JMB



Accurate Computer-based Design of a New Backbone Conformation in the Second Turn of Protein L

Brian Kuhlman¹[†], Jason W. O'Neill²[†], David E. Kim¹[†], Kam Y. J. Zhang² and David Baker^{1*}

¹Department of Biochemistry University of Washington Seattle, WA 98195, USA

²Division of Basic Sciences Fred Hutchinson Cancer Research Center, 1100 Fairview Ave., North Seattle WA 98109, USA The rational design of loops and turns is a key step towards creating proteins with new functions. We used a computational design procedure to create new backbone conformations in the second turn of protein L. The Protein Data Bank was searched for alternative turn conformations, and sequences optimal for these turns in the context of protein L were identified using a Monte Carlo search procedure and an energy function that favors close packing. Two variants containing 12 and 14 mutations were found to be as stable as wild-type protein L. The crystal structure of one of the variants has been solved at a resolution of 1.9 Å, and the backbone conformation in the second turn is remarkably close to that of the *in silico* model (1.1 Å RMSD) while it differs significantly from that of wild-type protein L (the turn residues are displaced by an average of 7.2 Å). The folding rates of the redesigned proteins are greater than that of the wildtype protein and in contrast to wild-type protein L the second β -turn appears to be formed at the rate limiting step in folding.

© 2002 Academic Press

*Corresponding authors

Keywords: computational protein design; β -hairpin design; protein folding; protein L

Introduction

There has been significant recent progress in the area of computational protein design. In most cases the goal has been to find new sequences for an already existing protein backbone;¹⁻⁴ there have been only a few studies in which computational methods have been used to create new protein backbones.^{5,6} A very notable success was the design of an α -helical coiled-coil with a novel twist and backbone conformation.⁷ In this case the backbone was moved along a parameterized surface, a method that is only suitable for symmetric systems like coiled coils.

We have developed a method that should be generally applicable towards designing new turns and loops in the context of full-length proteins. We searched the Protein Data Bank (PDB) for alternative hairpins with termini that superimpose with the region of interest, and then used a design program to identify sequences that are optimal for these backbone conformations in the context of the full-length protein. Low energy sequence-structure combinations are candidates for viable turns. In previous experiments on the IgG-binding domain of protein *G*, we found that this method could identify sequence changes that enhance stability: two design variants, containing 11 mutations each, were ~4 kcal mol⁻¹ more stable than wild-type protein *G*.⁸

Here, we focus on the structural accuracy of this backbone redesign strategy by redesigning the second turn of protein L, determining the crystal structure of the redesigned protein, and comparing it to the design model. We find that the method can accurately design a new backbone conformation: the root-mean-square deviation (RMSD) between the design model and the crystal structure for backbone atoms in the redesigned turn is 1.1 Å, and for buried residues the side-chain chi1 angles in the design model differ by only 9° on average from the chi1 angles in the crystal structure.

The 62 residue IgG-binding domain of protein L consists of a single α -helix packed on a fourstranded β -sheet formed by two symmetrically opposed β -hairpins⁹ (Figure 1). The first hairpin forms one of the four canonical turn types, a type I turn, while the second hairpin forms a four-residue

[†]B.K. and J.W.O'N. contributed equally to this work. Abbreviations used: PDB, Protein Data Bank; RMSD, root-mean-square deviation.

E-mail address of the corresponding author: dabaker@u.washington.edu



(1 1)

Figure 1. The successful design of a new backbone conformation in the second turn of protein L. C^{α} traces are shown for WT protein L (green), the crystal structure of L2 (blue), and the design model for L2 (yellow). The new turn is a two residue type I' turn while the WT turn contains four residues (labels are shown for turn residues).

turn that contains three consecutive residues with positive phi angles, two of which are non-glycine residues. Two lines of evidence suggest that the second turn has low intrinsic stability: the hairpin forms after the rate-limiting step in folding,¹⁰ and single point mutants in or near this turn result in a strand-swapped dimer of protein L in which the second turn is extended and the C-terminal strand inserts into the β -sheet of the partner.¹¹

