## Loops in Globular Proteins: A Novel Category of Secondary Structure

JACQUELYN F. LESZCZYNSKI AND GEORGE D. ROSE\*

The protein loop, a novel category of nonregular secondary structure, is a segment of contiguous polypeptide chain that traces a "loop-shaped" path in three-dimensional space; the main chain of an idealized loop resembles a Greek omega ( $\Omega$ ). A systematic study was made of 67 proteins of known structure revealing 270 omega loops. Although such loops are typically regarded as "random coil," they are, in fact, highly compact substructures and may also be independent folding units. Loops are almost invariably situated at the protein surface where they are poised to assume important roles in molecular function and biological recognition. They are often observed to be modules of evolutionary exchange and are also natural candidates for bioengineering studies.

The secondary structure of proteins falls into three classes:  $\alpha$ -helices,  $\beta$ -sheet, and reverse turns (1-4). Helices and sheet are termed "regular" structures because their residues have repeating main-chain torsion angles, and their backbone N-H and C=O groups are arranged in a periodic pattern of hydrogen bonding (1). In contrast, turns are "nonregular" structures with nonrepeating backbone torsion angles and, at most, one internal N-H...O=C hydrogen bond (1-4). Remaining residues, by subtraction, are often classified as "random coil," although, as Richardson has pointed out, they are neither random nor coil (1).

In this article we examine another category of nonregular secondary structure—the loop. A loop may be described as a continuous chain segment that adopts a "loop-shaped" conformation in threedimensional space, with a small distance between its segment termini. The main-chain trace of an idealized loop resembles a Greek omega ( $\Omega$ ). Backbone torsion angles for such a structure are nonrepeating, and there are few, if any, backbone hydrogen bonds. A simple loop subsumes no proper subsets that are also loops, while a compound loop contains at least one smaller embedded loop. Only simple loops are considered in the following discussion.

Loops have been discussed in relation to specific structures, such as the conspicuous loops in superoxide dismutase (5) and in immunoglobulin domains (6), the autolysis loop in serine protease zymogens (7), and the calcium-binding loops in parvalbumin (8). However, there has been no systematic study of these structures. Kuntz alluded to larger loops in his definitive paper on peptide chain turns (2), and the topic is mentioned briefly in a recent review (3). Looped-out regions are also evident in schematic representations of protein structure such as those of Richardson (1) or Lesk and Hardman (9). While such examples are clearly "looplike," their description has been only qualitative.

In our study, loops are defined explicitly. Stringent defining

criteria are chosen deliberately to exclude those structurally ambiguous examples containing substantial amounts of regular secondary structure. The definition was implemented in the form of a computer algorithm and used to identify all loops in 67 proteins of known structure. The set of identified loops was then characterized with respect to residue composition, size and shape, compactness, accessibility to solvent, and role in protein taxonomy.

Our survey reveals an abundant population of loops, on the order of four per protein molecule. Almost always, they are situated at the molecular surface; often, they are implicated in molecular function. Most of these loops are highly compact, globular structures, with low x-ray temperature factors and a packing efficiency that rivals that of  $\beta$ -sheet. The observed compactness is a consequence of loop sidechain atoms that pack tightly within the loop core. In view of such characteristics, the description of these chain segments as "random coil" warrants revision.

Loops are choice candidates for protein bioengineering studies. The catalog of loops presented here should be useful for the design of such experiments, as well as in the further study of nonregular protein secondary structure.

Identification of loops from x-ray coordinates. A loop is a continuous segment of polypeptide chain that is defined in terms of its (i) segment length, (ii) absence of regular secondary structure, and (iii) distance between segment termini. These criteria are now specified in detail.

The segment length must be between 6 and 16 residues. The lower length limit serves to eliminate reverse turns. Superficially, it might seem that a turn is merely a small loop, but an important characteristic distinguishes the two. Turns, which range from three to five residues in length, have backbone groups that pack together closely, forcing side chains to project outward (3). This stereochemical restriction is relaxed in larger segments where side-chain atoms can pack within the loop's own core. The upper length limit imposes a practical threshold that eliminates most of the compound loops.

A loop may contain no regular secondary structure. This criterion excludes adjacent strands of antiparallel  $\beta$ -sheet as well as structurally ambiguous cases. Secondary structure assignments for the residues were taken from the Kabsch and Sander (K&S) dictionary of protein secondary structure (10). However, two minor exceptions to the K&S classification were adopted: two-residue strands and single turns of helix (four or five residues) are not counted as regular secondary structure. Although they are ignored in most classification schemes, K&S includes  $\beta$ -strands that are just two residues in length. Strand lengths are distributed in a statistically well-behaved fashion, with the exception of these two-residue strands. Four- and

The authors are members of the Department of Biological Chemistry, Milton S. Hershey Medical Center, Pennsylvania State University, Hershey, PA 17033.

<sup>\*</sup>To whom correspondence should be addressed.

