

parate eukaryotic expres-
es for the heavy and light
nd V_L expression vectors
ATCC CCL10) using the
methotrexate and G418 as
clones that produce up to
er bottle cultures without

ied CEA and to CEA on
the original mouse MAb
ies, although it has not yet
of the mouse MAb, which

ne murine BW431/26 MAb
version. Surprisingly, 17 out
as well as the mouse MAb;
ill detectable. This finding
tween the variable domains
Bosslet *et al.*, unpublished

cid sequences of the donor
th the acceptor V_H (human
egree of sequence similarity
s are clustered in the regions
domains. Even in positions
each other the amino acid
have contributed to the suc-
ody. The crystal structures of
id we use their respective V
any murine MAb. Computer
potential framework amino
act with the CDRs or directly
ferred to the human fram-
ing scheme described here we
y particular monoclonal anti-

ne 73, 419 (1988).

Queen *et al.*³¹ have reshaped an antibody by selecting human acceptor V domains from the Kabat database to match the framework sequences of the murine donor V domains as closely as possible. This alternative reshaping strategy raises the possibility that any human antibody framework can be used as acceptor for the murine CDRs. Humanization of murine MABs is a technique that is now widely applied to MABs of potential use in diagnosis and therapy in humans. Whether one of the reshaping schemes will offer any particular advantage in successful humanization of MABs will become clearer as the number of reshaped antibodies increases.

Acknowledgments

Initial work was done by D. G. at the MRC/LMB, CB2 2QH Cambridge, England. We thank Dr. Greg Winter and P. T. Jones for their advice and help in establishing the humanization technique in our laboratory, and for their technical help with the computer analyses, Mrs. K. Müller and R. Bier for excellent technical assistance, Dr. K. Bosslet for helpful discussion, and Dr. H.-H. Sedlacek for support.

³¹ C. Queen, W. P. Schneider, H. E. Selick, P. W. Payne, N. F. Landolfi, J. F. Duncan, N. M. Avdalovic, M. Levitt, R. P. Junghans, and T. A. Waldmann, *Proc. Natl. Acad. Sci. U.S.A.* **86**, 10029 (1989).

[6] Molecular Modeling of Antibody Combining Sites

By ANDREW C. R. MARTIN, JANET C. CHEETHAM, and
ANTHONY R. REES

Introduction

Both the variable and constant domains of an antibody Fab consist of two twisted antiparallel β sheets which form a β -sandwich structure. The constant regions have three- and four-stranded β sheets arranged in a Greek key-like motif,¹ while variable regions have a further two short β strands producing a five-stranded β sheet.

The two β sheets of the variable domain are inclined at 30° to one another,² with a conserved disulfide bridge in each domain linking the two β sheets. Lesk and Chothia³ have shown the relative orientation of the two

¹ J. S. Richardson, *Adv. Protein Chem.* **34**, 168 (1981).

² C. Chothia and J. Janin, *Proc. Natl. Acad. Sci. U.S.A.* **78**, 4146 (1981).

³ A. M. Lesk and C. Chothia, *J. Mol. Biol.* **160**, 325 (1982).

β sheets to vary by up to 18° in variable domains and up to 10° in constant domains. This characteristic "immunoglobulin fold"⁴⁻⁶ is also seen in a number of molecules of related and unrelated function,⁷ one of the most recent examples being a bacterial protein, PapD, which mediates assembly of pili in *Escherichia coli*.⁸

The V_L and V_H domains interact via the five-stranded β sheets to form a nine-stranded β barrel of about 8.4 Å radius, with the strands at the domain interface inclined at approximately 50° to one another.⁹ The loops linking the β strands form the CDRs with the domain pairing bringing these loops from both V_L and V_H into close proximity. The CDRs themselves form some 25% of the V_L/V_H domain interface.¹⁰

The six hypervariable loops [complementarity determining regions (CDRs)] that are supported on the β -barrel framework form the antigen combining site (ACS). While their sequence is hypervariable in comparison with the rest of the immunoglobulin structure,¹¹ some of the loops show a relatively high degree of both sequence and structural conservation. In particular, CDR-L2 and CDR-H1 are highly conserved in conformation.¹² Analysis of conserved key residues has led Chothia and co-workers to define canonical ensembles into which the CDRs may be grouped.^{13,14}