Results

In order to increase the intrinsic stability of the second turn of protein L and change the backbone

conformation, we aimed to replace the wild-type four-residue turn with a canonical two residue turn (I, II, I' or II'). Two-residue turns are the most prevalent in the PDB and, in particular, I' and II' turns show a strong preference for β -hairpins.¹² In order to include a two-residue turn and still preserve the hydrogen-bond pairing in the β -sheet it was necessary to add or remove two residues from the protein L sequence. Alternate backbone conformations were identified by searching the PDB for canonical turns that have the appropriate number of residues and have termini that overlay well (RMSD < 0.5 Å) with strands 3 and 4 of protein L. The new turns, minus their side-chains, were grafted onto protein L using a conjugate gradient minimization procedure as described in Materials and Methods. A total of 150 turns containing two extra residues and 180 turns containing two fewer residues overlay well with strands 3 and 4. Six different groups of backbones were considered: type I, II, I' and II' turns with two more residues than protein L, and type I and II' turns with two fewer residues.

Sequences optimal for the new turns in the context of protein L were identified with a protein design procedure that uses a Monte Carlo search procedure to select amino acid rotamers that form well-packed structures.¹³ Unlike the wild-type backbone, the backbones with canonical turn types do not curl up towards the α -helix but rather bend away from the helix (Figures 1 and 2). Thus, in order to create good packing between the new turns and the helix it was necessary to consider new amino acid residues in the helix as well as in the turn. Four residues were redesigned in the helix and 8-12 residues were redesigned in the turn (Table 1). The lowest-energy sequence-structure combination for each backbone group was selected for experimental study. The sequences selected by



Figure 2. Detailed comparison of the crystal structure of L2 (blue) and WT protein L (green). In the wild-type structure the second turn bends up towards the helix while in the redesigned protein the second turn bends away from the helix.

WT Protein	L 26-FEKATSEAY-34	49-VDVADKG-YTL-60	
L1	AL	IEKVVS DN KYIF	I
		(YQS WR YS)	
L2	KLVL	IDKRVT NG VIIL	I'
		(VVV NG IR)	
L3	YL	IDKRYT PG ALIL	II
		(QVY KG LT)	
L4	LVL	IDKRQD GN VLVL	II'
		(KRL GD TF)	
S1	VL	IDRT DT RF	I
		(L SG R)	
S2	VL	IDRD GY LF	II'
		(E GD K)	

Table 1. Amino acid sequences of protein L variants

All residues in the turn (47-58) and residues 26, 30, 33 and 34 in the helix were allowed to vary in the design procedure. A dot indicates identity and a dash gap. Turn types are indicated to the right. Shown in brackets are the original sequences of the fragments that were used to generate the backbone coordinates of the new turns.

the design procedure differ considerably from the wild-type protein L sequence (between 8 and 14 mutations) and from the sequences of the fragments used to generate the new backbone coordinates (Table 1). At the two turn positions the design procedure picked amino acids that are commonly seen in the canonical turn types (N/D-G type I',G-N/D type II', PG type II).¹²

All six protein L variants expressed successfully and guanidine denaturation was used to evaluate their stability (Figure 3). Five of the variants are folded (S1 is only partially folded) and two of them, L2 and L4, have wild-type stability. The four variants containing two extra residues are more



Figure 3. The equilibrium denaturation of wild-type (WT) protein L compared to that of the six design varaints. L4 and L2 have stabilities comparable to that of wild-type protein L while the other four variants are destabilized. The free energies of unfolding (kcal mol⁻¹) are 4.6 (WT), 4.5 (L4), 4.4 (L2), 3.4 (L1), 2.5 (L3), 1.3 (S2) and ~0.5 (S1).

stable than the two variants with two fewer residues (S1 and S2), probably because they pack a larger surface area against the helix. In addition, the variants designed to have I' and II' turns are more stable than the variants of corresponding length that were designed to have type I and II turns. Unlike type I' and II' turns, type I and II turns do not have a strong preference to be in β -hairpins and it has been postulated that they are not as stable in hairpins as type I' and II' turns.¹² Our results are consistent with this hypothesis, but there may be other factors such as packing differences that account for the varying stability of the designed proteins as well.

