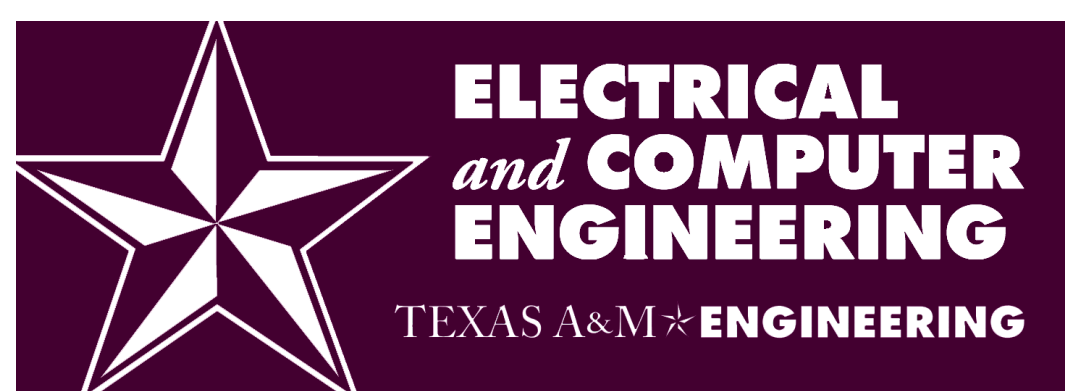


Interconnected Cost Function Networks (iCFN): an efficient exact algorithm for multistate protein design

MOSTAFA KARIMI^{1,2}, YANG SHEN^{1,2}

¹Department of Electrical and Computer Engineering, ²TEES-AgriLife Center for Bioinformatics and Genomic Systems Engineering, Texas A&M University, USA.