Table 1. Summary of 270 omega loops in 67 x-ray elucidated proteins.\*

PROT	FIRST NU	M SEQUENCE	PROT	FIRST NU	M SEQUENCE	PROT	FIRST	r Ni	UM SEQUEN	CE PRO	FIRST N	UM SEQUENCE
1ABP	93 7	V NK P	2CAB	78 10	V LD S	2GCH	112	7	A ST	V 3PG	4 11 15	SEUDV
1ABP	142 7	A N T A	2CAB	98 7	G SH G	2GCH	165	12	N TK	T 3PG	v 98 12	A O. KF
1ABP	203 6	G MS T	2CAB	108 7	т VК У	2GCH	217		5 5 5	T 3PG	x 109 12	FN PP
1ABP	236 13	A VG F	2CAB	128 13	Y SD G	1GPD	47	6	D SG	V 3PG	4 123 B	TD FS
1ABP	289 6	I TN F	2CAB	197 8	S Loo L Y	1GPD	76	7	E M N		N 132 14	KG VL
1ABP	299 6	FK LG	2CAB	230 11	LS. VP	1GPD	121	ģ	DS F		4 200 16	
2ACT	8 6	RS AV	1040	5 12	WG HW	1GPD	128	10	FV. K	V JrG	v 69 10	
2007	58 7		1040	17 7	нк та	1620	183	16	кт р	G 200	N 04 0	GENE
2ACT	89 15		1040	98 6		2685	103	7		G 2PT	N 94 9	Y NN D
2101	139 6		1010	108 7		2010	139	á		V 2PT	N 112 /	A 5R V
2801	141 16		1040	120 13		2010	160	11		1 2PT	N 142 11	G N Y P
2801	192 11		1000	166 7		2010	220		I PA	5 2PT	N 184A 8	G IK D
2801	102 11	N SE G	1CAC	107 9		2GRS	239		E NE	V 2PT	N 217 8	5 GN K
2801	190 0	R NGI	1040	197 0	5 LL L	2GR5	250	7	K TG	L IRE	1 91 6	Y QP Y
ANDU	203 7	A GI A	2CUA	232 8		2GRS	200			P IRH	5 34 10	S WE A
	100 12		2011A	70 9		2GRS	215	8	L NQ	T IRH.	43 15	A R5 F
4 10 10	115 0		2011A	114 6		2GRS	212	7	v DQ	N IRH	0 60 14	1 EV M
4600	113 8		2CHA	217 0	r 5v 5	2GR5	331	12		L IRH	0 85 6	G SI S
ANDU	122 /	1 M1 S	2014	217 8	5551	ZGRS	404	12	T PK	T IRH.	99 /	N GG S
AADH 280K	202 0	C Q G	JCNA	13 9	P NP 5	2665	405	7	A 15	E IRH	185 /	G RT Q
	217 0	G ED N	JONA	97 0 110 0			20	14	N QK	S IRH	193 /	E PG L
TAPA	21/ 8	N VN N	JCNA	116 8	K 5Q T	IHIP	28	14	R VE	Q IRH	D 216 8	L TE K
ZAPP	41 15	F SS V	3CNA	14/ 9	T TL E		43		C AF	M 1RH	D 284 10	P EK G
ZAPP	129 8	N TS Q	3CNA	160 6	S SS P	ITHIP	44	16	A DD	E IRN	5 36 6	т кс к
ZAPP	139 11	F DQ P	3CNA	199 11	I KD G	4LDH	1/3	16	R YG	V 1RN	5 87 10	T GC A
ZAPP	184 9	V DW S	JCNA	222 14	P SP D	4LDH	192	9	I GV	P 2RX	N 18 11	G XG Т
IAPR	8 10	T DY Y	3CNA	229 9	L LA N	4LDH	203	16	W SL	G 2RX	N 38 8	V CV G
IAPR	18 14	G QN L	3CPA	128 14	к тG V	T4LDH	212	14	L HD	W 1SB	r 17 6	Н SУ Т
1APR	43 16	G SD K	3CPA	142 15	D AG A	4LDH	219	8	N KW	K ISB	r 37 8	S SK V
1APR	61 9	P SK A	3CPA	156 11	A SY H	4LDH	239	.8	V IY	T ISB	r 74 13	A LA P
1APR	76 8	I GS A	3CPA	205 9	P YS I	4LDH	275	11	V KN	V 1SB	r 96 6	L GG S
1APR	90 14	D TG P	3CPA	231 7	к 5т 5	ILDX	70	9	S LK	I ISB	r 157 8	G SS T
1APR	129 10	D TS S	3CPA	244 7	I TQ A	ILDX	/9	8	V GS	L ISB	r 181 7	D SR A
1APR	189 9	I DW A	3CPA	272 14	R DS Q	ILDX	102	7	Q Qs	R ISB	r 257 10	L GK G