Requirement for Modeling

Since the first X-ray crystallographic structure determinations, sequences of immunoglobulin light and heavy chain variable regions have accumulated at an ever increasing rate. X-Ray structures, however, are accruing at a relatively slow pace with no more than three or four immunoglobulin structures being published each year and fewer than that being deposited in the Brookhaven Protein Databank.¹⁵ Protein crystallography

⁴ R. J. Poljak, L. M. Amzel, H. P. Avey, B. L. Chen, R. P. Phizackerley, and F. Saul, *Proc. Natl. Acad. Sci. U.S.A.* **70**, 3305 (1973).

⁵ M. Schiffer, R. L. Girling, K. R. Ely, and A. B. Edmunds, *Biochemistry* **12**, 4620 (1973).

⁶ C. Chothia, J. Novotný, R. E. Bruccoleri, and M. Karplus, *J. Mol. Biol.* **186**, 651 (1985).

⁷ A. F. Williams and A. N. Barclay, *Annu. Rev. Immunol.* **6**, 381 (1988).

⁸ A. Holmgren and C.-I. Brändén, *Nature (London)* **342**, 248 (1989).

⁹ J. Novotný, R. E. Bruccoleri, J. Newell, D. Murphy, E. Haber, and M. Karplus, *J. Biol. Chem.* **258**, 14433 (1983).

¹⁰ P. M. Colman, *Adv. Immunol.* **43**, 99 (1988).

¹¹ T. T. Wu and E. A. Kabat, *J. Exp. Med.* **132**, 211 (1970).

¹² P. de la Paz, B. J. Sutton, M. J. Darsley, and A. R. Rees, *EMBO J.* **5**, 415 (1986).

¹³ C. Chothia and A. M. Lesk, *J. Mol. Biol.* **196**, 901 (1987).

¹⁴ C. Chothia, A. M. Lesk, A. Tramontano, M. Levitt, S. J. Smith-Gill, G. Air, S. Sheriff, E. A. Padlan, D. Davies, W. R. Tulip, P. M. Colman, S. Spinelli, P. M. Alzari, and R. J. Poljak, *Nature (London)* **342**, 877 (1989).

¹⁵ F. C. Bernstein, T. F. Koetzle, G. J. B. Williams, E. F. Meyer, M. D. Brice, J. R. Rodgers, O. Kennard, T. Shimanouchi, and M. Tasumi, *J. Mol. Biol.* **112**, 535 (1977).

is limited by two major factors: the difficulty of refining an X-ray data set and the limited availability of high-resolution X-ray data. The use of computer graphics analysis (for example, molecular models) is available in limited quantities.

Currently, the only technique available for the study of antibody structure is nuclear magnetic resonance (NMR) spectroscopy. This technique is applicable to proteins of molecular weight (<20 kDa).¹⁸ Thus, the importance of antibody structure and thus circumventing the need for a reliable modeling technique. The effects of site-directed mutagenesis on the antigen binding site (CDR region) and metal binding sites at antigen sites.

Antibody Modeling—

Various workers have been using, in the main, X-ray crystallography in attempts to develop a reliable modeling technique. The effects of site-directed mutagenesis on the antigen binding site (CDR region) and metal binding sites at antigen sites.

Modeling of antibody classes:

1. modeling panels of antibody classes either by crystallography or NMR

¹⁶ R. M. Cooke and I. D. Car

¹⁷ K. Wüthrich, "NMR of Pro

¹⁸ G. M. Clore and A. M. Gr

¹⁹ M. J. Darsley, P. de la Paz

and Exploitation of Antib

and Analysis Vol. 15 (E. R

1985).