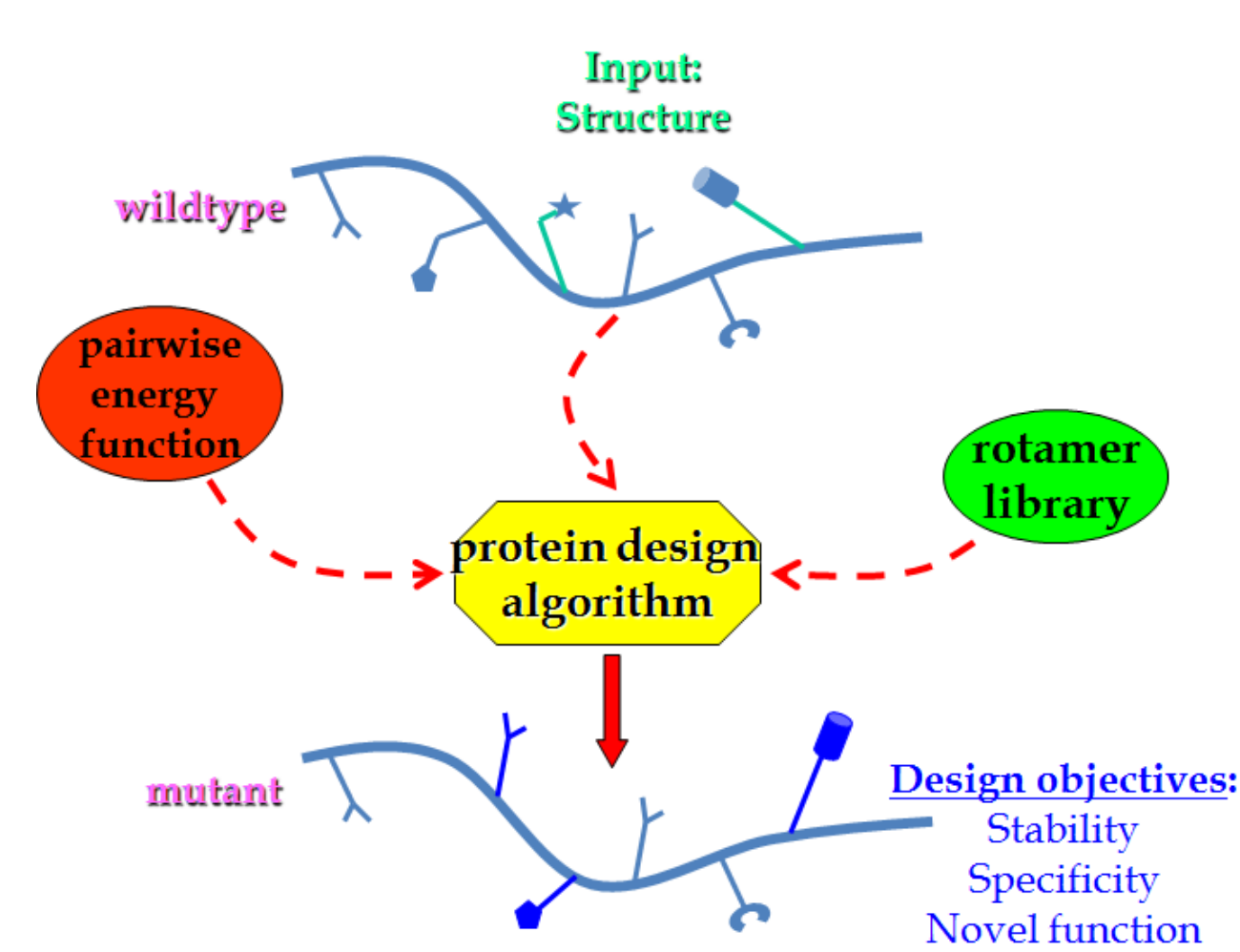


ABSTRACT

- We work on **multistate protein design problems** that address both positive and negative states and consider an ensemble of biophysical **substates** such as a protein being in various backbone conformers, unbound or bound to a target, or bound to various (off-)targets.
- The **generic formulation** allows for many applications such as stability, affinity, and specificity design.
- iCFN is an **exact algorithm** that guarantees the optimal solutions and near-optimal ensembles thus enable informative interaction with experiments. Its efficiency makes large-scale designs more tractable.
- Its **application** to T-cell receptor (TCR) design for specificity generates experimentally-agreeing results and reveals underlying mechanisms.

Availability: <https://shen-lab.github.io/software/iCFN>

PROBLEM DEFINITION



Protein design (Figure Credit: Ivelin Georgiev)

Energy model (Assumption: pairwise additive)

$$f(\mathbf{r}) = c + \sum_i E(i_r) + \sum_{i < j} E(i_r, j_s)$$

A generic formulation for multi-state protein design:

$$s^* = \arg \min_{s \in \mathcal{S}} \left(\min_{p \in \mathcal{P}} \min_{\mathbf{r} \in \mathcal{R}^+(s)} f_p^+(\mathbf{r}) - \min_{q \in \mathcal{Q}} \min_{\mathbf{r} \in \mathcal{R}^-(s)} f_q^-(\mathbf{r}) \right)$$