1APR	203 9	A TL G	lCPV	18 6	с кр s	1LDX	193	8	G RG	V 25G	A 16 16	I AS L
1APR	216 11	A IL I	1CPV	64 14	K LA L	1LDX	207	7	N NL	Q 25G	A 93 7	S FD Y
1APR	227 6	L PA A	1CRN	33 12	I ID Y	1 LDX	211	6	N LG	M 2SG	A 218 7	G NG G
1APR	233 16	V GL G	1CTX	1 15	I RC P	1LDX	218	7	W EE	G 3SG	3 16 16	I SS L
1APR	243 8	Q DG F	1CTX	26 10	C DG K	1LDX	236	1	ΑΥΥ	E 3SG	3 48 8	V RY Y
1APR	261 13	S IE I	1CYC	18 15	H TN L		276	8	к Ек	E 3SG	3 66 9	W AT V
1APR	280 8	А ЕС Т	1CYC	30 14	P NQ A	T1LH1	41	13	K DE	V 3SG	3 93 7	S FD Y
1APR	291 9	G AA I	1CYC	40 15	T GK S	1LH1	47	8	L KV	P 3SG	3 118 7	T VD I
1AZU	97	G NQ F	1CYC	70 15	N PA G	1LHB	46	14	P AL	T 3SG	3 167 14	A TG M
1AZU	35 12	Н РС Н	3CYT	18 15	H TN L	1LHB	55	10	F KE	L 3SG	3 190 12	V CP L
1AZU	67 6	G LD Y	3CYT	34 10	g lQ A	7LYZ	18	8	D NS	L 3SG	3 199 9	L YI G
1AZU	73 11	L KA H	3CYT	40 15	T GK S	7LYZ	36	7	s nQ	A 3SG	3 235 6	L VG V
lazu	84 9	Т КЕ К	3CYT	70 15	N PA G	7LYZ	44	9	N RT	D 2SN	5 43 10	E TV E
1azu	112 7	СТН S	1ECD	33 10	S IF A	7LYZ	60	16	S RN	L 2SN	5 114 6	V YN N
2B5C	32 16	L TL R	1ECD	41 9	F AS I	1LZM	134	6	A KW	Y 25N	5 136 6	K LW S
1BP2	23 8	N NC G	lest	69 12	G ET E	1MBN	40	8	L EF	K 2501	50 9	D NS A
†18P2	25 15	Y GV D	lest	94 11	W NY D	1MBS	37	14	P EL	K 2501	67 12	к к Я Н
1BP2	56 11	к кV D	lest	112 7	V TY V	1MBS	49	6	L KD	D 2501	0 103 7	S LY S
2BP2,	23 8	N NC G	1EST	142 10	G LL A	1MBS	78	7	к ке	A 250	0 122 16	D DG N
†2BP2	25 15	Y GV D	1EST	165 14	Y AT V	2MHB	40A	9	К ТD	L †250	0 132 6	S TG N
2BP2	61 8	С К Р	lest	216 11	V SR K	2MHB	39B	16	Q RA	v 255	197	G VT A
156B	16 10	V IK A	3FAB	24L 6	G SN I	2MHB	<b>4</b> 7B	11	D LG	N 1TI	4 67A 13	Y KI S
156B	47 12	Т РР М	3FAB	122L 11	P SK A	1NXB	6	8	Q НQ	T ITH	4 169A 6	A IG K
351C	16 11	Н АР А	3FAB	168L 6	к QN К	2PAB	49A	6	T SG	E 3TLI	1 24 8	Y SL Q
351C	51 12	G SM P	3fab	182L 6	L TQ W	8PAP	8	6	R QA	V 3TL	N 32 7	D ND G
155C	21 8	I QT D	3FAB	72H 6	N TN Q	8PAP	60	8	S YG	Y 3TL	44 10	A KG S
155C	47 8	А ЅК Ү	3FAB	99H 7	L II D	8PAP	86	15	Y PE	K 3TLI	55 16	W AP A
155C	83 13	K PG A	3FAB	132H 9	S KT A	8PAP	138	16	G KP	C 3TL	91 7	L SN N
155C	128 6	J JJ J	1FDX	12 12	G AI I	8PAP	175	11	N SN	G 3TLI	125 6	G DT F
2C2C	18 16	H TL F	1FDX	30 12	I DS C.	8PAP	191	8	R GY	G 3TLI	N 188 16	I GL R
2C2C	30 14	Р №Н К	1FDX	39 12	G SA P	8PAP	198	6	G VL	Y 3TL	204 10	S MG D
2C <b>2</b> C	41 16	А НМ К	3FXN	54 8	S AV L	1PCY	6	8	G AL	A 3TL	214 6	P DS K
2C2C	74 16	P KK S	2GCH	70 9	E FS E	1PCY	41	16	F DI	S 3TL	221 13	Y TI N
2CAB	67	G YN G	2GCH	94 8	Y NN N	1PCY	63	6	L NG	E 3TL	248 8	G TS V
2CAB	17 8	S KA N				1PCY	84	9	C SG	м		