²⁰ C. R. Mainhart, M. Potter,

²¹ M. E. Snow and L. M. Am

up to 10° in constant
4-6 is also seen in a
ion,⁷ one of the most
ch mediates assembly

nded β sheets to form
th the strands at the
e another.⁹ The loops
main pairing bringing
ity. The CDRs them-
e.¹⁰

determining regions
work form the antigen
variable in comparison
ne of the loops show a
tural conservation. In
ved in conformation.¹²
nia and co-workers^{13,14}
ay be grouped.

re determinations, se-
1 variable regions have
ructures, however, are
three or four immuno-
fewer than that being
Protein crystallography

zackerley, and F. Saul, *Proc.*

chemistry 12, 4620 (1973).

Mol. Biol. 186, 651 (1985).

1 (1988).

989).

er, and M. Karplus, *J. Biol.*

BO J. 5, 415 (1986).

th-Gill, G. Air, S. Sheriff,
inelli, P. M. Alzari, and R. J.

r, M. D. Brice, J. R. Rodgers,
112, 535 (1977).

is limited by two major factors: the time required to collect, process, and refine an X-ray data set and the intractability of certain proteins to crystallographic analysis (for example, certain proteins prove impossible to crystallize, or do not diffract to acceptable resolutions once crystallized). In addition to the technical complexity of crystallography itself, the supporting biochemistry can be difficult. The protein may be difficult to purify, available in limited quantities, or unstable.

Currently, the only experimental alternative to X-ray techniques is nuclear magnetic resonance (NMR).^{16,17} However, at the moment, the technique is applicable only to the solution of relatively small proteins (<20 kDa).¹⁸ Thus, with the experimental limitations and the extreme importance of antibodies in therapy and research, there have been many attempts to develop methods by which to model antibody combining sites and thus circumvent experimental procedures. The availability of accurate and reliable modeling procedures would also allow the prediction of the effects of site-directed mutagenesis (SDM) experiments and allow the intelligent application of SDM as well as larger modifications to the combining site (CDR replacement, introduction of catalytic activity, and metal binding sites) and, eventually, tailoring of combining sites to new antigens.

Antibody Modeling—Approaches to Date

Various workers have attempted to model antibody combining sites using, in the main, knowledge-based approaches.^{12,19-24} More recently, attempts have been made to use *ab initio* methods²⁵⁻²⁷ where the conformational space available to a loop is saturated, followed by a screening procedure. Generally this involves selecting the lowest energy conformation calculated using an empirically derived potential function.^{28,29}

Modeling of antibody combining sites may be divided into three classes:

1. modeling panels of mutant antibodies given a structure (determined either by crystallography or by molecular modeling) for the parent antibody

¹⁶ R. M. Cooke and I. D. Campbell, *BioEssays* 8, 52 (1988).

¹⁷ K. Wüthrich, "NMR of Proteins and Nucleic Acids." Wiley, New York, 1986.

¹⁸ G. M. Clore and A. M. Gronenborn, *Protein Eng.* 1, 275 (1987).

¹⁹ M. J. Darsley, P. de la Paz, D. C. Phillips, A. R. Rees, and B. J. Sutton, in "Investigation and Exploitation of Antibody Combining Sites" Methodological Surveys in Biochemistry and Analysis Vol. 15 (E. Reid, G. M. W. Cook, and D. J. Morré, eds.), Plenum, New York, 1985.

²⁰ C. R. Mainhart, M. Potter, and R. J. Feldmann, *Mol. Immunol.* 21, 469 (1984).

²¹ M. E. Snow and L. M. Amzel, *Proteins: Struct. Funct. Genet.* 1, 276 (1986).

2. modeling insertions, deletions, and CDR replacements given a structure for the parent antibody
3. the larger problem of modeling an unknown antibody structure from its amino acid sequence alone. This involves two stages—building the framework region and building the CDRs

Modeling Single-Site Mutations

This involves the replacement of one or more amino acid residues by others. In modeling such replacements, it is necessary to assess the possibility of backbone conformational changes in addition to side chain placement. The work of Sibanda and Thornton,³⁰ Thornton *et al.*,³¹ and of Greer³² indicates that changes in loops tend to be accommodated locally and one can thus be reasonably confident that the overall structure of the antibody will not be modified. Support for this assumption also comes from X-ray structures of mutant hemoglobins,³³ T4 lysozyme³⁴ and crambins,³⁵ the isomorphous crystallization of mutant proteins,³⁶ and solution data for phage λ -repressor mutants.^{37,38} If the residue replacements are conservative (especially if they are surface residues not making inter-CDR interactions), it may only be necessary to consider the placement of the side chain, but if nonhomologous changes are being made, changes to

