



Overview

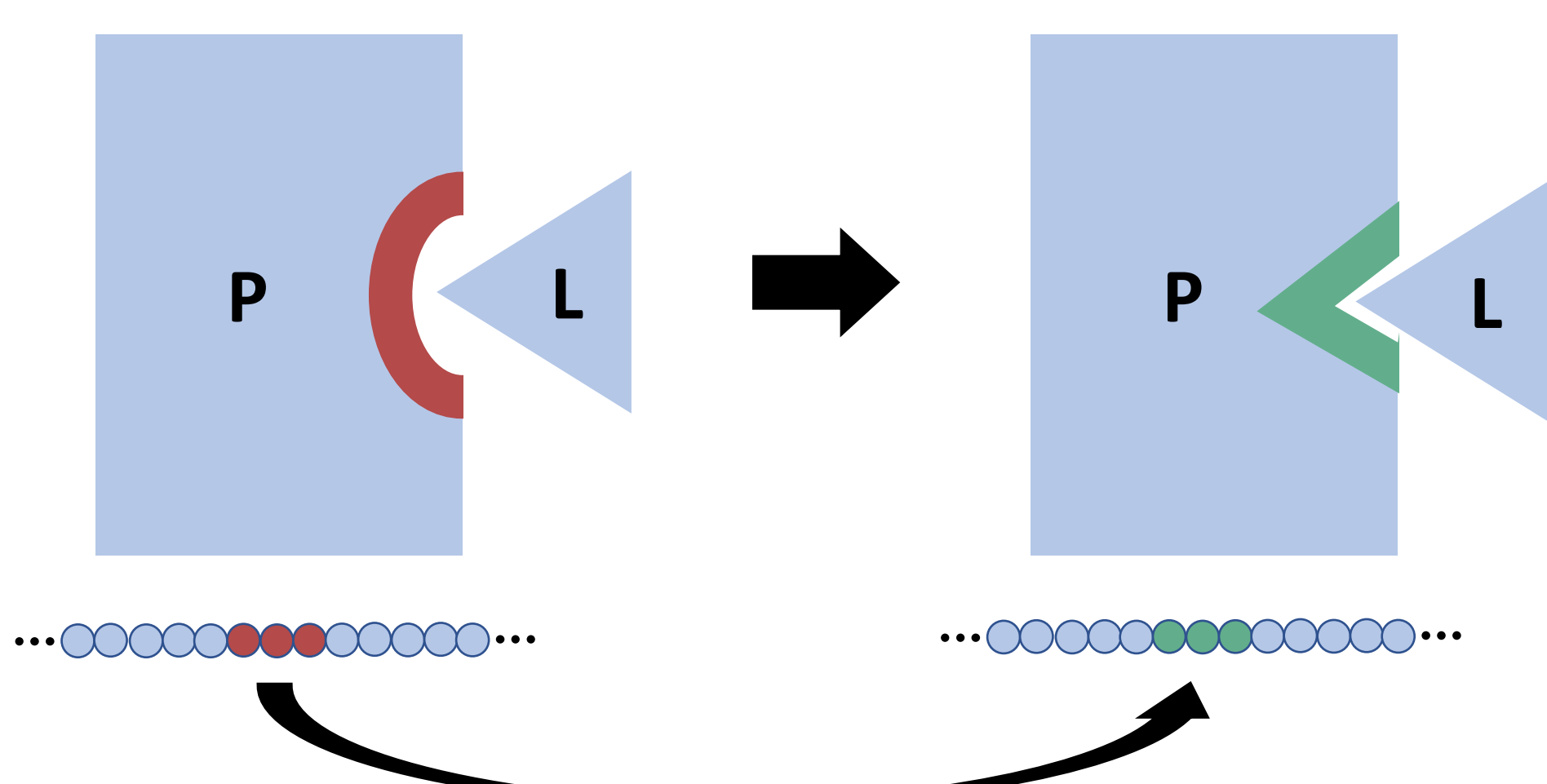
The objective of this work is to (1) create a framework for future advancements for protein design via K* optimization (ie. K*MAP) leveraging powerful MMAP algorithms, (2) explore new related bounded heuristics, and (3) create a foundation for efficient algorithms solving such protein design problems.

