

Engineering approach to protein design

Most of biochemistry research has been focused on *reverse* engineering proteins, pathways, and cells. We would like to turn this around and start *forward* engineering new functions into proteins.

Airplanes, buildings, and circuits are all designed in a computer using accurate models. Can we do the same thing for protein design?



Computer model of an airplane flying at Mach 0.9. Warmer colors denote higher air pressures; the transparent region indicates where local airflow velocity is Mach 1.

Steven Ashley. (Nov. 2003) "Flying on flexible wings" *Scientific American.* 84-91. Image credit: Richard Snyder, Wright-Patterson Air Force Base AFRL.

Application: Biosensor

Subcellular fluorescence imaging of signaling molecules

IP₃, cAMP, leukotrienes

Drugs with a low therapeutic index



Detecting bacteria





o-aminoacetophenone (P. aeruginosa)

dipicolinic acid (Bacillis)

Cancer diagnosis vanillylmandelic acid (presence in urine indicates pheochromocytoma or neuroblastoma)



Application: Custom enzymes

Binding to the transition state of a reaction catalyzes that reaction:



- Custom proteases and restriction enzymes
- Chemical synthesis
- Degrade toxins, bacterial biofilm matrix



Outline

Protein design technology

Initial tests

- Structural prediction
- Energetic prediction
- Redesigning a natural binding site

Application: Designing new binding sites

Protein design technology

Setting up a design calculation



• Pick scaffold with known crystal structure; pick design positions.

Protein design technology

Setting up a design calculation



• Construct possible ligand poses and sidechain conformations for each amino acid at each position.

Protein design technology

Setting up a design calculation

Ligand conformation \rightarrow



6 Calculate matrix of interaction energies between possible conformations.



properties) and structural optimization (to determine the properties of each proposed sequence). In a typical protein design calculation, we search a space of 5000 alternative conformations at 10 positions $\approx 10^{37}$ possibilities.





Continuum solvent model

Generalized Born radius



1.3 Born radius (Å) 14.7

Charged atoms closer to the protein's surface have:

- more favorable solvation energy
- smaller charge-charge interactions

Hydrophobic effect is proportional to exposed surface area.



Examples of protein behaviors treated by our model

We explicitly model the bound, unbound, and unfolded states. This allows us to model conformational changes, and also to optimize for stability, binding, and specificity separately.



Outline

Protein design technology

Initial tests

- Structural prediction
- Energetic prediction
- Redesigning a natural binding site

Application: Designing new binding sites

Predicting binding site coordinates

	Ligand poses	Sidechain rotamers	crystal structure / predicted structure	RMS error
ABP- arabinose	4111 ligand poses			
RBP- ribose	4639 ligand poses			

Predicting binding site coordinates

	Ligand poses	Sidechain rotamers	crystal structure / predicted structure	RMS error
ABP- arabinose	4111 ligand poses	6028 rotamers / position		
RBP- ribose	4639 ligand poses	5449 rotamers / position		

Predicting binding site coordinates

	Ligand poses	Sidechain rotamers	crystal structure / predicted structure	RMS error
ABP- arabinose	4111 ligand poses	6028 rotamers / position		0.677 Å
RBP- ribose	4639 ligand poses	5449 rotamers / position		0.148 Å





Highest resolution rotamer library has a rotamer within 0.3 Å of 99.9% of rotamers seen in high resolution crystal structures.

Energetic predictions: Scrambled sequences



Energetic predictions: Relative binding energies of mutants



Mutants of ABP binding arabinose

Why are structures easier to predict than energies?



Redesigned ribose binding protein



Positions identical to the native are highlighted in yellow.

Redesigned ribose binding protein

Out of $17^{10} = 2.0 \times 10^{12}$ possible sequences, the design algorithm picked a point mutant of the native as the top sequence, and the native as the second best sequence.

To give an idea of the size of the sequence space considered, this is equivalent to locating a 4.6 m^2 region (red square) in the entire United States.