- ²² E. A. Padlan, D. R. Davies, I. Pecht, D. Givol, and C. Wright, *Cold Spring Harbor Symp. Quant. Biol.* **41**, 627 (1976).
- ²³ S. J. Smith-Gill, C. R. Mainhart, T. B. Lavoie, R. J. Feldmann, W. Drohan, and B. R. Brooks, *J. Mol. Biol.* **194**, 713 (1987).
- ²⁴ E. A. Padlan and E. A. Kabat, *Proc. Natl. Acad. Sci.* **85**, 6885 (1988).
- ²⁵ R. M. Fine, H. Wang, P. S. Shenkin, D. L. Yarmush, and C. Levinthal, *Proteins: Struct. Funct. Genet.* **1**, 342 (1986).
- ²⁶ R. E. Brucoleri, E. Haber, and J. Novotný, *Nature (London)* **335**, 564 (1988).
- ²⁷ J. Moulton and M. N. G. James, *Proteins: Struct. Funct. Genet.* **1**, 146 (1986).
- ²⁸ B. Brooks, R. E. Brucoleri, B. D. Olafson, D. J. States, S. Swaminathan, and M. Karplus, *J. Comput. Chem.* **4**, 187 (1983).
- ²⁹ J. Åqvist, W. F. van Gunsteren, M. Leifonmark, and O. Tapia, *J. Mol. Biol.* **183**, 461 (1985).
- ³⁰ B. L. Sibanda and J. M. Thornton, *Nature (London)* **316**, 170 (1985).
- ³¹ J. M. Thornton, B. L. Sibanda, M. S. Edwards, and D. J. Barlow, *BioEssays* **8**, 63 (1988).
- ³² J. Greer, *J. Mol. Biol.* **153**, 1027 (1981).
- ³³ G. Fermi and M. Perutz, in "Atlas of Molecular Structures in Biology" (D. C. Phillips and F. M. Richards, eds.), Oxford Univ. Press (Clarendon) London and New York, 1981.
- ³⁴ M. G. Grütter, R. B. Hawkes, and B. W. Matthews, *Nature (London)* **277**, 667 (1979).
- ³⁵ W. A. Hendrickson and M. M. Teeter, *Nature (London)* **290**, 107 (1981).
- ³⁶ M. Knossow, R. S. Daniels, A. R. Douglas, J. J. Skehel, and D. C. Wiley, *Nature (London)* **311**, 678 (1984).
- ³⁷ M. H. Hecht, H. C. M. Nelson, and R. T. Sauer, *Proc. Natl. Acad. Sci.* **80**, 2676 (1983).
- ³⁸ M. A. Weiss, M. Karplus, D. J. Patel, and R. T. Sauer, *J. Biomol. Struct. Dynamics* **1**, 151 (1983).

backbone conformation
philic residue to a hydro
bury the hydrophobic gr

A number of appro
tions. One of the s
(MOP).^{12,20,23,38a,39-41} Ar
residue is constructed w
to the atoms in the par
dihedrals are inherited b
subjected to energy mi
protein.

Sidechain replacem
PLACE and REFI com
gram FRODO.⁴² REFI
side chain conformatio
"molten atom" or "loc
time such that each mo
the atom by decreasing
angles, and dihedral an
lated and applied to ea
moving one atom will a
50 cycles are required to
procedures may be imp
Waltham, MA).