Contributions:

1. Formulation of K*MAP as a graphical model
2. wMBE-K*, a weighted mini-bucket scheme for K*MAP enhanced with a domain partitioning scheme
3. AOBB-K*, a depth-first branch-and-bound algorithm over AND/OR search spaces for solving K*MAP
4. A thresholding scheme introducing and exploiting determinism with correctness guarantees
5. Extensive analysis comparing these schemes to state-of-the-art BBK* illustrating their potential

Problem

Redesign of proteins to form higher affinity complexes



Find new amino acid assignment to residues of interest that optimize affinity between interacting subunits

K* Objective

An approximation of binding affinity between molecules (based on the biological association constant known as K_a)

$$K^*(r) = \frac{Z_{PL}(r)}{Z_P(r) Z_L(r)}$$

$$Z_\gamma(r) = \sum_{c \in C_\gamma(r)} \exp\{-E_\gamma(c)/RT\}$$

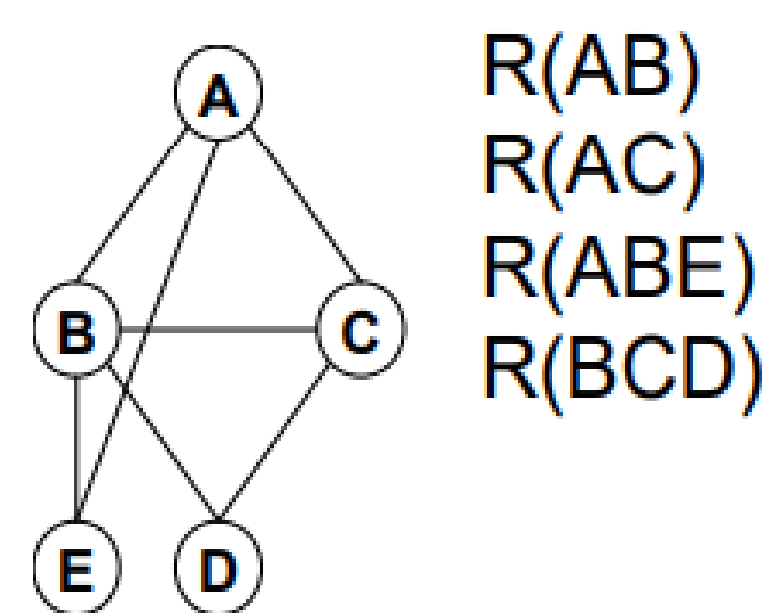
r = amino acid assignments to residues
 $C(r)$ = possible conformations given r
 $E(c)$ = energy given conformation c
 \mathcal{R} = universal gas constant
 \mathcal{T} = absolute temperature (Kelvin)

captures the goodness of the subunit(s) in form $\gamma \in \{P, L, PL\}$

AND/OR Search

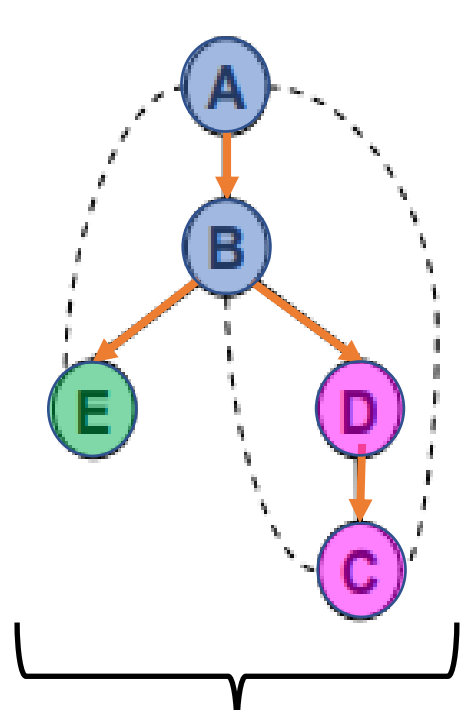
Compact search space taking advantage of conditional independences present in the model