Protein design vs. in vitro evolution



Redesigned ribose binding protein

Essential elements of the design algorithm:

- High resolution rotamer library
- Final gradient-based local minimization step
- Accurate solvation model

Design	Rotamers /	Local	Solvation	Rank	Identity	K _d	Sequence (10
calc.	position	minimization	model		to native	(exptl.)	primary contacts)
		no		1	3	> 330 mM	<mark>N</mark> MMMIM <mark>FN</mark> AN
А	2800		Brooks	∠	2		<mark>n</mark> mmmlm <mark>f</mark> tan
				3	4		<mark>nf</mark> mlvm <mark>fn</mark> an
B 5449				[1	8	690 mM	<mark>NFFDRRF</mark> SS <mark>Q</mark>
	5449	no	Brooks	≺ 2	9		<mark>NFFDRRFN</mark> S <mark>Q</mark>
				3	8		<mark>N</mark> MFDRRFN <mark>S</mark> Q
		yes	Still	[1	6	> 900 mM	<mark>N</mark> YY <mark>DRR</mark> Y <mark>N</mark> AQ
С	5449			┤ 2	6		<mark>nymdrryn</mark> s <mark>q</mark>
				3	7		NY <mark>FDRR</mark> Y <mark>N</mark> AQ
				[1	9	17.2 µM	L <mark>FFDRRFNDQ</mark>
★ D	5449	yes	Brooks	┤ 2	10	210 nM	<mark>NFFDRRFNDQ</mark>
				5	9		NT <mark>FDRRFNDQ</mark>
Native					10	210 nM	NFFDRRFNDQ

Note: The Brooks solvation model is only 2% off from the Poisson-Boltzmann equation. The Still solvation model is less accurate.

High resolution is critical for protein design



Outline

Protein design technology

Initial tests

- Structural prediction
- Energetic prediction
- Redesigning a natural binding site

Application: Designing new binding sites

Design targets

Redesign ribose binding protein to bind ...



D-xylose он----он в он он







HC

Designed binding proteins

	9	13	15	16	89	90	103	105	141	164	190	215	235
RBP (native)	SER	ASN	PHE	PHE	ASP	ARG	SER	ASN	ARG	PHE	ASN	ASP	GLN
RBP→arabinose 2	GLN	ASN	MET	TYR	VAL	MET	GLN		MET	PHE	ASN	SER	VAL
RBP→xylose 1		ASN	PHE	PHE	GLN	GLN			MET	PHE	ASN	SER	MET
RBP→xylose 2		MET	TYR	PHE	GLN	HIS			MET	PHE	ASN	SER	GLN
RBP→estradiol 4	SER	ASN	VAL	MET	ALA	ASN	ASN		MET	PHE	ASN	SER	ILE
RBP→IAA 1		ARG	THR	MET	VAL	MET	HIS	TYR	MET	PHE	ASN	ALA	SER
RBP→IAA 2		ARG	THR	MET	ALA	MET	HIS	TYR	MET	PHE	ASN	SER	SER
RBP→IAA 3		ARG	THR	MET	VAL	ASN	HIS	TYR	MET	PHE	ASN	ALA	SER
Sequence is only show	vn at p	ositior	ns beir	na des	ianed.								



Experimental characterization of designed binding proteins

Protein	Ligand	K _d	Protein stability (kcal/mol)
RBP→arabinose 2	arabinose	250 mM	
RBP→xylose 1	xylose	160 mM	4.2
RBP→xylose 2	xylose	270 mM	
RBP→estradiol 4	estradiol	46 mM	2.0
RBP→IAA 1	IAA	11 mM	2.5
RBP→IAA 2	IAA	14 mM	1.0
RBP→IAA 3	IAA	16 mM	1.9

Designed estradiol binding protein

human estrogen receptor (1A52)



Shape complementarity = 0.72



RBP->estradiol 4

Shape complementarity = 0.75



Designed indole-3-acetic acid binding protein

RBP→IAA 1



Shape complementarity = 0.81



Designed arabinose binding protein



RBP→IAA library

We generated a library of RBP variants based on amino acid frequencies in the top 72 sequences from the RBP \rightarrow IAA design.