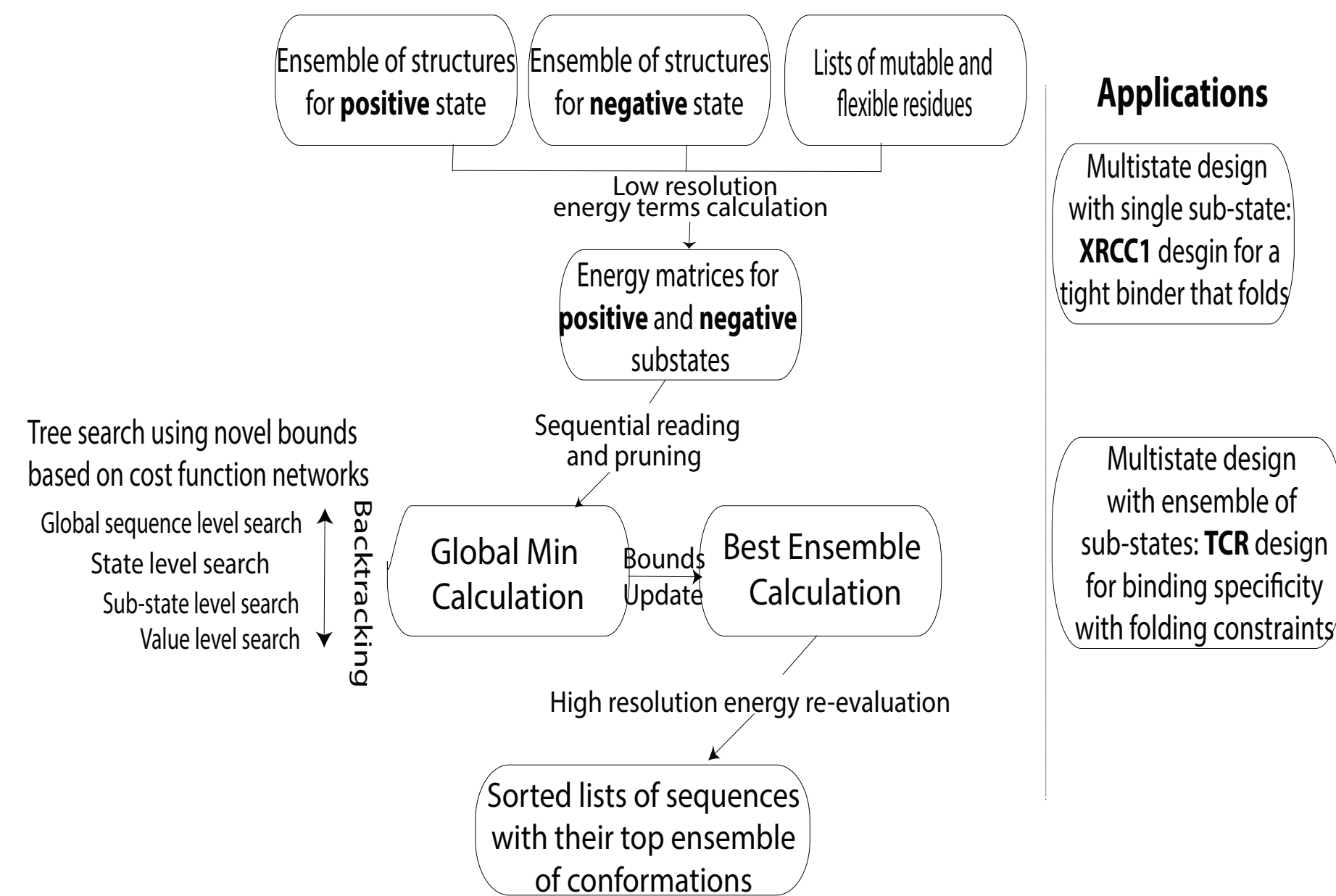
s.t. Constraints on substate functions $f_p^+(\mathbf{r})$ & $f_q^-(\mathbf{r})$

iCFN's approach to solving the NP-hard problem:

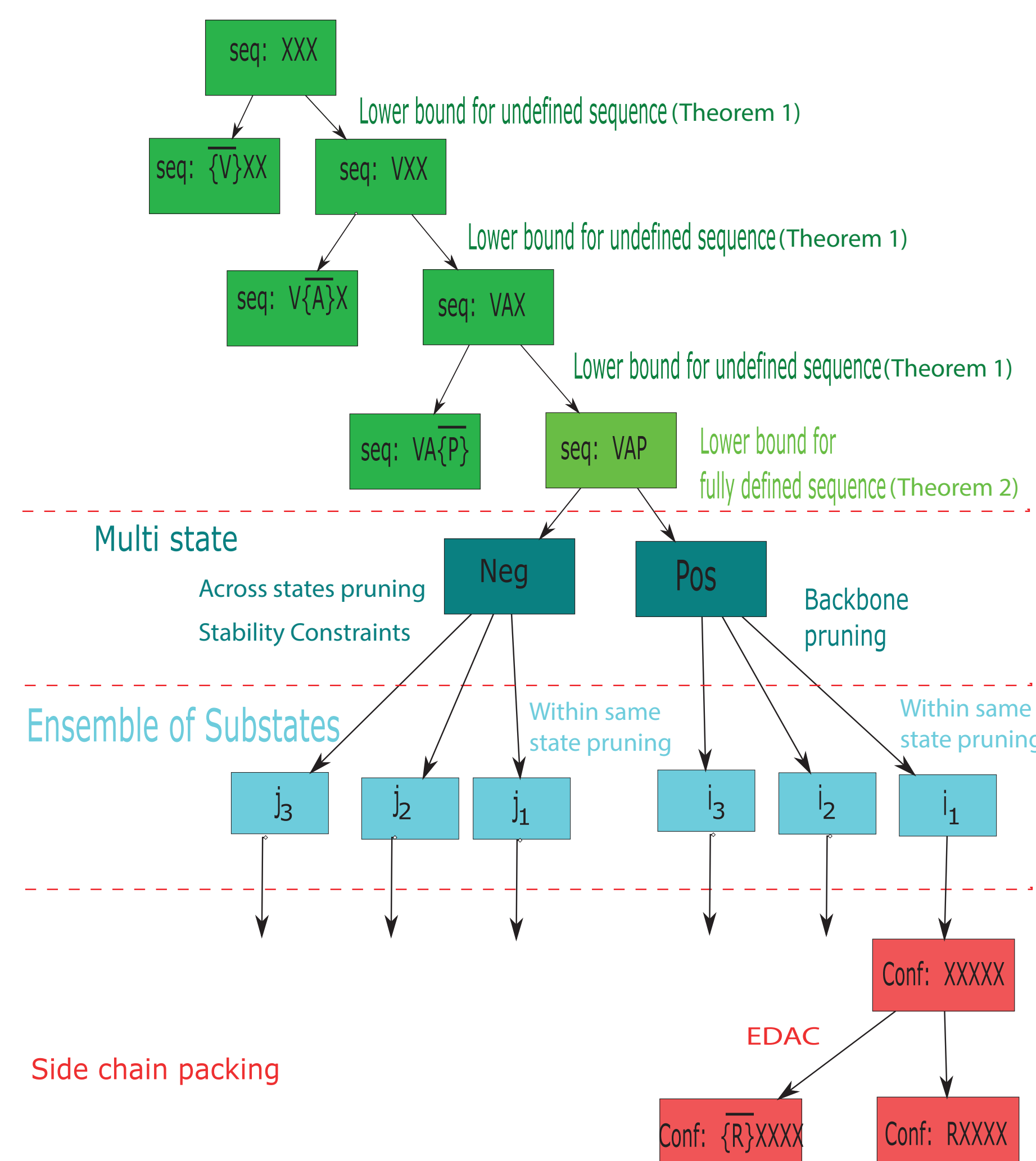
- Each substate design is formulated as a **Weighted Constraint Satisfaction Problem (WCSP)** and modeled by a **Cost Function Network (CFN)** ($\mathcal{X}, \mathcal{D}, \mathcal{C}$);
- The coupled WCSPs are represented as CFNs interconnected over a **tree** of sequences, substates, and rotamers (values);
- Novel lower **bounds** are developed for the sequence (variable) space (Thm. 1,2) and Existential Directional Arc Consistency (EDAC) is exploited for the rotamer (value) space;
- Depth First Branch and Bound (DFBB)**-based tree search allows positive and negative designs to inform each other and substates within and across states to prune each other.

METHODS

Overall scheme



Tree search



Theorem 1. Lower bound for any undefined sequence S :

$$\min_{(k,l)} \left(\Delta c_{kl} + \sum_i \min_{a \in S(i)} \min_{(r,r')} (\Delta E_{kl}(i_r, r') + \sum_{j > i} \min_{a' \in S(j)} \min_{(s,s')} \Delta E_{l,k}(i_r, r', j_s, s') \right)$$

with complexity $O(n^2 R^2 a^2 r)$ (where n is the number of positions, R the average number of rotamers per position, a the average number of substates per state, and r the average number of rotamers per amino acid).

Theorem 2. Lower bound for any defined sequence S :

$$\min_{k \in \mathcal{P}} L_k(S) - \min_{l \in \mathcal{Q}} U_l(S)$$

in which $L_k(S)$ is EDAC for sequence S in the k^{th} substate and U_l is LDS for sequence S in the l^{th} substate.

Additional bounds for each substates, across substates of the same state, and across substates of different states.

RESULTS

Multi-state XRCC1 design with a single substate per state:

N _{seq}	d(A)	d(N)	e = 0.5 kcal/mol		1 kcal/mol		1.5 kcal/mol				
			Pre-DEE Size	Post-DEE Size (Ensemble)	COMETS	Reduced iCFN	COMETS	Reduced iCFN	COMETS	Reduced iCFN	
1	3	9	2.98 × 10 ¹¹	2.94 × 10 ¹⁰	4.03	0.82	0.82	2.27	0.13	2.27	0.05
1	6	16	5.24 × 10 ¹⁰	2.75 × 10 ¹⁰	6.84	0.16	0.12	6.97	0.17	0.12	0.19
2	3	10	1.1 × 10 ¹²	2.73 × 10 ¹¹	6.65	0.28	0.15	6.36	0.16	0.29	0.16
2	6	19	5.88 × 10 ¹¹	1.58 × 10 ¹¹	19.46	2.63	1.41	29.07	2.67	1.44	29.27
3	3	11	7.94 × 10 ¹⁰	4.31 × 10 ¹⁰	M	12.64	6.16	M	12.67	6.51	M
3	6	20	1.81 × 10 ¹¹	1.44 × 10 ¹¹	M	62.17	32.9	M	62.22	33.26	M
4	3	14	6.54 × 10 ¹⁰	4.31 × 10 ¹⁰	M	401.26	264.26	M	403.51	268.3	M
4	6	25	7.94 × 10 ¹⁰	4.16 × 10 ¹⁰	M	2000	1458	M	2156	1493	M
5	3	15	5.24 × 10 ¹⁰	2.74 × 10 ¹⁰	M	4803	3738	M	4913	3605	M
5	6	26	1.65 × 10 ¹¹	1.38 × 10 ¹¹	M	48373	37398	M	48974	37223	M