The most stable variants, L2 and L4, were selected for structural studies. L4 precipitated out of solution at concentrations above 1 mM, and therefore we were not able to get suitable crystals. L2 was soluble above 3 mM, and we were able to obtain crystals that diffracted to 1.9 Å (Table 2). Except for the turn region, the structure of L2 is very similar to that of wild-type protein L (Figures 1 and 2). In the turn region the structures are dramatically different. In place of the wild-type four-residue turn, the mutant contains a two-residue (Asn53 and Gly54) type I' turn (Figure 2). Instead of curling up towards the helix, the turn bends away from the helix as designed. A superposition of L2 with the design model shows an excellent agreement (Figure 4). To more precisely compare L2, the design model and wild-type protein L, the three structures were superimposed using the backbone coordinates from residues that were not redesigned, residues 6 to 47 and residues 60 to 65 (58 to 63 in wild-type protein L). Using this structural alignment the backbone RMSD between the turn residues in L2 (54-57) and the turn residues in the design model (54-57) is 1.1 Å, while the backbone RMSD between the turn residues in L2 and the turn residues in wild-type pro-

Table 2. Data	statistics	for	L2
---------------	------------	-----	----

Unit cell	
a (Å)	38.62
b (Å)	55.54
c (Å)	67.38
Space group	$P2_{1}2_{1}2_{1}$
Resolution (Å)	1.9
Completeness (%)	96.9
$R_{\text{merge}}(\%)^{\mathbf{a}}$	4.9
Refinement statistics	
R _{crvst} ^b	19.7
R _{free}	22.7
Test size (%) ^c	10.0
No. molecules in asymmetric unit	2
No. of non-hydrogen atoms	
Protein	1004
Water	80
B-factor(Å ²)	20.6
RMSD from ideal values ^d	
Bond lengths (Å)	0.009
Bond angles (deg.)	1.5
Ramachandran plot (%) ^e	
Most favored regions	99.2
Additional allowed regions	0.8
Disallowed regions	0

^a $R_{\text{merge}} = \sum_{hkl} \sum_i (|I_{hkl}^i - \langle I_{hkl} \rangle|) / \sum_{hkl} (I_{hkl})$, where I_{hkl}^i is the intensity of an individual measurement of the reflection with Miller indicies h, k and l, and $\langle I_{hkl} \rangle$ is the mean intensity of that reflection.

^b $R_{cryst} = \sum_{hkl} (|F_{hkl}^o - F_{hkl}^c| / F_{hkl}^o)$ where F_{hkl}^o and F_{hkl}^c are the observed and calculated structure factor amplitudes.

 $^{\rm c}$ $R_{\rm free}^{23}$ is equivalent to $R_{\rm cryst}$ but calculated with reflections omitted from the refinement process. The $R_{\rm free}$ reflections were extracted using the CCP4 program, FreeRflag.

^d Calculated with the program CNS.²⁴

^e Calculated with the program PROCHECK.²⁰

tein L (53-56) is 7.2 Å. The all-atom RMSD between L2 and the design model for residues that were redesigned is 1.4 Å. Most of the deviations between L2 and the design model occur in the side-chains of solvent-accessible amino acids. For buried residues (less than 30% of the surface area solvent accessible) the phi, psi and chi1 angles in the design model differ by less than 9° on average from those in L2 (Table 3).

Although four out of the six design variants are less stable than wild-type protein L, all six of them

fold more quickly than the wild-type protein (Figure 5). In contrast, it was found previously that single point mutations in the second turn of wildtype protein L have little effect on the folding rate, indicating that the second turn of wild-type protein L is not formed in the folding transition state.¹⁰ The increase in folding rates in the new designed proteins suggests that the second turn may now be formed in the transition state. In order to investigate this possibility further we made point mutants in the context of the L4 sequence. Mutations in the second turn of L4 (G55A) and the first turn of L4 (G15A) both affect the folding rate equally (220 s^{-1} *versus* 340 s^{-1} for L4). This result can be contrasted with wild-type protein L where the G15A mutation reduces the folding rate tenfold and mutations in the second turn reduce the folding rate by less than twofold. It appears that the redesigned turn in L4 is now partially formed in the folding transition state. These results suggest that the redesigned hairpins have more intrinsic stability and make more favorable local interactions than the wild-type turn. The overall stabilities of the design variants are not greater than wild-type protein L, probably because the new turns do not pack against the rest of the protein as well as the wildtype turn.

