Computational Protein Design Using AND/OR Branch-and-Bound Search

Yichao Zhou Yuexing Zhou Jianyang Zeng

Institute for Interdisciplinary Information Sciences Tsinghua University

Apr, 2015

Yichao Zhou, Yuexing Zhou, Jianyang Zeng Computational Protein Design Using AND/OR Branch-and-Bound

What is Structure-Based Protein Design?





Protein structure template

1D amino acid sequence

Yichao Zhou, Yuexing Zhou, Jianyang Zeng Computational Protein Design Using AND/OR Branch-and-Bound

Why We Need Protein Design? Applications



Figure: New Drug Discovery



Figure: Enzyme Optimization



Figure: Drug Resistance Prediction



Figure: New Biosensor Design

Protein Design as an Optimization Problem



Protein Design as an Optimization Problem



Search Algorithm



Yichao Zhou, Yuexing Zhou, Jianyang Zeng Computational Protein Design Using AND/OR Branch-and-Bound

Search Algorithm



Yichao Zhou, Yuexing Zhou, Jianyang Zeng Computational Protein Design Using AND/OR Branch-and-Bound



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree rooted at x.

A (B) > A (B) > A (B)



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.



Heuristic function

Heuristic function h(x) returns a lower bound of the energy of leaf nodes in the subtree root at x.

Approximation

Simplify the energy model by only considering the interaction between residues whose energy term is greater than a threshold.













AND/OR Search Space





(a) Residue Interaction Network

(b) AND/OR Search Space

イロト イポト イヨト イヨ

Node value

$$v_x = \begin{cases} \sum_{y \in \text{child}(x)} v_y, & \text{if } x \text{ is an AND node;} \\ \min_{y \in \text{child}(x)} e(y) + v_y, & \text{if } x \text{ is an OR node.} \end{cases}$$

energy terms between y and ancestors of y

AND/OR Search Space





(a) Residue Interaction Network

(b) AND/OR Search Space

伺い イヨト イヨト

Node value

$$v_x = \begin{cases} \sum_{y \in \text{child}(x)} v_y, & \text{if } x \text{ is an AND node;} \\ \min_{y \in \text{child}(x)} e(y) + v_y, & \text{if } x \text{ is an OR node.} \end{cases}$$

1 y 51 IIS DCL ancest

AND/OR Search Space





(a) Residue Interaction Network

(b) AND/OR Search Space

・ 同 ト ・ ヨ ト ・ ヨ ト

Node value

$$v_x = \begin{cases} \sum_{y \in \text{child}(x)} v_y, & \text{if } x \text{ is an AND node;} \\ \min_{y \in \text{child}(x)} e(y) + v_y, & \text{if } x \text{ is an OR node.} \end{cases}$$

energy terms between y and ancestors of y

• AND/OR branch-and-bound search can find GMEC efficiently

- Protein design models are not perfect:
 - errors in energy functions and structure of protein templates
 - assumptions of rigid backbone and discrete side-chain conformations
- The GMEC solution is not sufficient in practice
- Sub-optimal (e.g., *second-best*, ...) conformations are required

Finding the *k* best solutions

- add a threshold to prevent from pruning sub-optimal solutions
- maintain k node values for each node, i.e., array
 v_x[1], v_x[2],..., v_x[k]

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ ・

- AND/OR branch-and-bound search can find GMEC efficiently
- Protein design models are not perfect:
 - errors in energy functions and structure of protein templates
 - assumptions of rigid backbone and discrete side-chain conformations
- The GMEC solution is not sufficient in practice
- Sub-optimal (e.g., *second-best*, ...) conformations are required

Finding the *k* best solutions

- add a threshold to prevent from pruning sub-optimal solutions
- maintain k node values for each node, i.e., array
 v_x[1], v_x[2],..., v_x[k]

