i

Detect HH-contact, if $Xdiff_{ij} + Ydiff_{ij} + Zdiff_{ij} = 1$ for $s_i = s_j = H$. Then add 1 to *HHContacts*.

(Technically, use **reified constraints**)

Solving the First Model

- Model is a COP (Constraint Optimization Problem)
- Branch and Bound Search for Minimizing *Energy*
- Combined with Symmetry Breaking
- How good is the propagation?
- Main problem of propagation: bounds on contacts/energy
From a partial solution, the solver cannot estimate the maximally possible number of HH-contacts well.

Solving the First Model

- Model is a COP (Constraint Optimization Problem)
- Branch and Bound Search for Minimizing *Energy*
- Combined with Symmetry Breaking
- How good is the propagation?
- Main problem of propagation: bounds on contacts/energy
From a partial solution, the solver cannot estimate the maximally possible number of HH-contacts well.

Solving the First Model

- Model is a COP (Constraint Optimization Problem)
- Branch and Bound Search for Minimizing *Energy*
- Combined with Symmetry Breaking
- How good is the propagation?
- Main problem of propagation: bounds on contacts/energy
From a partial solution, the solver cannot estimate the maximally possible number of HH-contacts well.

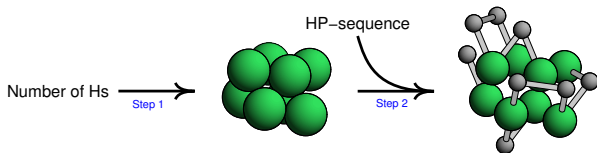
Solving the First Model

- Model is a COP (Constraint Optimization Problem)
- Branch and Bound Search for Minimizing *Energy*
- Combined with Symmetry Breaking
- How good is the propagation?
- Main problem of propagation: bounds on contacts/energy
From a partial solution, the solver cannot estimate the maximally possible number of HH-contacts well.

Solving the First Model

- Model is a COP (Constraint Optimization Problem)
- Branch and Bound Search for Minimizing *Energy*
- Combined with Symmetry Breaking
- How good is the propagation?
- Main problem of propagation: bounds on contacts/energy
From a partial solution, the solver cannot estimate the maximally possible number of HH-contacts well.