SCIENCE, VOL. 234

five-residue helices are classified as type III reverse turns. Again, the distribution of helical segment lengths is statistically well behaved, except for these single-turn helices.

The distance between segment termini, that is, the end-to-end distance, is measured as the distance from the first  $\alpha$ -carbon to the last  $\alpha$ -carbon in the segment. The end-to-end distance must be less than 10 angstroms and may not exceed two-thirds the maximum distance between any two  $\alpha$ -carbons within the segment under consideration. This criterion selects as loops those segments with termini that "neck in" like an omega ( $\Omega$ ). The set of discovered loops is not overly sensitive to the coefficient of two-thirds. In practice, the end-to-end distance varies between 3.7 and 10.0 Å.

Loops identified as described above are frequently members of small families. Such families arise whenever a range of segments, all of similar length, satisfies the definition. For example, if residues i through j comprise a loop, then it often happens that residues i through j + 1 also comprise a loop. In our study, each family is represented by its most compact member. To choose these individual representatives, we evaluated the compactness of each segment in every family, and the most compact loops were selected.

The coefficient of compactness of Zehfus and Rose (11) was used to assess compactness. This coefficient, Z, is a sensitive single-value figure of merit that identifies those segments with the smallest solvent-accessible surface area for their volume. Explicitly,

$$Z = \frac{\text{accessible surface area of segment}}{\text{accessible surface area of sphere of equal volume}}$$
(1)

Solvent-accessible surface areas and volumes were calculated by the methods of Lee and Richards (12) and Pavlov and Federov (13), respectively.

Z is a dimensionless ratio and should show no dependence on unit size. However, the configurational freedom of very small segments is restricted in comparison to larger ones, and is thus biased toward more compact arrangements. To adjust for this apparent size dependency, a compensating correction term was applied to all Z values used in our study. This term,

$$0.488 \times e^{(-.068 \times N_r)} + 0.970$$

(where  $N_r$  = number of residues in the segment), is chosen to yield

Fig. 1. Histograms showing the distribution of loop sizes for the 270 loops in Table 1 by number of residues (A) and by end-to-end distance (B). The mean values are: (A) 9.8 and (B) 6.4.

a standard normalized value when multiplied by Z values of the most compact units of all sizes (11).

To identify loops, all continuous segments in 67 proteins from the Brookhaven database (14) were screened. (The Brookhaven database, a U.S. government supported resource, maintains atomic coordinates of x-ray elucidated proteins.) Only proteins from the K&S dictionary (10) were used in order to ensure self-consistency in the secondary structure assignments. Those continuous segments 6 to 16 residues in length that satisfied the end-to-end distance criterion were retained, if devoid of regular secondary structure. Coefficients of compactness were then calculated for all survivors and the most compact segment was chosen to represent each family cluster.

**Characterization of loops**. Compact loops are common structures in proteins. In the 67 proteins included in our survey, 270 loops were found, an average of more than four per molecule. The distribution of loop sizes is shown in Fig. 1 and the full set of loops is listed in Table 1.

Of those examined, only six proteins are without loops entirely: glucagon, insulin, mellitin, ovomucoid, avian pancreatic polypeptide, and pancreatic trypsin inhibitor. These proteins are all less than 60 amino acid residues in length and are among the nine smallest proteins in the database. Three, glucagon, mellitin, and pancreatic polypeptide, are also nonglobular, as determined by their axial ratios (as is discussed below).

Structures perceived intuitively as loops may not satisfy our stringent definition; the calcium-binding loops of parvalbumin (8) are an example. When the defining criteria are relaxed slightly to avoid elimination of loops containing three-residue  $\beta$ -strands, 22 additional loops are identified (Table 2), including several of these common examples.

Figures 2 and 3 illustrate typical loops from cytochrome c (15) and thermolysin (16). In each case, the loop main chain surrounds an internal cavity that is packed with side-chain atoms from loop residues. This kind of arrangement results in a highly compact chain fold for the segment. Occasionally, a metal ligand is also included, as shown in the thermolysin loop. In proteins with multiple loops, the individual loops occur frequently in spatial clusters, as they do in superoxide dismutase (Fig. 4).

Some loops with irregular tails satisfy our defining criteria both with and without their tail segment; in these cases the most compact representative is chosen, as previously stated. Occasionally, however, the larger version of such a loop, although more compact, is excluded because it exceeds the 16-residue upper limit. For example, the superoxide dismutase loop 67-78 is part of a larger, more compact loop 61-80. Despite this restriction, the upper size limit rarely eliminates a loop entirely; when the threshold is extended to 30 residues, only one de novo loop is found (phosphoglycerate mutase 191–211). As was mentioned above, the 16-residue cutoff serves to eliminate compound loops. If we use this upper bound,

<sup>\*</sup>Key. PROT, Brookhaven name of protein (14). FIRST, residue number of loop NH<sub>2</sub>terminus; the Brookhaven numbering system, which need not correspond to the numbering in the Kabsch and Sander dictionary (10), was used. NUM, number of residues in loop. SEQUENCE, first, second, penultimate, and ultimate residues in loop. Single letter abbreviations for the amino acid residues are: A, Ala; C, Cys; D, Asp; E, Glu: F, Phe; G, Gly; H, His; I, Ile; K, Lys; L, Leu; M, Met; N, Asn; P, Pro; Q, Gln; R, Arg; S, Ser; T, Thr; V, Val; W, Trp; Y, Tyr; X, Asx; J, unknown. – †Indicates a compound loop. Proteins used (and their parenthesized Brookhaven file names) are: 1arabinose-binding protein (1ABP), actinidin (2ACT), alcohol dehydrogenase (4ADH), adenylate kinase (2ADK), alphalytic protease (1ALP), penicillopepsin (2APP), rhizopuspepsin (1APR), azurin (1AZU), cytochrome b5 (2B5C), phospholipase A2 (1BP2), prophospholipase A2 (2BP2), cytochrome b52 (156B), cytochrome c551 (351C), cytochrome c550 (155C), cytochrome c (2C2C), carbonic anhydrase B (2CAB), carbonic anhydrase C (1CAC), alpha chymotrypsin (2CHA), concanavalin A (3CNA), carboxypeptidase (3CPA), calcium-binding parvalbumin (1CPV), crambin (1CRN), alpha cobratoxin (1CTX), ferrocytochrome c (1CYC), cytochrome c (3CYT), erythrocruorin (1ECD), tosylelastase (1EST), lambda immunoglobulin Fab NEW (3FAB), *Peptoaccus* ferredoxin (1FDX), flavodoxin (3FXN), gamma chymotrypsin (2GCH), glucagon (1GCN), glyceraldehyde-3-phosphate dehydrogenase (1GPD), glutathione reductase (2GRS), high potential iron protein (1HIP), insulin (1INS), apolactate dehydrogenase (4LDH), lactate dehydrogenase isoenzyme (1LDX), acetatemet-leghemoglobin (1LH1), lamprey methemoglobin (2MAB), papain (8PAP), plastocyanin (1PCY), phosphoglycerate mutase (3FGM), avian pancreatic polypeptide (1PPT), trypsin (2PTN), pancreatic trypsin inhibitor (4PTI), Bence-Jones immunoglobulin REI (1REI), rhodanese (1RHD), ribonuclease S (1RNS), rubredoxin (2RXN), subtilisin BPN (1SBT), *Streptomyces* proteinase A (2S