Graphical Model Network



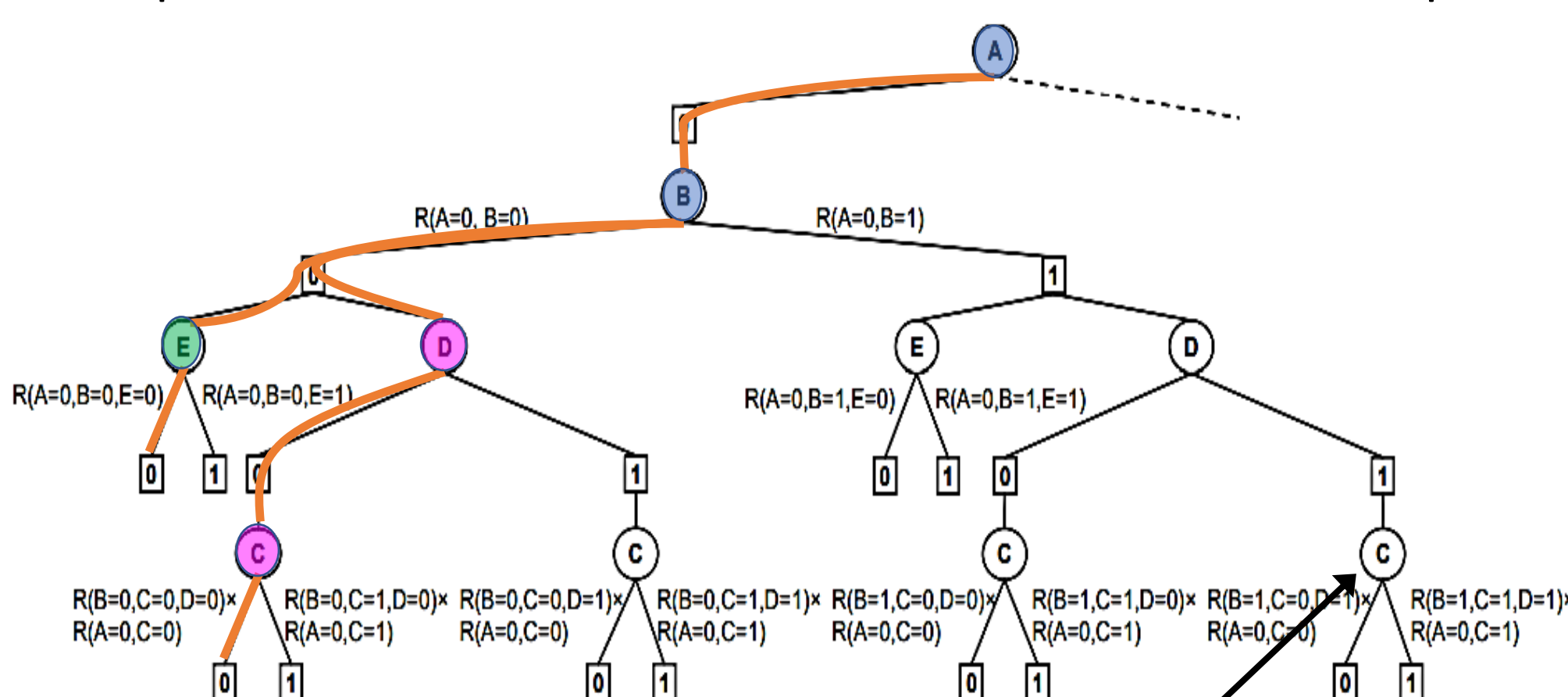
$R(AB)$
 $R(AC)$
 $R(ABE)$
 $R(BCD)$

Possible Pseudo Tree



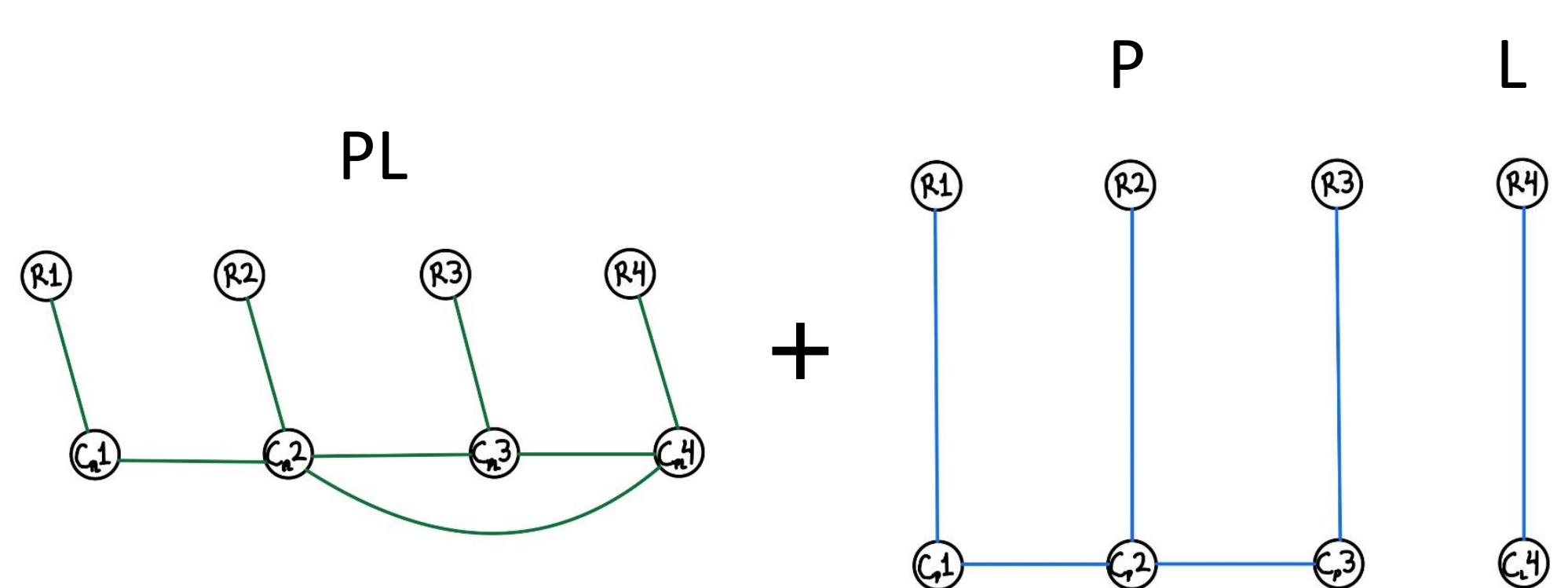