("M" indicates an error for being out of a 20Gb-memory limit whereas iCFN used at most 80Mb)

iCFN outperforms COMETS, the only other exact method for multistate design, in both memory usage and CPU time, which enables large designs in practice.

Multi-state TCR design with ensembles of substates (MD simulated)

(target peptide: AAG; off-target peptide: ELA)
Guaranteed near-optimum ensemble:

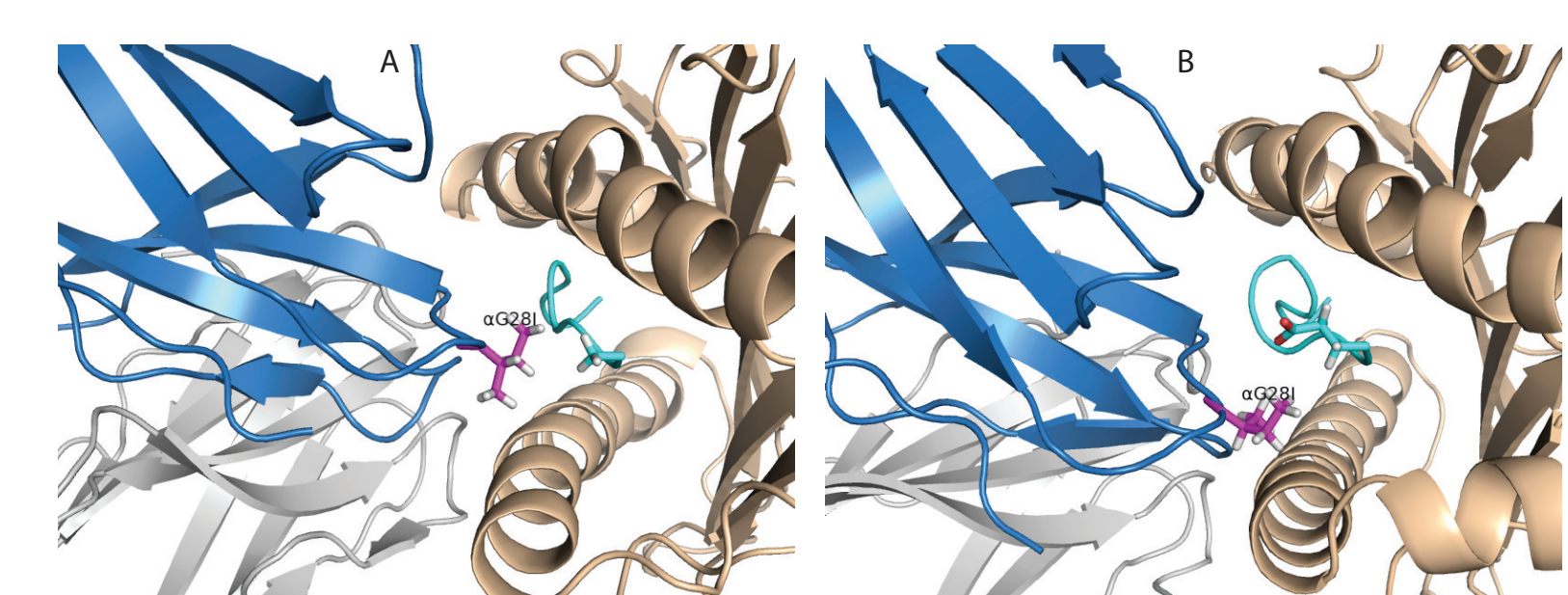
Position	Reduced iCFN				iCFN					
	Pre-DEE Size	Post-DEE Size	Nodes Expanded	Leaves Visited	Sequences	Time (s)	Nodes Expanded	Leaves Visited	Sequences	Time (s)
28	10 ¹⁰	10 ⁵	6.83 × 10 ⁵	6.20 × 10 ⁵	26	16.69	6.14 × 10 ⁵	6.00 × 10 ⁵	10	21.36
28	10 ¹⁰	10 ⁵	5.92 × 10 ⁵	5.70 × 10 ⁵	26	114.22	4.09 × 10 ⁵	4.00 × 10 ⁵	2	21.35
98	10 ¹⁰	10 ⁵	5.15 × 10 ⁵	5.00 × 10 ⁵	25	103.29	4.16 × 10 ⁵	4.00 × 10 ⁵	20	43.35
100	10 ¹⁰	10 ⁵	7.43 × 10 ⁵	7.11 × 10 ⁵	26	154.51	5.19 × 10 ⁵	5.00 × 10 ⁵	2	23.54
26,28	10 ¹⁰	10 ⁵	1.70 × 10 ⁶	1.62 × 10 ⁶	676	7454.95	9.44 × 10 ⁵	9.00 × 10 ⁵	4	1083.89
28,98	10 ¹⁰	10 ⁵	1.82 × 10 ⁶	1.73 × 10 ⁶	650	17440.96	2.51 × 10 ⁶	2.38 × 10 ⁶	108	2072.52
28,100	10 ¹²	10 ¹²	1.49 × 10 ⁶	1.75 × 10 ⁶	676	19780.66	2.62 × 10 ⁶	2.40 × 10 ⁶	10	2226.52
28,98	10 ¹⁰	10 ⁵	1.45 × 10 ⁶	1.37 × 10 ⁶	650	22378.53	3.13 × 10 ⁶	2.90 × 10 ⁶	15	2510.31
28,100	10 ¹²	10 ¹²	1.77 × 10 ⁶	1.60 × 10 ⁶	676	24631.34	4.22 × 10 ⁶	3.80 × 10 ⁶	19	2656.47
28,28,100	10 ¹⁰	10 ⁵	—	—	17576	—	1.48 × 10 ⁶	1.40 × 10 ⁶	6	950.12
28,98,100	10 ¹⁰	10 ⁵	—	—	16900	—	6.76 × 10 ⁵	6.00 × 10 ⁵	27	183.86
28,98,100	10 ¹⁰	10 ⁵	—	—	16900	—	3.73 × 10 ⁵	3.40 × 10 ⁵	12	158.98