Fig. 2. Stereoview of a typical loop from cytochrome c (residues 40–54).



only seven compound loops fail to be excluded; these loops are indicated by a dagger in Table 1.

Loops may contain one or more reverse turns; these facilitate the main-chain direction changes needed to bring segment termini together. If loop curvature is sufficiently gradual, the chain direction can be reversed without resorting to an explicit turn, but this is unusual. All of the loops contain at least one turn or bend residue, as defined by Kabsch and Sander (10), but not every loop contains a complete turn or bend. Because they include reverse turns, loops do not constitute a pure structural category. Nonetheless, a loop and a turn are distinct moieties. While both result in changes in the overall direction of the polypeptide chain, a loop cannot be viewed merely as an "overgrown" turn. A turn, unlike a loop, has backbone groups that pack together closely, forcing side chains to project outward (3). For steric reasons, a segment of main chain cannot circumscribe an interior cavity of polyatomic dimensions until it exceeds a length of five residues.

The residue composition of loops was assessed by calculating the normalized frequency of occurrence, f, for each residue type, X, such that

$$f = \frac{X_{\rm L} X_{\rm T}}{N_{\rm L}/N_{\rm T}} \tag{2}$$

where  $X_L$  is the number of residues of type X in loops,  $X_T$  is the total number of residues of type X,  $N_L$  is the total number of residues in loops, and  $N_T$  is the total number of residues in the database. A value of f = 1 implies that X is distributed randomly in loops. Values greater than unity imply that X is found preferentially in loops; conversely, f values less than unity imply a less than average frequency of occurrence of X in loops.

Examination of f values for all residues in loops reveals that residues present most often in reverse turns (1-4, 17) are also found most often in loops (Gly, Pro, Asp, Asn, and Ser) with the notable addition of Tyr (Table 3). All but Tyr have short side chains, and all but Pro are polar; Pro favors turns for steric reasons (3). Hydrophobic residues are strongly disfavored in loops; these include Val, Met, Ile, Leu, and Ala with aliphatic side chains, and His, Trp, and Phe with aromatic side chains.

To quantify the accessibility of a loop, we calculated the solventaccessible surface area (12) in each of three successive states: in the standard state (18), as an isolated secondary structure, and within the protein (Fig. 5). Approximately half of the area of regular secondary structure is lost upon the formation of the isolated secondary structure, and the remaining half is lost when that secondary structure is buried within the protein (19). In our sample of 67 proteins, the percentage of the area lost when the chain folds into an isolated loop (34 percent) is comparable to the area loss upon formation of an isolated helix (35 percent). However, the subsequent loss when the loop is incorporated into the protein (47 percent) is less than that of the helix (60 percent). These statistics indicate that loops tend to be somewhat more accessible to solvent than helices.

These fractional accessibilities reveal that loops are almost invariably situated at the molecular surface (Figs. 4 and 5). It should be noted that the definition of a loop does not require that it be at the surface. Moreover, the data on solvent accessibility are consistent with the residue composition; loops are found at the protein surface and contain a preponderance of hydrophilic residues.

Loops are as compact as the proteins that contain them (Fig. 6). The coefficient of compactness, Z, is used to assess compactness, as described above. The Z values of loop segments range between 1.43 and 1.86 with a mean ( $\pm$  standard deviation) of 1.61  $\pm$  0.07. In comparison, Z values for the 67 proteins range between 1.36 and 1.93 (with a single outlier at 2.09); the mean of this distribution is 1.67  $\pm$  0.13.

It is conceivable that the apparent compactness of loops is biased by the use of the most compact segment to represent a loop family. As a control, the largest member of each family was chosen instead and used as the representative member; the coefficient of compactness was then calculated for these largest representatives. The distribution for these largest representatives is similar to that for compact representatives, ranging from 1.48 to 1.86 with a mean ( $\pm$ standard deviation) of 1.64  $\pm$  0.07. This control demonstrates that loops are inherently compact structures.

As would be expected from their observed compactness, loops are not flat, but globular. This visual impression is confirmed when we calculate the principal moments of inertia for all loops, helices,

Table 2. Additional loops found in 67 x-ray elucidated proteins with defining criteria relaxed to allow three-residue strands of  $\beta$  sheet. The key is the same as that for Table 1.