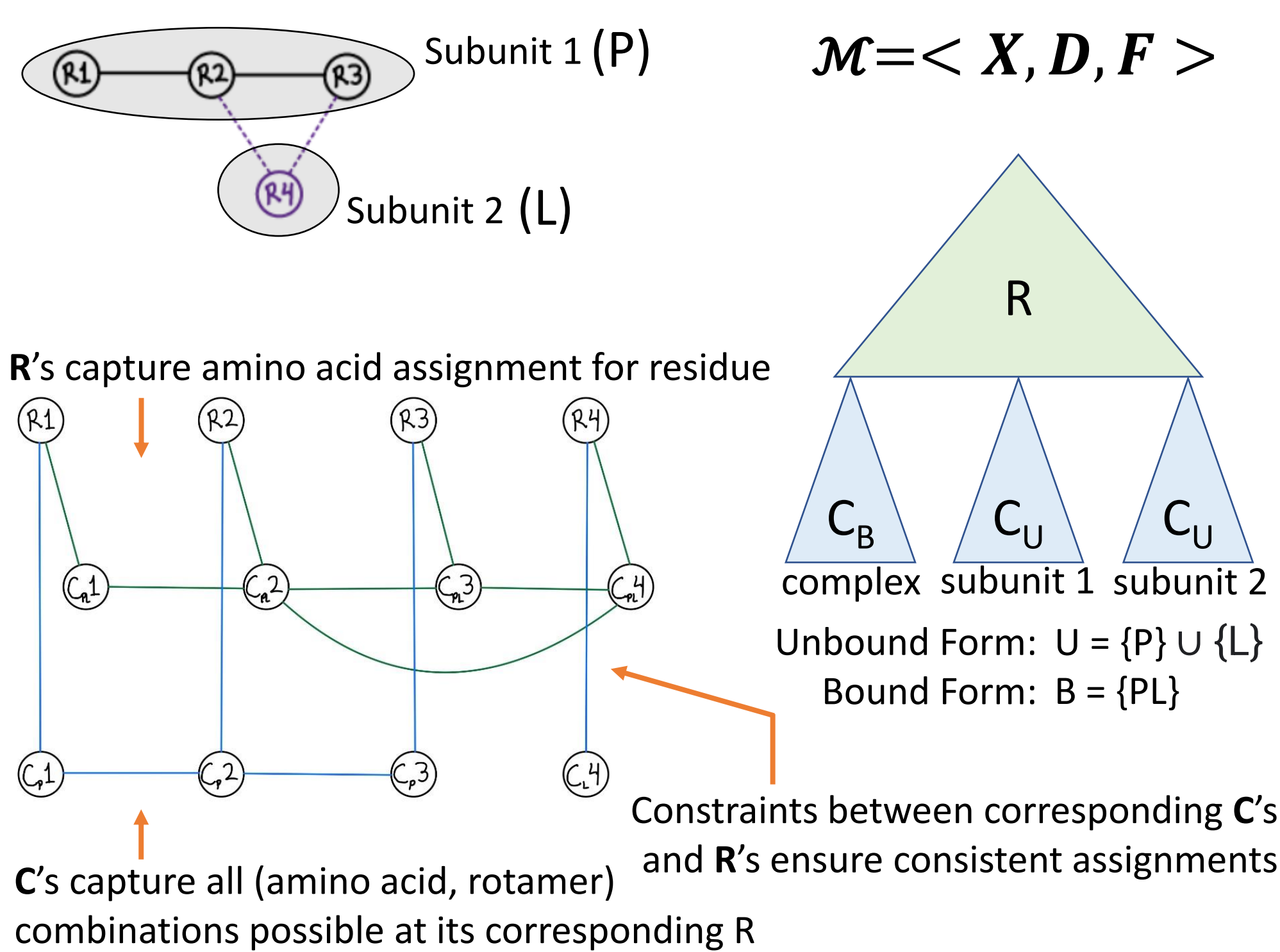
Directed tree (based on a variable ordering) that branches when conditional independences exist given assignments to ancestors.