To determine if we could more completely switch the folding mechanism of protein L, we destabilized the first hairpin with the G15A mutation and made further mutations to evaluate which interactions are present in the transition state. The effects of mutations can be summarized by the Φ -value notation developed by Fersht and coworkers¹⁴ (see Materials and Methods): a Φ value near 1 suggests that a residue is ordered in the folding transition state, while a value near 0 suggests that the residue is disordered in the folding transition state. G55A made in the context of L4/G15A has a Φ -value of 1.1 and N14A in the same context has a Φ -value of 0.4. In contrast, G55A in the context of wild-type protein L has a Φ-value of 0.2 and N14A in the context of wild-

Table 3.	Phi, Psi, and	Chi1	angles of modeled	amino acids c	ompared with a	ngles from	the X-ray	structure of L2
	, ,					a		

Replaced amino acid	Δ Phi	Δ Psi	Δ Chi1	Δ Chi2
Lys26	6	2	25	51
Leu30	0	2	20	14
Val33	0	7	10	
Leu34	1	0	17	4
Ile49	15	15	7	15
Lys51	51	14	8	9
Arg52	17	28	106	90
Val53	11	2	2	
Thr54	27	29	115	
Asn55	2	14	4	32
Gly56	10	10	0	
Val57	22	5	1	
Ile58	3	5	7	14
Ile59	0	6	2	2
Core Average	6	5	9	10

The values in the table are the absolute values of the changes in the dihedral angles. The core average is the average of the dihedral angle changes for those residues that have less than 30% of their surface area solvent exposed (30, 33, 34, 49, 57, 58, 59).



Figure 4. Detailed comparison between the crystal structure of L2 (blue) and the design model of L2 (yellow). Side-chain atoms are shown for mutated residues. The design model was successful at predicting the coordinates and rotamers of all buried or partially buried hydrophobic residues: L30, V33, L34, I49, V53, V57, I58 and I59. As would be expected more variation is seen in the exposed polar residues: K26, K51, R52, T54 and N55.

type protein L has a Φ -value of 0.7. These results suggest that in contrast to wild-type protein L, in L4/G15A the second β -turn is more formed in the folding transition state than is the first β -turn. Previously, we were able to introduce a similar change into a protein topologically related to protein L, protein G.⁸ These studies demonstrate that the intrinsic stability of substructures within a protein is an important determinant of folding pathways.

Conclusion

We have demonstrated that it is possible to rationally select sequences that will alter protein backbone conformation in a predictable manner with atomic level resolution. In order to accomplish this transformation, we used fragments from the PDB as a source of alternate backbone conformations. The advantage of this method is that it guarantees that the design should be feasible at the level of local structural preferences. Non-local interactions between the new fragment and the rest of the protein are optimized using the sequence design procedure. Our method should be generally applicable towards redesigning turns and loops in a variety of proteins. The main limitation of the current technique is that it requires fragments from the PDB that span the region of interest, and therefore will not be feasible for building large structures. We are currently testing methods that piece together protein fragments in order to create larger templates for protein design.

Materials and Methods

Design procedure

A non-redundant set of protein structures culled from the PDB was scanned for turns with termini that overlay well with strands 3 and 4 of protein L (http://www.fccc.edu/research/labs/dunbrack/culledpdb.html). The atoms used to check the overlap were all backbone atoms on residues 51 and 57, as well as the backbone nitrogen atoms on residue 50 and the carbonyl group of residue 56. The cutoff for a suitable match was a RMSD of 0.5 Å over these 11 atoms. The turns were also screened for the correct number of residues required to form the canonical two residue turns.

The turns were then grafted on to protein L using a multi-step process that began with a starting structure that was created by adding residues sequentially from the N terminus using wild-type protein L bond lengths, bond angles and dihedral angles in the portion of the structure derived from protein L and bond lengths, bond angles and dihedral angles derived from the turn residues in the turn portion of the structure. The connecting residues, two in each strand, were given ideal bond lengths and angles, and dihedral angles typical of resi



Figure 5. Refolding kinetics (k_{obs}) as a function of denaturant concentration (GuHCl). All six design variants fold more rapidly than wild-type protein L. The folding rate constants (s⁻¹) extrapolated to 0 M GuHCl are 60 (WT), 330 (L4), 280 (L2), 280 (L1), 180 (S2), and 100 (S1).

dues in β -strands. The starting structure was then input into a conjugate minimization procedure which favored low RMSD values between the protein L structure and residues with coordinates derived from the protein L structure (non-turn residues), low RMSD values between the turn residues and their original PDB coordinates, good hydrogen bonding between strands 1 and 4, and few bumps between atoms while preserving bond lengths and angles. Only backbone dihedral angles were allowed to vary in this procedure and the largest changes occurred in the residues connecting the turn with protein L. The hydrogen-bonding potential was a simple distance-based model that favored hydrogen-oxygen distances of 1.8 Å, hydrogen-carbonyl carbon distances of 2.8 Å, and oxygen-nitrogen distances of 2.8 Å.