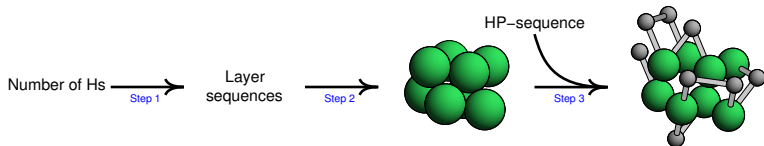
The Advanced Approach: Cubic & FCC



Steps

1. Core Construction
2. Mapping

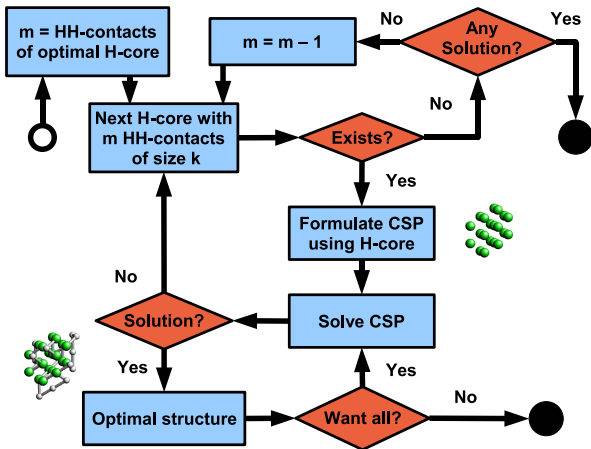
The Advanced Approach: Cubic & FCC



Steps

1. Bounds
2. Core Construction
3. Mapping

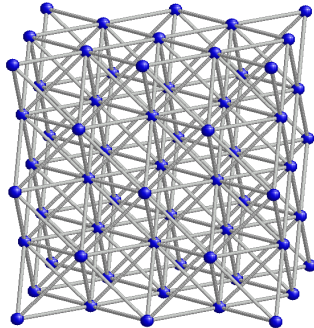
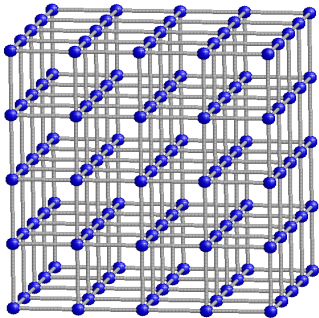
Workflow: Predict Best Structure(s) of HP-Sequence



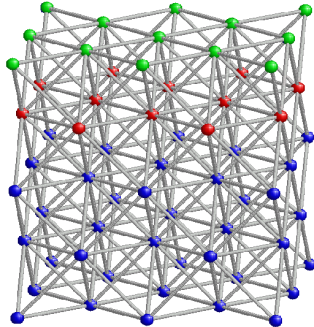
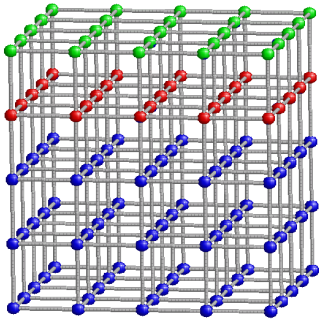
Computing Bounds

- Prepares the construction of cores
- How many contacts are possible for n monomers, if freely distributed to lattice points
- Answering the question will give information for core construction
- Main idea: split lattice into layers
consider contacts
 - within layers
 - between layers

Layers: Cubic & FCC Lattice



Layers: Cubic & FCC Lattice

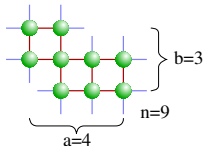


Contacts

Contacts =

Layer contacts + Contacts between layers

- Bound **Layer contacts**: $\text{Contacts} \leq 2 \cdot n - a - b$



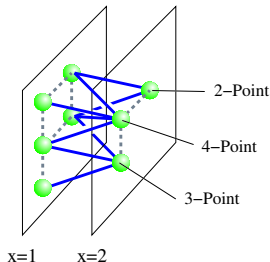
- Bound **Contacts between layers**

- cubic: **one** neighbor in next layer

$$\text{Contacts} \leq \min(n_1, n_2)$$

- FCC: **four** neighbors in next layer

i – points



Bounding Interlayer Contacts in the FCC

- **Needed:**

- upper bound for number of contacts between two successive layers in FCC
- NOTE: Layers only described by parameters $(n_1, a_1, b_1); (n_2, a_2, b_2)$

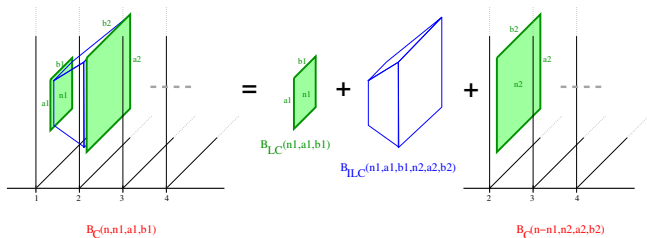
- **Method:**

- compute bounds for number of 1/2/3/4-points of first layer
- distribute n_2 points greedily
- technical difficulty: tight bounds of 1/2/3/4-points depend on further parameters

- **Result:** $B_{\text{ILC}}^{\text{FCC}}(n_1, a_1, b_1, n_2, a_2, b_2)$

Recall: $B_{\text{ILC}}^{\text{cubic}}(n_1, a_1, b_1, n_2, a_2, b_2) = \min(n_1, n_2)$