The pseudo tree is used to construct the AND/OR search space



OR nodes represent variables
 AND nodes represent assignments to their corresponding parent variable

Graphical Model Formulation



K*MAP Task

Find amino acid assignments to the residues that maximize K*

$$Z_\gamma(R_1 \dots R_N) = \sum_{C_1, \dots, C_N} \prod_{C_\gamma(i) \in \mathcal{C}} \mathcal{C}_{\gamma(i)}(R_i, C_{\gamma(i)}) \cdot \prod_{E_\gamma^{sb}(i) \in E_\gamma^{sb}} e^{-\frac{E_\gamma^{sb}(i)(C_{\gamma(i)})}{RT}} \cdot \prod_{E_\gamma^{pw}(i, j) \in E_\gamma^{pw}} e^{-\frac{E_\gamma^{pw}(i, j)(C_{\gamma(i)}, C_{\gamma(j)})}{RT}}$$

$$K^*(R_1, \dots, R_N) = Z_B(R_1, \dots, R_N) / Z_U(R_1, \dots, R_N)$$

$$\text{task: } K^*MAP = \max_{R_1, \dots, R_N} K^*(R_1, \dots, R_N)$$

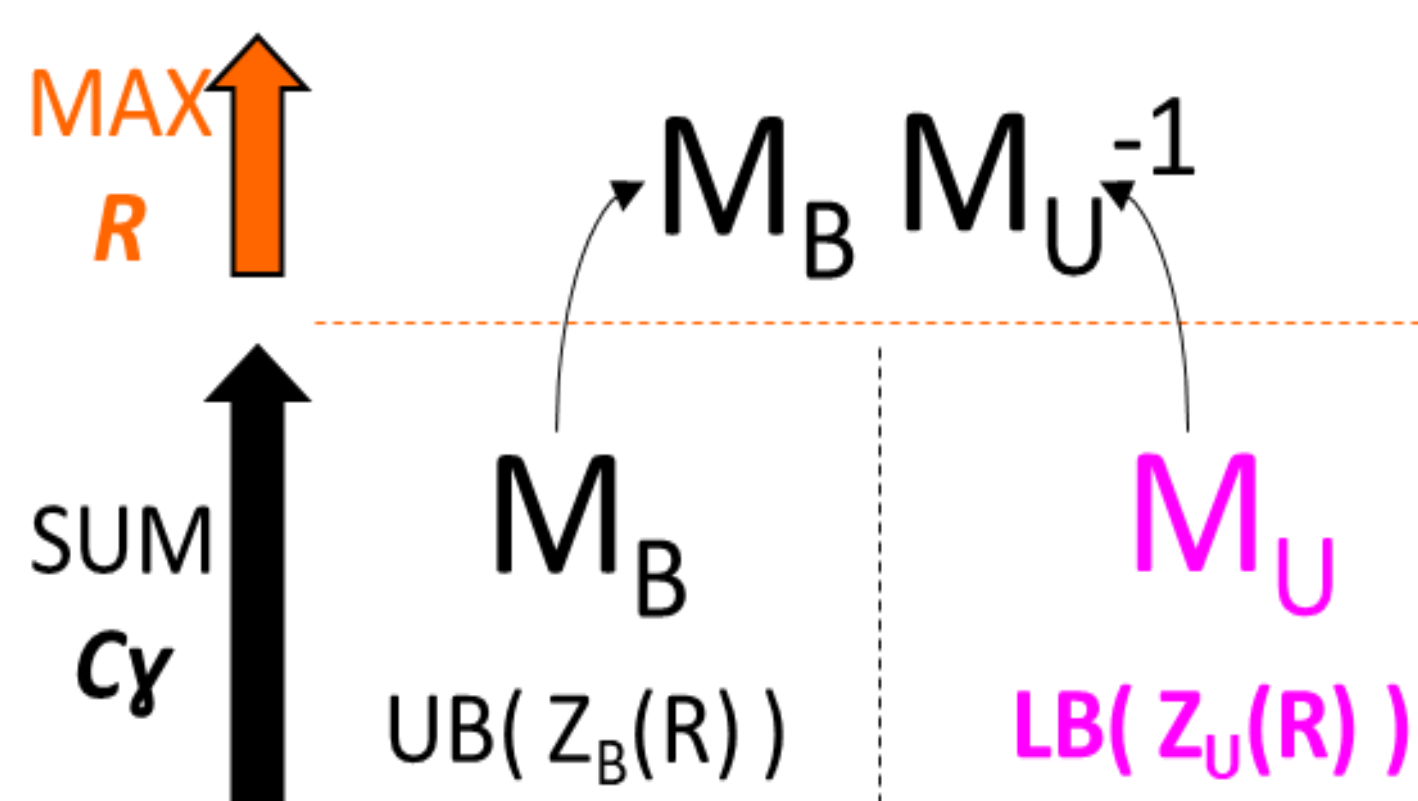
wMBE-K*

Precompiled K* heuristic to guide search, based on dynamic programming message-passing Mini Bucket Elimination

$$K^*(r) = \frac{Z_{\text{complex}}(r)}{Z_{\text{subunit 1}}(r) Z_{\text{subunit 2}}(r)}$$

$$\sum_r [\prod(\psi_r)] \leq \prod(\sum_x)$$

$$\sum_x f(x) \triangleq [\sum_x f(x)^{\frac{1}{w}}]^w \quad w = \sum_r w_r$$



Domain-Partitioned wMBE-K*

Strategy to improve lower bounds generated by wMBE-K*

- Let X, Y , and Z be three variables
- Let $obj = \sum_X f(x, y) \cdot g(x, z)$