However, when usin
in a local energy minim
perturbation procedure
making the appropriate
performed about the di
conformations are ident
tion.²⁸ Restraints are a
sumed to be local to the

- ^{38a} The terms MOP (maxim
dure), and CPP (coupled pe
³⁹ R. J. Feldmann, M. Potter,
⁴⁰ M. K. Swenson, A. W. Bu
Biology" (B. Pulman, ed.),
⁴¹ P. Warne, F. A. Momany,
13, 768 (1974).
⁴² T. A. Jones, this series, Vol
⁴³ J. Hermans and J. E. McQu
⁴⁴ R. E. Hubbard, *Proc. Comp
⁴⁵ H. L. Shih, J. Brady, and N*

acements given a
antibody structure
olves two stages—
CDRs

no acid residues by
to assess the possibil-
to side chain place-
ton *et al.*,³¹ and of
ommodated locally
erall structure of the
mption also comes
sozyme³⁴ and cram-
teins,³⁶ and solution
ue replacements are
t making inter-CDR
he placement of the
g made, changes to

Cold Spring Harbor Symp.

in, W. Drohan, and B. R.

1988).

Levinthal, *Proteins: Struct.*

35, 564 (1988).

146 (1986).

minathan, and M. Karplus,

pia, *J. Mol. Biol.* **183**, 461

(1985).

ow, *BioEssays* **8**, 63 (1988).

Biology" (D. C. Phillips and

and New York, 1981.

ondon) **277**, 667 (1979).

107 (1981).

. C. Wiley, *Nature (London)*

Acad. Sci. **80**, 2676 (1983).

Mol. Struct. Dynamics **1**, 151

backbone conformation are likely. For example, changing a surface hydrophilic residue to a hydrophobic one may result in a refolding of the loop to bury the hydrophobic group away from the solvent.

A number of approaches have been used to model single-site mutations. One of the simplest is the maximum overlap procedure (MOP).^{12,20,23,38a,39-41} Amino acid replacements are made such that the new residue is constructed with as many atoms as possible in identical positions to the atoms in the parent structure; that is, wherever possible, the parent dihedrals are inherited by the mutated residue. The entire structure is then subjected to energy minimization to obtain a prediction for the mutant protein.

Sidechain replacement by MOP is readily achieved using the REPLACE and REFI commands of the interactive molecular graphics program FRODO.⁴² REFI uses an optimization technique to fit a template side chain conformation to the atoms present in the parent structure. This "molten atom" or "local change" approach⁴³ moves only one atom at a time such that each movement improves the immediate environment of the atom by decreasing the differences from ideality of bond lengths, bond angles, and dihedral angles near this atom. Small displacements are calculated and applied to each atom in turn. The process is thus iterative, as moving one atom will affect the environment of another. Typically, up to 50 cycles are required to achieve convergence to less than 0.01 Å. Similar procedures may be implemented using HYDRA⁴⁴ or QUANTA (Polygen, Waltham, MA).

However, when using MOP, it is possible for the structure to be trapped in a local energy minimum resulting in an incorrect model. The minimum perturbation procedure (MPP)⁴⁵ helps to overcome this problem. After making the appropriate residue replacement(s), a conformational search is performed about the dihedral angles of the new side chain. Low-energy conformations are identified and subjected to restrained energy minimization.²⁸ Restraints are applied such that conformational changes are assumed to be local to the area of the residue replacement(s).

³⁶ The terms MOP (maximum overlap procedure), MPP (minimum perturbation procedure), and CPP (coupled perturbation procedure) were introduced by Snow and Amzel.²¹

³⁹ R. J. Feldmann, M. Potter, and C. P. J. Glaudemans, *Mol. Immunol.* **18**, 683 (1981).

⁴⁰ M. K. Swenson, A. W. Burgess, and H. A. Scheraga, in "Frontiers in Physicochemical Biology" (B. Pulman, ed.), p. 115. Academic Press, New York, 1978.

⁴¹ P. Warne, F. A. Momany, S. V. Rumball, R. W. Tuttle, and H. A. Scheraga, *Biochemistry* **13**, 768 (1974).

⁴² T. A. Jones, this series, Vol. 115, p. 157.

⁴³ J. Hermans and J. E. McQueen, *Acta Crystallogr. Sect. A* **30**, 730 (1974).

⁴⁴ R. E. Hubbard, *Proc. Comput.-Aided Mol. Des. Conf.* p. 99 (1984).