- iCFN visits on average 6.7 (7.4), 58.8 (110.8), and 455.1 (1397.2) times less sequences for the best single (ensemble of) sequence(s) in single, double, and tripe designs, respectively.
- iCFN runs 3.4 (4.2) and 5.9 (7.8) times faster than reduced iCFN does for global optimum (top ensemble) in an average single and double design, respectively; and it solves tripe designs within 1~2 CPU days whereas reduced iCFN could not within 1 CPU week.
- iCFN's relative computational gain increases as complexity increases!

Design accuracy:

Method	True Positive (TP)	FP	False Negative (FN)
Rosetta	G28I, G28L, G28Y, F100W	D26Y	D26W, F100Y
Rosetta Min	G28I, G28L, G28Y	N/A	D26W, F100W, F100Y
iCFN	D26W, G28I, G28L, G28Y, F100W, F100Y	D26Y	N/A

Molecular mechanisms of AAG-binding specificity for G28I:



Differential effects of G28I to (A) AAG-binding and (B) ELA-binding revealed in iCFN structural models. Cartoons: DMF5 α ; DMF5 β ; AAG/ELA peptides; MHC α chain. Stick: α G28I. Worse vdW packing and continuum electrostatics upon mutation for N-terminal glutamate of ELA but not for N-term alanine of AAG.

REFERENCES

- [1] Mostafa Karimi, Yang Shen. "iCFN: an efficient exact algorithm for multistate protein design", Bioinformatics 34 (17), i811-i820.

ACKNOWLEDGEMENTS

This project is in part supported by the National Institute of General Medical Sciences of the National Institutes of Health (R35GM124952 to YS) and the Defense Advanced Research Projects Agency (FA8750-18-2-0027 to ZW). Most of the computing time is provided by the Texas A&M High Performance Research Computing.