New sequences were selected for the turn in the context of full-length protein L using a design procedure described.¹³ In brief, the design procedure uses a Monte Carlo search procedure to identify amino acid rotamers that have good Lennard-Jones energies, have low desolvation energies and are favorable for the backbone phi and psi angles. During the design process all amino acids, except for cysteine, were considered at residues 47-58, as well as residues 26, 30, 33 and 34 in the helix.

Plasmid construction and protein expression

DNA cassettes for the redesigned portion of protein L were created from two overlapping DNA fragments that were annealed and extended using *Klenow*. The resulting double stranded DNA products were digested using *Eco*RI and *Mlu*I, gel-purified, and cloned into a modified Protein L pET 15b construct.¹⁵ Point mutations were made using the QuikChange site-directed mutagenesis kit (Stratagene). Protein expression and purification were carried out as described.¹⁵ All mutants were verified by DNA sequencing and mass spectrometry.

Stability and kinetic measurements

Protein solutions were made in 50 mM sodium phosphate, pH 7, and the temperature was kept at 295 K. Guanidine-induced denaturation was monitored using the circular dichroism (CD) signal at 220 nm as described.^{10,16} The kinetics of folding were followed by fluorescence on a BioLogic SFM-4 stopped-flow instrument. Folding data were well fit by single exponentials. Φ-values were computed using:

$$\Phi = -RT \ln(k_{\rm f}^{\rm ref}/k_{\rm f}^{\rm mut})/\Delta\Delta G$$

where the change in stability ($\Delta\Delta G$) was determined from equilibrium unfolding experiments and $k_{\rm f}$ are the folding rate constants.

Crystallization, data collection, and structure determination

The purified L2 variant of protein L was dialyzed to 100 mM NaCl and 2 mM EDTA and then concentrated to ~25 mg/ml. Crystals of L2 were grown by hanging drop diffusion from 1.0 M sodium citrate and 100 mM cacodylate at pH 6.5. X-ray data was collected on an RAXIS-IV image plate at room temperature under Cu Ka radiation generated by a RIGAKU rotating-anode generator. Diffraction data were processed with Denzo & Scalepack.¹⁷ The program EPMR¹⁸ was used to find the molecular replacement solution using the wild-type protein L structure⁹ with the second β -turn deleted as a template. Simulated annealing composite omit $2F_{o} - F_{c}$ maps were used for model rebuilding with Xfit.¹⁹ Structural refinement was performed using CNS.¹² The L2 structure is validated by PROCHECK,²⁰ VERIFY-3D²¹ and ERRAT.²² There were no violations detected by these methods.

Protein Data Bank accession number

The coordinates have been deposited in the RCSB Protein Data Bank, with accession code 1KH0.

Acknowledgments

We thank Tanja Kortemme and Carol Rohl for helpful comments on the manuscript. B.K. was supported by the Cancer Research Fund of the Damon Runyon-Walter Winchell Foundation fellowship. This work was also supported by a grant from the NIH and HHMI.

References

- Dahiyat, B. I. & Mayo, S. L. (1997). *De novo* protein design: fully automated sequence selection. *Science*, 278, 82-87.
- Desjarlais, J. R. & Handel, T. M. (1995). *De novo* design of the hydrophobic cores of proteins. *Protein Sci.* 4, 2006-2018.
- Marvin, J. S. & Hellinga, H. W. (2001). Manipulation of ligand binding affinity by exploitation of conformational coupling. *Nature Struct. Biol.* 8, 795-798.
- Shimaoka, M., Shifman, J. M., Jing, H., Takagi, J., Mayo, S. L. & Springer, T. A. (2000). Computational design of an integrin I domain stabilized in the open

high affinity conformation. Nature Struct. Biol. 7, 674-678.