Recursion Equation for Bounds

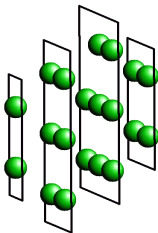


- $B_C(n, n_1, a_1, b_1)$: Contacts of core with n elements and first layer L_1 : n_1, a_1, b_1
- $B_{LC}(n_1, a_1, b_1)$: Contacts in L_1
- $B_{ILC}(n_1, a_1, b_1, n_2, a_2, b_2)$: Contacts between E_1 and E_2 : n_2, a_2, b_2
- $B_C(n - n_1, n_2, a_2, b_2)$: Contacts in core with $n - n_1$ elements and first layer E_2

Layer sequences

From Recursion:

- by Dynamic Programming: **Upper bound on number of contacts**
- by Traceback: **Set of layer sequences**



layer sequence = $(n_1, a_1, b_1), \dots, (n_4, a_4, b_4)$

Set of layer sequences gives distribution of points to layers in all point sets that possibly have maximal number of contacts

Core Construction

Problem

IN: number n , contacts c

OUT: all point sets of size n with c contacts

- Optimization problem
- Core construction is a hard combinatorial problem

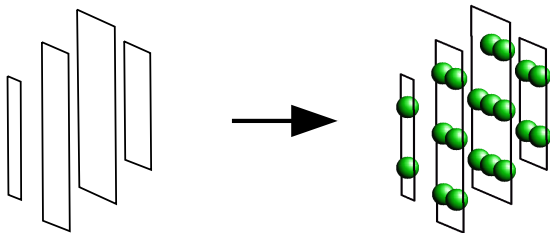
Core construction: Modified Problem

Problem

IN: number n , contacts c , set of layer sequences S_{ls}

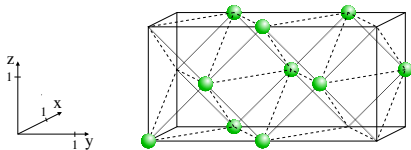
OUT: all point sets of size n with c contacts and layer sequences in S_{ls}

- Use constraints from layer sequences
- Model as **constraint satisfaction problem (CSP)**



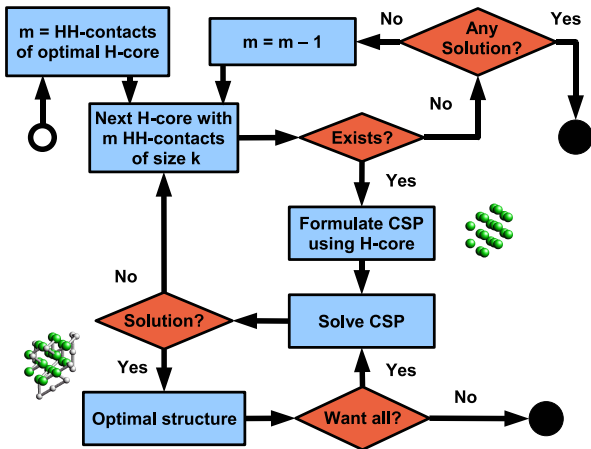
$(n_1, a_1, b_1), \dots, (n_4, a_4, b_4)$ Core = Set of lattice points

Core Construction — Details



- Number of layers = length of layer sequence
- Number of layers in x , y , and z : Surrounding Cube
- enumerate numbers of layers \Rightarrow fix cube \Rightarrow enumerate points

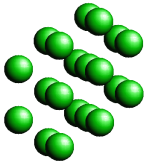
Workflow



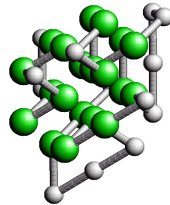
Mapping Sequences to Cores

find structure such that

- H-Monomers on core positions → hydrophobic core
- all positions differ → self-avoiding
- chain connected → walk



compact core



optimal structure

Mapping Sequence to Cores — CSP

Given: sequence s of size n and n_H H s
core $Core$ of size n_H

CSP Model

- Variables X_1, \dots, X_n
 X_i is position of monomer i
Encode positions as integers

$$I \begin{pmatrix} x \\ y \\ z \end{pmatrix} \equiv x + M * y + M^2 * z$$

(unique encoding for 'large enough' M)