PROT	FIRST	r nu	JM S	SEQUENC	CE	PROT	FIRST	NU	JM S	SEQUENC	CE
4ADH	130	8	F	тр	I	1CPV	51	12	D	QD	E
4ADH	158	16	Α	кі	G	1CPV	89	9	G	DK	I
1ALP	190	12	Α	cs	W	1GPD	279	15	v	SF	D
1ALP	200	8	I	ΤΑ	Q	2GRS	370	7	v	VP	Ρ
2APP	212	11	G	IL	L	lHIP	65	15	L	FA	s
2APP	290	8	N	SL	I	4LDH	289	14	L	PI	v
1APR	31	12	L	NW	v	3PGM	166	11	I	AM	I
1APR	161	7	Α	AS	D	2RXN	5	8	т	су	I
1APR	169	11	D	FN	к	2SGA	119	7	Y	LS	Y
1APR	310	10	v	vi	R	2SGB	138	9	R	RT	Н
155C	36	14	Ρ	NS	Ε	1000	23	9	v	ст	Y

SCIENCE, VOL. 234

strands of sheet, and protein monomers within the set of 67 proteins. Ratios of the largest to the smallest eigenvalues were formed. The axial ratios for loops resemble those for whole proteins (Fig. 7), while ratios for helices and strands are more rodlike.

The free energy change upon closing a protein segment into a loop consists of an unfavorable entropic contribution and a compensating enthalpic contribution. The enthalpy needed to counterbalance loop-closing entropy may be due to either intrasegment or extrasegment interactions, or both. At one extreme, a loop might be stabilized by interactions within its own core. At the other extreme, the rest of the protein might provide a stable framework that pinches together the termini of an intervening segment, forcing that segment to "loop out." Greater structural autonomy would be expected in the former case.

The number of noncovalent contacts between loops can be used to provide a rough estimate of loop enthalpy. Using united-atom radii (12), we plotted the number of noncovalent contacts as a function of the number of atoms in the loop, and obtained a line described by the equation ( $\pm$  standard error):

Number of contacts =  

$$3.0 (\pm 0.03) \times \text{number of atoms} - 11.3 (\pm 2.7)$$

The equation can be interpreted to mean that loop enthalpy is essentially a linear function of loop length. (The negative intercept is expected upon extrapolation to zero length because a threshold of several atoms would be required to establish any contacts.) If we

Table 3. Residue frequencies in loops, normalized with the use of equation

Gly	1.35	Glu	1.09	Trp	0.85
Pro	1.28	Thr	1.07	His	0.83
Tyr	1.28	Lys	1.02	Ala	0.77
Asp	1.22	<sup>1</sup> / <sub>2</sub> Cys	0.94	Leu	0.76
Asn	1.22	Gln	0.93	Ile	0.68
Ser	1.20	Arg	0.91	Met	0.67
Cys	1.16	Phe	0.90	Val	0.64

assume a binding energy of -0.03 kcal/mol per contact, the average loop enthalpy is on the order of -0.6 kcal/mol per residue.

The entropy of loop closure scales linearly with the logarithm of segment length (20), but a confident numerical estimate requires theory that takes into account heterogeneous loops of approximately one statistical segment in length. While such an estimate is beyond the scope of this article, it is evident that the loop-closing entropy for these compact loops of 6 to 16 residues is offset, at least in part, by extensive favorable contacts within the loop.

Protein secondary structure has often been codified into a small number of states on the basis of backbone dihedral angles and hydrogen-bonding patterns. The usual categories include helix, sheet, reverse turn, and random coil. Identification of these categories is not always straightforward, and a given segment may be classified differently by different investigators. Not surprisingly, estimates of the relative abundance within these categories vary



(3)

Fig. 3. Typical loops from cytochrome c (residues 40-54) and thermolysin (residues 188-203). (A and B) Space-filling representation of the cytochrome c loop, with the same orientation as Fig. 2. Backbone atoms are shown in red and side-chain atoms in blue. The loop main chain forms an

internal cavity that is filled by side-chain groups from loop residues. (C and D) Space-filling representation of the thermolysin loop with backbone atoms shown in red and side-chain atoms in blue; (C) without the metal ligand, and (D) with the metal ligand, shown in green.



Fig. 4. Stereoview showing clustering of loops in superoxide dismutase.



Fig. 5. Histograms showing distribution of the percentage of surface area lost for loops on folding from (A) the standard state (18) to the isolated secondary structure and (B) the isolated secondary structure to the native protein. The mean values ( $\pm$  standard deviation) are: (A) 34.4 percent ( $\pm$  6.5) and (B) 46.7 percent ( $\pm$  14.6).

somewhat, particularly in the case of turns (3). However, a consensus estimate finds that regular secondary structures—helices and sheet—make up slightly less than half of proteins, on average. Chou and Fasman (17) allocate another third of all residues to turns, although other estimates are closer to a quarter (10, 21).

Two comprehensive studies that codify residues into discrete states based on objective criteria can be found (10, 22). We used the assignment of Kabsch and Sander (10) to eliminate regular secondary structure prior to loop identification, as discussed above.