- Su, A. & Mayo, S. L. (1997). Coupling backbone flexibility and amino acid sequence selection in protein design. *Protein Sci.* 6, 1701-1707.
- 6. Desjarlais, J. R. & Handel, T. M. (1999). Side-chain and backbone flexibility in protein core design. *J. Mol. Biol.* **290**, 305-318.
- Harbury, P. B., Plecs, J. J., Tidor, B., Alber, T. & Kim, P. S. (1998). High-resolution protein design with backbone freedom. *Science*, 282, 1462-1467.
- Nauli, S., Kuhlman, B. & Baker, D. (2001). Computer based redesign of a protein folding pathway. *Nature Struct. Biol.* 8, 602-605.
- O'Neill, J. W., Kim, D. E., Baker, D. & Zhang, K. Y. J. (2001). Structures of the B1 domain of protein L from Peptostreptococcus magnus with a tyrosine to tryptophan substitution. *Acta Crystallog. sect. D*, 57, 480-487.
- Kim, D. E., Fisher, C. & Baker, D. (2000). A breakdown of symmetry in the folding transition state of protein L. J. Mol. Biol. 298, 971-984.
- O'Neill, J., Kim, D., Johnsen, K., Baker, D. & Zhang, K. (2001). Single site mutations induce 3D domain swapping in the B1 domain of Protein L from *Peptostreptococcus magnus. Structure (Camb)*, 9, 1017-1027.
- Sibanda, B. L. & Thornton, J. M. (1985). Beta-hairpin families in globular proteins. *Nature*, **316**, 170-174.
- 13. Kuhlman, B. & Baker, D. (2000). Native protein sequences are close to optimal for their structures. *Proc. Natl Acad. Sci. USA*, **97**, 10383-10388.
- 14. Fersht, A. (1999). *Structure and Mechanism in Protein Science*, W.H. Freeman and Company, New York.
- 15. Gu, H., Yi, Q., Bray, S. T., Riddle, D. S., Shiau, A. K. & Baker, D. (1995). A phage display system for

studying the sequence determinants of protein folding. *Protein Sci.* **4**, 1108-1117.

- Scalley, M. L., Li, Q., Gu, H., McCormack, A., Yates, J. R. & Baker, D. (1997). Kinetics of folding of the IgG binding domain of peptosteptococcal protein L. *Biochemistry*, 36, 3373-3382.
- Otwinowski, Z. & Minor, W. (1997). Processing of X-ray diffraction data collected in oscillation mode. In *Macromolecular Crystallography* (Charles, W., Carter, W. & Sweet, R. M., eds), vol. 276, pp. 307-326, Academic Press, San Diego.
- Kissinger, C. R., Gehlhaar, D. K. & Fogel, D. B. (1999). Rapid automated molecular replacement by evolutionary search. *Biol. Crystallog. sect. D*, 55, 484-491.
- McRee, D. E. (1992). A visual protein crystallographic software system for X11/XView. *J. Mol. Graph.* 10, 44-46.
- Laskowski, R. A., MacArthur, M. W., Moss, D. S. & Thornton, J. M. (1993). PROCHECK: a program to check the stereochemical quality of protein structures. J. Appl. Crystallog. 26, 283-291.
 Lüthy, R., Bowie, J. U. & Eisenberg, D. (1992).
- Lüthy, R., Bowie, J. U. & Eisenberg, D. (1992). Assessment of protein models with three-dimensional profiles. *Nature*, 356, 83-85.
- 22. Colovos, C. & Yeates, T. O. (1993). Verification of protein structures: patterns of nonbonded atomic interactions. *Protein Sci.* 2, 1511-1519.
- 23. Brünger, A. T. (1992). Free *R* value: a novel statistical quantity for assessing the accuracy of crystal structures. *Nature*, **355**, 472-475.
- Brünger, A. T., Adams, P. D., Clore, G. M., DeLano, W. L., Gros, P., Grosse-Kunstleve, R. W., Jiang, J. S. *et al.* (1998). Crystallography & NMR system: a new software suite for macromolecular structure determination. *Acta Crystallog. sect. D*, 54, 905-921.

Edited by F. Cohen

(Received 25 June 2001; received in revised form 12 October 2001; accepted 12 October 2001)