⁴⁵ H. L. Shih, J. Brady, and M. Karplus, *Proc. Natl. Acad. Sci.* **82**, 1697 (1985).

Simple approaches such as MOP and MPP can prove extremely successful, especially when the replacement is conservative and the residue makes few interactions with neighboring residues. Karplus's group⁴⁵ cites an example where they have modeled a Gly \rightarrow Asp mutation in influenza virus hemagglutinin by MPP. They accurately positioned the aspartate side chain in the region indicated by a difference density peak from crystallographic techniques.³⁶

The coupled perturbation procedure (CPP)²¹ extends the MPP approach by searching the conformational space not only of the replaced residue, but also of a "dependency set." The dependency set for each replaced residue constitutes those residues whose conformations are likely to be affected by the replacement, i.e., those which are capable of making interactions with any conformation of the parent or replacement residue. As in MPP, the structure is then subjected to restrained energy minimization. By using more comprehensive conformational search procedures such as CONGEN,⁴⁶ CPP can search side chain conformations throughout the whole structure, effectively extending the dependency set to all residues in the molecule. However, the computer time required for such a search is likely to become prohibitive.

All three methods (MOP, MPP, CPP), however, are likely to prove inadequate if a drastic change in the nature of a residue is made, since none accounts for the possibility of backbone conformational change.

Modeling Insertions, Deletions, and CDR Replacements

As with modeling single-site mutations, such changes are likely to be accommodated locally within the CDR. However, the magnitude of these local changes may be large and it is thus always necessary to consider conformational changes in the backbone of the loop. CDR replacement is thus an extension of the insertion/deletion problem as it is obviously necessary to consider the whole loop conformation. However, in contrast to modeling the entire combining site, the environment in which the loop is constructed is largely defined by the presence of the other five loops of the combining site.

To date, most modeling in this category has relied on the presence of CDRs of the required length in the database of known antibody structures. When these are not present, the closest available CDR has been modified manually using computer graphics.^{12,23,24} When loops of the correct length are available, such methods are likely to be quite successful. However, when manual insertions or deletions need to be made, the results are highly

⁴⁶ R. E. Bruccoleri and M. Karplus, *Biopolymers* 26, 137 (1987).

dependent on user intervention. It has been shown that, in a large number of searches,^{27,46} not only the vicinity of the replacement is important, but also the side and again require (see below).

Modeling the Framework

When modeling a CDR alone, it is first necessary to consider the packing of the V_L and V_H domains. The degree of variability in the framework is a function of the degree of homology to the structure of the parent antibody (see unpublished, 1991).

A single known high homology to the structure for maximum overlap with the parent antibody. Alternatively, light and heavy chain pairing structures.¹⁴ V_L and V_H sequence homology. If the framework is composed of V_L1/V_H2 and V_L1/V_H1 domains, it is necessary to leave the V_L1/V_H2 domain intact.

Problems have been encountered in modeling the heavy chain framework of the antibody HyHEL-5. The heavy chain framework is composed of the equivalent residues of the light chain. The temperature factor for the heavy chain is higher than the mean for the light chain. The identified five residues in the heavy chain are: PCA-1H, PCA-1H, PCA-1H, PCA-1H, PCA-1H. This suggests these residues are not well defined in the crystal structure. The atoms (N, C α , C, O) in the heavy chain are identified: PCA-1H, Va, Vb, Vc, Vd, Ve, Vf, Vg, Vh, Vi, Vj, Vk, Vl, Vm, Vn, Vo, Vp, Vq, Vr, Vs, Vt, Vu, Vv, Vw, Vx, Vy, Vz. The heavy chain is poorly defined in the crystal structure.