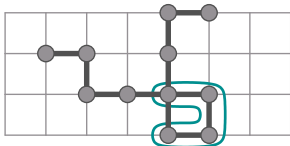
- Constraints
 1. $X_i \in Core$ for all $s_i = H$
 2. X_i and X_{i+1} are neighbors
 3. X_1, \dots, X_n are all different

Constraints for Self-avoiding Walks

- Single Constraints “self-avoiding” and “walk” weaker than their combination
- no efficient algorithm for consistency of combined constraint “self-avoiding walk”
- relaxed combination: stronger and more efficient propagation

k-avoiding walk constraint

Example: 4-avoiding, but not 5-avoiding



Putting it together

Predict optimal structures by combining the three steps

1. Bounds
2. Core Construction
3. Mapping

Some Remarks

- Pre-compute optimal cores for relevant core sizes
Given a sequence, only perform Mapping step
- Mapping to cores may fail!
We use suboptimal cores and iterate mapping.
- Approach extensible to HPNX
HPNX-optimal structures at least nearly optimal for HP.
- Approach extensible to side chains
H side chains form core.

Putting it together

Predict optimal structures by combining the three steps

1. Bounds
2. Core Construction
3. Mapping

Some Remarks

- Pre-compute optimal cores for relevant core sizes
Given a sequence, only perform Mapping step
- Mapping to cores may fail!
We use suboptimal cores and iterate mapping.
- Approach extensible to HPNX
HPNX-optimal structures at least nearly optimal for HP.
- Approach extensible to side chains
H side chains form core.

Putting it together

Predict optimal structures by combining the three steps

1. Bounds
2. Core Construction
3. Mapping

Some Remarks

- Pre-compute optimal cores for relevant core sizes
Given a sequence, only perform Mapping step
- Mapping to cores may fail!
We use suboptimal cores and iterate mapping.
- Approach extensible to HPNX
HPNX-optimal structures at least nearly optimal for HP.
- Approach extensible to side chains
H side chains form core.

Putting it together

Predict optimal structures by combining the three steps

1. Bounds
2. Core Construction
3. Mapping

Some Remarks

- Pre-compute optimal cores for relevant core sizes
Given a sequence, only perform Mapping step
- Mapping to cores may fail!
We use suboptimal cores and iterate mapping.
- Approach extensible to HPNX
HPNX-optimal structures at least nearly optimal for HP.
- Approach extensible to side chains
H side chains form core.

Time efficiency

Prediction of **one** optimal structure

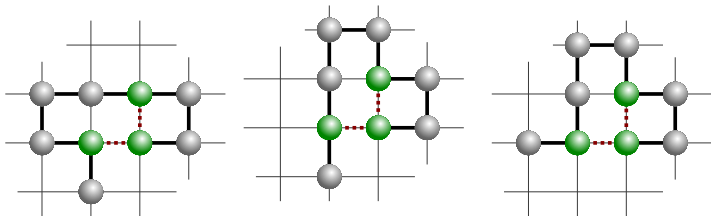
(“Harvard Sequences”, length 48 [Yue *et al.*, 1995])

CPSP	PERM
0,1 s	6,9 min
0,1 s	40,5 min
4,5 s	100,2 min
7,3 s	284,0 min
1,8 s	74,7 min
1,7 s	59,2 min
12,1 s	144,7 min
1,5 s	26,6 min
0,3 s	1420,0 min
0,1 s	18,3 min

- **CPSP**: “our approach”, constraint-based
- **PERM** [Bastolla *et al.*, 1998]: stochastic optimization

Many Optimal Structures

Sequence HPPHPPPHP



...?