- Let $X' = \{x \in X | g(x, z) \neq 0\}$ s.t. $\epsilon_{X'} = \min_{x \in X'} g(x, z)$
- Since $\epsilon_{X'} > 0$, we can derive...

$$obj = \sum_X f(x, y) \cdot g(x, z)$$

$$obj = \sum_{x \in X'} f(x, y) \cdot g(x, z) + \sum_{x \in X \setminus X'} f(x, y) \cdot g(x, z)$$

$$obj = \sum_{x \in X'} f(x, y) \cdot g(x, z) \geq \epsilon_{X'} \cdot \sum_{x \in X'} f(x, y) > 0$$

AOBB-K*

- Branch-and-bound algorithm over AND/OR search spaces
 - AOBB-K* is exact
- Can use wMBE-K* to guide search
- Exploits determinism by using constraint propagation
- Incorporates a global constraint enforcing biologically relevant solutions

Subunit Stability Constraints

Condition to enforce the stability of each subunit to be no less than a given threshold from that of the wild-type stability

$$Z_{\text{subunit } i}(r) > Z_{\text{subunit } i}(r^{wt}) * \exp\{-5/RT\}$$

Stability of naturally occurring version Constant factor for thresholding

Infusing Determinism: τ -Underflows

Replace extremely unfavorable assignments with hard constraints taking exploiting the strength of constraint propagation

- Let f be a non-negative function
- Consider $\tau \in \mathbb{R}^+$
- Then the τ underflow of f is...

$$f(x) = \begin{cases} f(x), & f(x) \geq \tau \\ 0, & \text{otherwise} \end{cases}$$