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ ・

- AND/OR branch-and-bound search can find GMEC efficiently
- Protein design models are not perfect:
 - errors in energy functions and structure of protein templates
 - assumptions of rigid backbone and discrete side-chain conformations
- The GMEC solution is not sufficient in practice
- Sub-optimal (e.g., second-best, ...) conformations are required

Finding the *k* best solutions

- add a threshold to prevent from pruning sub-optimal solutions
- maintain k node values for each node, i.e., array
 v_x[1], v_x[2],..., v_x[k]

ヘロン 人間 とくほ とくほ とう

- AND/OR branch-and-bound search can find GMEC efficiently
- Protein design models are not perfect:
 - errors in energy functions and structure of protein templates
 - assumptions of rigid backbone and discrete side-chain conformations
- The GMEC solution is not sufficient in practice
- Sub-optimal (e.g., second-best, ...) conformations are required

Finding the *k* best solutions

- add a threshold to prevent from pruning sub-optimal solutions
- maintain k node values for each node, i.e., array
 v_x[1], v_x[2],..., v_x[k]

ヘロン 人間 とくほ とくほ とう

- AND/OR branch-and-bound search can find GMEC efficiently
- Protein design models are not perfect:
 - errors in energy functions and structure of protein templates
 - assumptions of rigid backbone and discrete side-chain conformations
- The GMEC solution is not sufficient in practice
- Sub-optimal (e.g., second-best, ...) conformations are required

Finding the *k* best solutions

- add a threshold to prevent from pruning sub-optimal solutions
- maintain k node values for each node, i.e., array
 v_x[1], v_x[2],..., v_x[k]

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ ・



- Find *k* best rotamer assignments for *x*;
- Sort all the *v* values for the child nodes of *x*;
- Keep k smallest node values



- Find *k* best rotamer assignments for *x*;
- Sort all the *v* values for the child nodes of *x*;
- Keep k smallest node values

A (B) > A (B) > A (B)



- Find *k* best rotamer assignments for *x*;
- Sort all the *v* values for the child nodes of *x*;
- Keep k smallest node values

A (1) > A (1) > A



• Want to find *k* best ways to merge the sub-problems;

- $v_x[i]$ should be the sum of $v_a[i_a] + v_b[i_b] + v_c[i_c]$;
- Obviously $v_x[1] = v_a[1] + v_b[1] + v_c[1]$, hard for $v_x[2]$ and $v_x[3]$;
- Brute force algorithm: $O(k^d)$
- Our algorithm: $O(dk \log(dk))$
 - Develop a new method to reduce the complexity exponentially
 - Use binary heap as the data-structure for efficiency



- Want to find *k* best ways to merge the sub-problems;
- $v_x[i]$ should be the sum of $v_a[i_a] + v_b[i_b] + v_c[i_c]$;
- Obviously $v_x[1] = v_a[1] + v_b[1] + v_c[1]$, hard for $v_x[2]$ and $v_x[3]$;
- Brute force algorithm: $O(k^d)$
- Our algorithm: $O(dk \log(dk))$
 - Develop a new method to reduce the complexity exponentially
 - Use binary heap as the data-structure for efficiency



- Want to find *k* best ways to merge the sub-problems;
- $v_x[i]$ should be the sum of $v_a[i_a] + v_b[i_b] + v_c[i_c]$;
- Obviously $v_x[1] = v_a[1] + v_b[1] + v_c[1]$, hard for $v_x[2]$ and $v_x[3]$;
- Brute force algorithm: $O(k^d)$
- Our algorithm: $O(dk \log(dk))$
 - Develop a new method to reduce the complexity exponentially
 - Use binary heap as the data-structure for efficiency

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ ・



- Want to find *k* best ways to merge the sub-problems;
- $v_x[i]$ should be the sum of $v_a[i_a] + v_b[i_b] + v_c[i_c]$;
- Obviously $v_x[1] = v_a[1] + v_b[1] + v_c[1]$, hard for $v_x[2]$ and $v_x[3]$;
- Brute force algorithm: $O(k^d)$
- Our algorithm: O(dk log(dk))
 - Develop a new method to reduce the complexity exponentially
 - Use binary heap as the data-structure for efficiency