- There can be many ...
- HP-model is **degenerated**
- Number of optimal structures = **degeneracy**

Completeness

Predicted number of **all** optimal structures
("Harvard Sequences")

CPSP	CHCC
10.677.113	1500×10^3
28.180	14×10^3
5.090	5×10^3
1.954.172	54×10^3
1.868.150	52×10^3
106.582	59×10^3
15.926.554	306×10^3
2.614	1×10^3
580.751	188×10^3

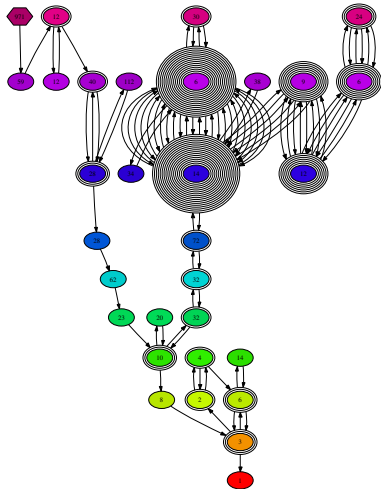
- **CPSP**: "our approach"
- **CHCC** [Yue *et al.*, 1995]: complete search with hydrophobic cores

Unique Folder

- HP-model degenerated
- Low degeneracy \approx stable \approx protein-like
- Are there protein-like, unique folder in 3D HP models?
- How to find out?

Unique Folder

- HP-model degenerated
- Low degeneracy \approx stable \approx protein-like
- Are there protein-like, unique folder in 3D HP models?
- How to find out?

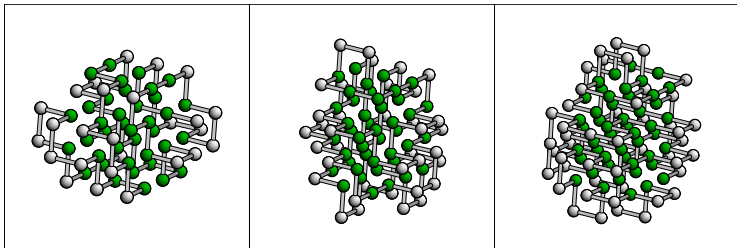


MC-search through sequence space

Unique Folder

- HP-model degenerated
- Low degeneracy \approx stable \approx protein-like
- Are there protein-like, unique folder in 3D HP models?
- How to find out?

Yes: many, e.g. about 10,000 for $n=27$



Software: CPSP Tools

<http://csp.informatik.uni-freiburg.de>

CPSP Tools

Menu

[Home](#)

[HPstruct](#)

structure pred.

[HPconvert](#)

PDB, CML, ...

[HPview](#)

3D visualization

[HPdeg](#)

degeneracy

[HPnnet](#)

neutral network

[HPdesign](#)

seq. design

[LatFit](#)

PDB to lattice

[Results](#)

direct access

[Help](#)

[FAQ](#)

CPSP Tools

Constraint-based Protein Structure Prediction

[Bioinformatics Group](#)

[Albert-Ludwigs-University Freiburg](#)

web-tools version 1.1.1 (06.04.2011)

The CPSP-tools package provides programs to solve exactly and completely the problems typical of studies using 3D lattice protein models. Among the tasks addressed are the prediction of globally optimal and/or suboptimal structures as well as sequence design and neutral network exploration.

Choose a tool from the left for ad hoc usage

(CPSP-tools version 2.4.2) (LatPack version 1.7.2)

or

Download the full [CPSP-tools](#) or [LatPack package](#) for local usage!

If you use the CPSP-tools please cite the following publications:

- Martin Mann, Sebastian Will, and Rolf Backofen.

[CPSP-tools - Exact and Complete Algorithms for High-throughput 3D Lattice Protein Studies.](#)

In *BMC Bioinformatics* 9: 230, 2008.