The 67 proteins in our study contain 11,885 residues: 26 percent in helix, 19 percent in sheet, 26 percent in turns, and 21 percent in loops. On the basis of the K&S assignments, calculation of the percentage of residues in helix and sheet is straightforward, but assignment of residues to turns or loops is confounded because loops contain reverse turns. In that our statistics count such residues among loops, the percentage of residues in reverse turns must be reduced accordingly. When subdivided, 11 percent of all residues are found in turns within loops and 15 percent in turns external to loops. A further correction, although slight, should be made for those two-residue strands and single-turn helices that are counted both in loops and as regular secondary structure; in combination, these two minor categories contain less than 1 percent of all residues. Subject to these adjustments, helices, sheet, reverse turns, and loops account for approximately 80 percent of all residues in the 67 proteins of our study. When our conservative criteria for defining loops are relaxed only slightly, more than 90 percent of all residues are included in the accounting.

Loops as a definitive category for structure analysis. The principal question raised by our findings is whether loops comprise a distinct class of secondary structure. Of course, the classical definition of secondary structure as hydrogen-bonded backbone structure (23) automatically excludes loops. In practice, however, secondary structure has come to be synonymous with the conformation of continuous segments of the polypeptide chain (24).

The fact that loops exist in a range of conformations would seem to argue against their classification as a discrete category. Yet, the situation is not entirely dissimilar to that of reverse turns which can range between three and four residues and adopt multiple conformations (1-4). This range of variability increases exponentially with segment length, and, in loops, it is extremely large.

The conspicuous compactness of most loops is the factor that most convincingly underwrites their classification as a discrete entity. They are autonomously well-folded structures because their observed compactness does not depend on interactions with the rest of the protein. Indeed, were loops amorphous, their location at the molecular surface would render them ready targets for indiscriminate proteolysis (25), leading to rapid protein turnover, but there is no evidence for this.

Because loops contain reverse turns, they should perhaps be viewed as structural composites, akin to supersecondary structure (1, 26). In any event, loops can be identified objectively in x-ray elucidated proteins, and they occur with a frequency comparable to that of the  $\beta$ -sheet. We propose that protein segments which satisfy our defining criteria be called omega  $(\Omega)$  loops.



Fig. 6. Histograms showing distribution of the compactness coefficient (Z) for (A) the proteins used in this study and (B) the 270 loops from Table 1. The mean values ( $\pm$  standard deviation) are: (A) 1.67 ( $\pm$  0.13) and (B) 1.61 ( $\pm$  0.07).



Fig. 7. Histograms showing the distribution of axial ratios for (A) the proteins used in this study and (B) all loops from Table 1. The mean values ( $\pm$  standard deviation) are: (A) 1.77 ( $\pm$  0.66) and (B) 2.17 ( $\pm$  0.51). The three outliers in (A) are glucagon, mellitin, and avian pancreatic polypeptide, three of the six proteins without loops.

SCIENCE, VOL. 234

The term loop has also been used to describe a chain segment that is cross-linked by a disulfide bond. Although one can envision an  $\Omega$ loop with a disulfide bond between its ends, none are observed among the loops in Table 1. There are interloop disulfide bonds, but these are not situated between loop termini. In addition, the cystines in these omega loops form loop-loop and loop-protein disulfide bridges.

Some loop residues have been implicated in antibody binding. Lysozyme loops 18 to 25, 44 to 52, and 60 to 75 contain antigenic residues 19, 21, 45 to 48 (27), and 64 to 80 (28). Since the entire protein surface is thought to be potentially antigenic (29), loop involvement in antigenic sites may be a consequence of the fact that loops are on the protein surface.

It is unclear whether an isolated loop segment will fold independently in solution (30), but the proposition is testable. The questions are analogous to those raised by Bierzynski et al. (31) and Kim and Baldwin (31) in their studies of the independent stability of the C-peptide helix from ribonuclease. The isolated segment could be monitored for nativelike interactions by nuclear magnetic resonance. Alternatively, the conformation of loops containing suitable ligands might be probed with metal ions (32).

Isolated loop segments are also attractive peptides for use in model studies. They are highly solvated, both within the protein and alone in solution; and they are readily cyclized by addition of cysteines at their termini. An analogous use of cyclized peptides to model reverse turns has been successfully exploited by Gierasch and co-workers (3-4).

Omega loops are appealing candidates for bioengineering studies. From the data in Table 1, experiments could be designed to test the hypothesis that loops function as integral units and have the potential for modular exchange between proteins. Thus, a loop might be "swapped" or excised entirely, and the consequences for protein stability and enzymatic activity can be assessed.

Evolution appears to have established precedents for loop swap experiments. Several protein families in our database such as the cytochromes c and the serine proteases contain loops at homologous locations. Some of those homologous loops are structurally similar, while others have conserved end points, but differing overall structures.

Macromolecular recognition is a hallmark of biological systems. Recognition sites for glycosylation, phosphorylation, supramolecular assembly, and transport all reside on the protein surface. It is plausible that omega loops assume a central role in such processes.

## REFERENCES AND NOTES

- J. S. Richardson, Adv. Protein Chem. 34, 167 (1981).
   C. M. Venkatachalam, Biopolymers 6, 1425 (1968); I. D. Kuntz, J. Am. Chem. Soc. 94, 8568 (1972); P. N. Lewis, F. A. Momany, H. A. Scheraga, Proc. Natl. Acad. Sci. U.S.A. 68, 2293 (1971).