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ ・



- Want to find *k* best ways to merge the sub-problems;
- $v_x[i]$ should be the sum of $v_a[i_a] + v_b[i_b] + v_c[i_c]$;
- Obviously $v_x[1] = v_a[1] + v_b[1] + v_c[1]$, hard for $v_x[2]$ and $v_x[3]$;
- Brute force algorithm: $O(k^d)$
- Our algorithm: $O(dk \log(dk))$
 - Develop a new method to reduce the complexity exponentially
 - Use binary heap as the data-structure for efficiency

• □ • • □ • • □ • • □ •

Experiment

- OSPREY platform from *Donald Lab* (Duke University)
- daoopt AOBB framework from UC Irvine
- Native sequence recovery (core redesign)
- Large rotamer library for better accuracy
- DEE/A* as the comparison

Environment Information

CPU: Intel Xeon[™] E5-1620 3.6GHz

Memory: 4G Memory Limitation

Experiment

- OSPREY platform from *Donald Lab* (Duke University)
- daoopt AOBB framework from UC Irvine
- Native sequence recovery (core redesign)
- Large rotamer library for better accuracy
- DEE/A* as the comparison

Environment Information

CPU: Intel Xeon[™] E5-1620 3.6GHz

Memory: 4G Memory Limitation

Dataset

- 23 previously unsolvable protein redesign cases
- 5 previously solvable protein redesign cases
- 6 full protein redesign cases

Results

- solve 21 out of 23 problems that was previously unsolvable;
- provide a large speedup by several orders of magnitude for previously solvable problems;
- solve all full protein redesign problems while DEE/A* exceeds 4G memory limitation.

Dataset

- 23 previously unsolvable protein redesign cases
- 5 previously solvable protein redesign cases
- 6 full protein redesign cases

Results

- solve 21 out of 23 problems that was previously unsolvable;
- provide a large speedup by several orders of magnitude for previously solvable problems;
- solve all full protein redesign problems while DEE/A* exceeds 4G memory limitation.

Dataset

- 23 previously unsolvable protein redesign cases
- 5 previously solvable protein redesign cases
- 6 full protein redesign cases

Results

- solve 21 out of 23 problems that was previously unsolvable;
- provide a large speedup by several orders of magnitude for previously solvable problems;
- solve all full protein redesign problems while DEE/A* exceeds 4G memory limitation.

Running time of AND/OR branch-and-bound algorithm				
Running time of our optimized OSPREY				
Running time of traditional OSPREY				
Conformation Space				
PDB	Space size	OŜPREY	cOSPREY	AOBB
1IQZ	7.11e+17	1,824,235	40,217	117
2C0V	1.14e+10	317	21	1
3FGV	6.44e+12	59,589	5,091	< 1
3DN J	5.11e+12	7,469	570	3
2FHZ	1.83e+18	3,475,716	70,783	13

Time is measured in millisecond.

Conclusion and Future Work

Our Contribution

- Apply AND/OR branch-and-bound search to protein design
- Design an algorithm to find sub-optimal solutions
- Test AOBB search in native sequence recovery experiment

Results

- Memory efficient
- Speedup by several orders of magnitude comparing to A*/DEE

Future Work

- Using AND/OR best-first search for sub-optimal solutions
- Parallelizing AOBB on a super computer

Discussion

- Prof. Alexander Ihler
- Dr. Lars Otten

Funding

- National Basic Research Program of China
- National Natural Science Foundation of China
- China's Youth 1000-Talent Program.

Thank you!

Yichao Zhou, Yuexing Zhou, Jianyang Zeng Computational Protein Design Using AND/OR Branch-and-Bound

▲■▼ ▲ 国 ▼ ▲ 国 ▼