RBP→IAA library

We screened 279 sequences from the library for binding to IAA, and found sequences with 10-fold improved affinity 10-fold (highlighted):

Protein	K _d (IAA)	13	15	16	89	90	103	105	141	164	190	215	235
RBP→IAA 1	11 mM	ARG	THR	MET	VAL	MET	HIS	TYR	MET	PHE	ASN	ALA	SER
RBP→IAA 2	14 mM	ARG	THR	MET	ALA	MET	HIS	TYR	MET	PHE	ASN	SER	SER
RBP→IAA 3	16 mM	ARG	THR	MET	VAL	ASN	HIS	TYR	MET	PHE	ASN	ALA	SER
RBP→IAA 101A-F11	1.4 mM	ARG	SER	MET	GLY	CYS	HIS	TYR	MET	PHE	ASN	ALA	SER
RBP→IAA 95A-C1	1.1 mM	ARG	SER	MET	ILE	CYS	HIS	TYR	MET	PHE	ASN	ALA	SER





Designing for specificity



Space for galactose CH₂OH is seen in the native and arabinose binding design (both of which bind galactose), but not in the specificity design.







ligand



Acknowledgements

Harbury Lab

Erica Raffauf Jim Havranek Jarrett Wrenn Lance Martin Rebecca Weisinger Becky Fenn Kierstin Schmidt Dan Herschlag Axel Brunger Tom Wandless

Michael Levitt Buzz Baldwin Loren Looger Pavel Strop

End

Examples of computational protein design

Serotonin receptor created from an arabinose receptor



Looger LL, Dwyer MA, Smith JJ, Hellinga HW. (2003) Nature 423: 185-90.

Protein with a new α/β fold



Calculated / Observed Kuhlman B, Dantas G, Ireton GC, Varani G, Stoddard BL, Baker D. (2003) *Science*. 302:1364-8.

Specific coiled-coil interactions



Havranek JJ, Harbury PB. (2003) Nature structural biology. 10(1): 45-52.

Finding a protein's low energy conformations

Molecular dynamics is slow because crossing large energetic barriers is a rare event. Can we skip these barriers and just find the low energy conformations?



Low energy loop conformations

- Database search
- Systematic search
- Randomly perturb and splice together existing conformations



Low energy sidechain conformations

Rotamers seen in crystal structures



Rotamer library



Probabilistic description of protein conformation

We represent the protein/ligand system as a probabilistic ensemble of different backbone, side chain, and ligand conformations. This allows us to model conformational changes and thermal fluctuations.

Loop conformations





Brooks' empirical expression with a $1/r^5$ term gives much more accurate Born radii for atoms in a test protein (protein tyrosine phosphatase 1B).





Sidechain conformations

For a set of small molecules, peptides, and proteins, solvation energy calculated using the $1/r^5$ term closely matches the goldstandard Poisson-Boltzmann equation (each point corresponds to a single structure).

Design positions for RBP



For ribose binding protein, we only picked primary contacts.

Primary contacts with ligand:

h-bonding	13	Asn
	89	Asp
	90	Arg
	141	Arg
	190	Asn
	215	Asp
	235	Gln
hydrophobic	15	Phe
	16	Phe
	164	Phe

Primary contacts: 13,15,16,89,90,141,164,190,215,235 Secondary contacts: 9,64,103,132,137,214

Computational point mutagenesis

Predicted dissociation energy (kcal/mol, relative to native) of RBP-ribose

							Res	idue					
		9	13	15	16	89	90	103	141	164	190	215	235
Native seq	uence	SER	ASN	PHE	PHE	ASP	ARG	SER	ARG	PHE	ASN	ASP	GLN
Mutation	ALA		-3.5		-9.8	-50	-24	-0.1				-9.9	-20
	ARG	-71	-48	-37			0	-52	0				-37
	ASN	-20	0	-11	-13	-30	-25	-17			0	-16	-30
	ASP			-12	-14	0	-19					0	-29
	GLN	-21	-4.8		-5.4	-44	-32	-36			-22		0
	GLU		-15	-13		-44		-30			-44	-30	
	HIS	-17			-17	-62	-64	3.46			-13	-75	-36
	ILE		-8.2	-20		-52	-24	-0.2				-41	-27
	LEU		1.57		-5.2	-67	-28					-40	-22
	LYS		-39	-47	-47	-115	-44	-82					-36
	MET		-4.1		-1.2	-53	-29						
	PHE		-18	0	0	-101	-60			0		-109	
	SER	0	-2.4		-11	-44	-26	0				-5.3	-22
	THR			-0.6		-55		1			-11	-19	-27
	TRP		-26			-160	-72						
	TYR		-24	-22	-14	-129	-65				-48	-81	
	VAL		-3			-56	-17					-23	-23