- G. D. Rose, L. M. Gierasch, J. A. Smith, Adv. Protein Chem. 37, 1 (1985). J. A. Smith and L. G. Pease, CRC Crit. Rev. Biochem. 8, 315 (1980). J. A. Tainer, E. D. Getzoff, K. M. Beem, J. S. Richardson, D. C. Richardson, J. Mol. Biol. 160, 181 (1982) 6
- D. R. Davies, E. A. Padlan, D. M. Segal, *Annu. Rev. Biochem.* 44, 639 (1975); J. S. Richardson, D. C. Richardson, K. A. Thomas, E. W. Silverton, D. R. Davies, *J.* Mol. Biol. 102, 221 (1976).
- A. A. Kossiakoff, J. L. Chambers, L. M. Kay, R. M. Stroud, Biochemistry 16, 654 7 (197

- (1977).
  8. R. H. Kretsinger and C. E. Nockolds, J. Biol. Chem. 248, 3313 (1973).
  9. A. M. Lesk and K. D. Hardman, Science 216, 539 (1982).
  10. W. Kabsch and C. Sander, Biopolymers 22, 2577 (1983).
  11. M. H. Zchfus and G. D. Rose, Biochemistry 25, 5759 (1986).
  12. B. K. Lee and F. M. Richards, J. Mol. Biol. 55, 379 (1971). Probe radius that we used was 1.4 Å. The atomic radii, in Å, are: tetrahedral C, 2.0; trigonal carbon, 1.7; carbonyl O, 1.4; hydroxyl O, 1.6; carboxyl O, 1.5; tetrahedral N, 2.0; trigonal N, 1.7; divalant S. J. St. and en ulthurdus S. 2005.

- M. Y. Pavlov and B. A. Federov, *Biopolymers* 22, 1507 (1983).
   M. Y. Pavlov and B. A. Federov, *Biopolymers* 22, 1507 (1983).
   F. C. Bernstein et al., J. Mol. Biol. 112, 535 (1977).
   T. Takano and R. E. Dickerson, *Proc. Natl. Acad. Sci. U.S.A.* 77, 6371 (1980).
   M. A. Holmes and B. W. Matthews, J. Mol. Biol. 160, 623 (1982).
   P. Y. Chou and G. D. Fasman, *ibid.* 115, 135 (1977); *Annu. Rev. Biochem.* 47, 251 (1978); M. Lawitt, Biochemistry 17, 4277 (1978).
- (1978); M. Levitt, *Biochemistry* 17, 4277 (1978). The standard state surface area of a residue, X, is the average area of that residue in a 18. representative ensemble of Gly-X-Gly tripeptides [see table 1 in G. D. Rose, A. R. Geselowitz, G. J. Lesser, R. H. Lee, M. H. Zehfus, *Science* **229**, 834 (1985)]. The standard state surface area of a segment is taken as the sum of its residue standard states
- States.
  C. Chothia, J. Mol. Biol. 105, 1 (1976), F. M. Richards, Carlsberg Res. Commun. 44, 47 (1979).
  P. J. Flory, Statistical Mechanics of Chain Molecules (Wiley, New York, 1969).
  G. D. Rose and J. P. Seltzer, J. Mol. Biol. 113, 153 (1977); G. D. Rose and D. B. Wetlaufer, Nature (London) 268, 769 (1977). 19.
- 20.
- M. Levitt and J. Greer, J. Mol. Biol. 114, 181 (1977).
- M. Levitt and C. Chothia, Nature (London) 261, 552 (1976).
  M. Levitt and C. Chothia, Nature (London) 261, 552 (1976). 23.
- 25
- 26. 27.
- S. J. Smith-Gill et al., J. Immunol. 128, 314 (1982); S. J. Smith-Gill, T. B. Lovoie, C. R. Mainhart, *ibid.* 133, 384 (1984). R. Arnon, Immunochemistry of Enzymes and Their Antibodies, M. R. J. Salton, Ed.
- 28. (Wiley, New York, 1977)
- (Wiley, New York, 1977).
  D. C. Berjamin et al., Annu. Rev. Immunol. 2, 67 (1984); D. W. Fanning, J. A. Smith, G. D. Rose, Biopolymers 25, 863 (1986); D. J. Barlow, M. S. Edwards, J. M. Thornton, Nature (London) 322, 747 (1986).
  D. B. Wetlaufer, Adv. Protein Chem. 34, 61 (1981).
  P. S. Kim and R. L. Baldwin, Annu. Rev. Biochem. 51, 459 (1982); A. Bierzynski, P. S. Kim, R. L. Baldwin, Proc. Natl. Acad. Sci. U.S.A. 79, 2470 (1982). 29.
- 31.
- F. S. Kini, K. L. Baldwill, Frot. Natl. Acta. Sci. U.S.A. 79, 2470 (1982).
  W. D. Horrocks, Jr., Progr. Inorg. Chem. 31, 1 (1984).
  We thank Micheal Zehfus for suggestions, Lyndon Hibbard and William Young for discussion during preliminary stages of this work, and an anonymous reviewer for valuable comments. Supported by NIH grants GM29458 and AGO6084.

6 May 1986; accepted 25 September 1986



At The AAAS Annual Meeting , 14-18 February 1987, in Chicago (See 12 December issue of SCIENCE for details)