Empirical Analysis

Competing Scheme: BBK*

[Ojewole et al., 2018]

State-of-the-art protein redesign algorithm, part of the software package, OSPREY, developed for over ten years for protein design

- A*-like best first
- utilizes dynamic optimistic greedy heuristic
- approximate scheme w/ tightness parameter

Empirical Results

Dataset (#instances)	AOBB-K*		any-AOBB-K*	
	K* ≥ (F1,F2)	K* > (F1,F2)	time < (F1,F2)	K* ≥ (F1,F2)
Orig. (30)	30,30	2,2	23,28	30,30
Expand. (12)	6,11	0,4	2,4	11,11

benchmarks (3 MAP Variables)	ω	τ	(203 ≤ Dmax ≤ 206)			[ω, τ]-AOBB-K*		BBK*				
			iB w*	d X	UB	pre-t	search	time	K*	time	K*	
1gwc_00021*	-	-	4	4	7	13	28.8	123.8	81.3	205.1	551.3	11.7
	0.001	-	4	4	7	13	28.8	124.3	12.1	136.4		
	-	1E-05	4	4	7	13	28.8	117.1	3.5	120.7		
2hmv_00025*	-	-	4	6	9	17	42.3	109.8	44.0	153.8	880.5	13.6
	0.001	-	4	6	9	17	42.3	109.3	12.2	121.5		
	-	1E-05	4	6	9	17	42.3	100.8	1.7	102.4		
2rf9_00013*	-	-	4	6	9	17	37.7	83.0	17.8	100.8	39.2	15.0
	0.001	-	4	6	9	17	37.7	82.8	1.6	84.4		
	-	1E-05	4	6	9	17	37.7	71.4	0.4	71.8		
2rfe_00012*	-	-	4	5	8	15	34.1	58.3	2.6	60.9	11.8	13.9
	0.001	-	4	5	8	15	34.1	58.6	0.3	58.9		
	-	1E-05	3	5	8	15	5.5	14.9	20.3			
2rfe_00014*	-	-	4	5	8	15	35.1	58.2	2.4	60.6	44.9	14.4
	0.001	-	4	5	8	15	35.1	58.2	1.3	59.4		
	-	1E-05	3	5	8	15	4.9	15.2	20.0			
2rfe_00017*	-	-	5	5	8	15	26.4	166.8	167.8	334.6	78.0	10.8
	0.001	-	4	5	8	15	27.4	89.1	5.0	94.1		
	-	1E-05	4	5	8	15	27.4	85.7	7.1	92.8		
2rfe_00030*	-	-	4	5	8	15	31.3	101.4	0.5	101.9	275.4	11.0
	0.001	-	4	5	8	15	30.4	105.4	3.5	108.9		
	-	1E-05	4	5	8	15	30.4	105.4	3.5	108.9		
2xgy_00020*	-	-	5	5	8	15	26.2	81.8	278.9	360.7	1388.1	11.0
	0.001	-	4	5	8	15	27.2	60.4	8.2	68.6		
	-	1E-05	4	5	8	15	27.2	60.1	23.8	83.9		
3u7y_00009*	-	-	4	4	7	13	11.4	62.6	36.8	99.5	215.8	4.5
	0.001	-	4	4	7	13	11.4	62.5	2.1	64.7		
	-	1E-05	4	4	7	13	61.3	5.1	66.4			
3u7y_00011*	-	-	4	4	7	13	28.3	74.0	2.1	76.1	26.6	11.9
	0.001	-	4	4	7	13	28.3	83.4	0.1	83.5		
	-	1E-05	4	4	7	13	28.3	80.7	0.5	81.2		
4wwi_00019*	-	-	5	5	8	15	37.0	169.2	12.1	181.3	34.0	15.0
	0.001	-	4	5	8	15	38.0	62.0	7.2	69.2		
	-	1E-05	4	5	8	15	54.2	8.4	62.7			

Summary

- AOBB-K* shows promise vs. state-of-the-art BBK*
- Competitive run-times
- Can find better solutions
- τ -AOBB-K* can greatly improve runtimes
- AOBB-K* has scalability issues as the number of MAP variables increase

Acknowledgements

Special Thanks:

- Thomas Schiex [INRAE, France]
- Bruce Donald and members of his lab [Duke University] including: Graham Holt, Jonathan Jou and Nathan Guerin

Supported in part by: NSF Grant IIS-2008516