Stability reduced by more than 5 kcal/mol

-1.2 Affinity reduced by more than 5 kcal/mol

^{1.57} High affinity mutant

Effect of softening the van der Waals energy

The van der Waals energy is frequently softened so as not to penalize the small steric clashes resulting from limited sampling resolution. A side effect of this is to make hydrogen bonds appear stronger than they actually are.



Effect of softening the van der Waals energy

A softened van der Waals energy encourages the design algorithm to bury charges and polar residues at a designed hydrophobic interface:

	13	15	16	89	90	141	164	190	215	235
(ARG	GLU	MET	SER	ALA	MET	PHE	ASN	SER	SER
D d $(\hat{\mathbf{e}})$	ARG	GLU	MET	TYR	SER	MET	TYR	ASN	GLN	ASN
ib BB	ARG	GLU	LEU	ALA	LEU	ILE	PHE	ASN	ASN	SER
h =	ARG	GLU	THR	ALA	ASN	MET	ASN	GLU	ALA	ALA
, stelle	ARG	GLU	THR	ALA	ASN	MET	ASN	ASN	ASN	VAL
) si	HIS	GLU	LEU	MET	ASN	GLU	PHE	ASN	GLU	THR
igi∧	ARG	GLU	MET	TYR	SER	MET	TYR	ASP	GLN	ASN
DV	ARG	ASP	ALA	MET	ASN	MET	ASN	ASN	ASN	THR
	ARG	GLU	THR	ALA	ASN	MET	ASN	GLU	ASN	ALA
(ARG	PHE	LEU	ALA	THR	MET	PHE	ASN	SER	VAL
(ASN	ALA	MET	ALA	ASN	MET	PHE	ASN	ALA	ALA
	ASN	VAL	MET	ALA	ASN	MET	PHE	ASN	ALA	ALA
BB O	ASN	VAL	MET	ALA	ASN	MET	PHE	ASN	ALA	SER
h =	ASN	SER	MET	ALA	ASN	MET	PHE	ASN	ALA	ALA
, stelle	ASN	VAL	MET	SER	ASN	MET	PHE	ASN	ALA	ALA
) si	MET	VAL	MET	ALA	ALA	MET	PHE	ASN	ALA	ALA
≥ igr	ASN	VAL	MET	ALA	ASN	MET	PHE	ASN	SER	ALA
D	ASN	ILE	MET	ALA	ASN	MET	PHE	ASN	ALA	ALA
	SER	VAL	MET	ALA	ASN	MET	PHE	ASN	ALA	ALA
	ASN	VAL	LEU	ALA	ASN	MET	PHE	ASN	ALA	ALA

Tryptophan stacking in galactose binding proteins



Designing for specificity



(Each point represents a single sequence)

Assays for measuring K_d

Technique	Protein concentration	Limitations
tryptophan fluorescence	< 10 × K _d	Small signal, can't use absorbant ligand
calorimetry	$5-500 imes K_d$	∆H must be non-zero
spin concentrator	> K _d	Need radioligand
solid phase binding assay	40 µM to detect mM binders	Need radioligand

Assays for measuring K_d

Spin concentrator (quick alternative to radioligand equilibrium dialysis)



Future directions

Scaffold selection

Potential energy function

- explicit water
- lone pairs
- quantum effects
- better hydrogen bond model
- protein polarization

Structural sampling

- backbone flexiblity
- more design positions

Future application: Therapeutic proteins

Limitations of monoclonal antibodies

- Unable to bind to some targets, such as deep grooves in proteins
- Some antibodies are unstable or aggregationprone
- Non-human antibodies are immunogenic in humans
- Therapeutic antibodies are glycosylated and thus more expensive to manufacture
- Large size of whole antibodies may limit their tissue distribution