⁴⁷ P. M. Alzari, M.-B. Lascon

⁴⁸ A. C. R. Martin, D.Phil. Th

ve extremely suc-
e and the residue
lus's group⁴⁵ cites
tation in influenza
d the aspartate side
eak from crystallo-

ends the MPP ap-
nly of the replaced
endency set for each
ormations are likely
e capable of making
replacement residue.
ed energy minimiza-
al search procedures
ormations throughout
ncy set to all residues
d for such a search is

r, are likely to prove
ue is made, since none
onal change.

ents

anges are likely to be
the magnitude of these
necessary to consider
p. CDR replacement is
lem as it is obviously
n. However, in contrast
ment in which the loop
f the other five loops of

elied on the presence of
own antibody structures.
CDR has been modified
ops of the correct length
ite successful. However,
ade, the results are highly

87).

dependent on user interaction and their reproducibility is thus low. Alternatively, it has been possible to employ complete conformational searches,^{27,46} not only of the side chains, but also of the backbone in the vicinity of the replacement. Such methods are extremely computer intensive and again require user interaction (see Conformational Search Methods, below).

Modeling the Framework

When modeling a completely unknown antibody from its sequence alone, it is first necessary to build the framework. This is relatively straightforward as it is highly conserved,⁴⁷ although differences do occur in packing of the V_L and V_H domains with respect to one another. A higher degree of variability is seen around the takeoff region of CDR-H3 and an analysis of this effect is currently under way (J. Pedersen and A. R. Rees, unpublished, 1991).

A single known high-resolution antibody structure with high sequence homology to the structure being modeled may be used as a starting point for maximum overlap (MOP)-type correction of sequence differences.¹² Alternatively, light and heavy chains may be modeled from separate starting structures.¹⁴ V_L and V_H domains are selected separately on the basis of sequence homology. If the two domains are chosen from different antibodies composed of V_L1/V_H1 and V_L2/V_H2 , respectively, the composite is formed by least squares fitting V_L1 onto V_L2 before removing V_H1 and V_L2 to leave the V_L1/V_H2 composite which inherits the V_H/V_L packing of V_L1/V_H1 .

Problems have been identified with poorly defined framework regions influencing the construction of CDRs.⁴⁸ For example, the N terminus of the heavy chain packs against CDR-H3 and in the crystal structure of the antibody HyHEL-5 the first two residues of V_H are placed differently from the equivalent residues in Gloop2 and other antibody crystal structures. The temperature factors of the HyHEL-5 atoms were examined and those higher than the mean plus three standard deviations ($\bar{x} + 3\sigma$) noted. This identified five residues with temperature factors falling outside the normal distribution: PCA-1H, Val-2H, Gln-3H, Arg-40H, and Asp-102H. This suggests these residues are either extremely mobile, or difficult to place in the crystal structure. The analysis was repeated examining just main chain atoms (N, C α , C, O) and, again, the N-terminal residues of V_H were identified: PCA-1H, Val-2H, and Gln-3H. Thus, these residues seem to be poorly defined in the crystal structure, with the side chains of the first two

⁴⁷ P. M. Alzari, M.-B. Lascombe, and R. J. Poljak, *Annu. Rev. Immunol.* 6, 555 (1987).

⁴⁸ A. C. R. Martin, D.Phil. Thesis University of Oxford, Oxford (1990).

residues and the main chain of the third being particularly poor.⁴⁹ There is, however, no evidence for any other conformation being more correct.

Main-chain temperature factors should thus be examined to identify poorly defined regions of the structure and such regions should be replaced with consensus conformations from other known crystal structures with lower temperature factors in these regions. An automated framework-building procedure has now been developed (J. Pedersen and A. R. Rees, unpublished, 1991).

Modeling the CDRs of an Entire Antibody Fv

Modeling the CDRs themselves is a more difficult task. By their very nature, these loops are hypervariable, showing, in some cases, extreme variability in sequence, length, and conformation.

Most approaches to the problem so far have used the available CDRs from antibody crystal structures as starting models for those in the antibody to be modeled.^{12,13,20,23,24,39,50}

MOP-Based Methods. Procedures based on the MOP method for single-residue replacements have been used independently by at least two groups.^{12,20,23} Their approaches differ in the way in which they select loops as starting points. The groups of Feldmann and Smith-Gill^{20,23} choose a single antibody on which to model all the hypervariable loops, selecting loops from other antibodies only when the required loop length is not present in that starting antibody structure. In contrast, the group of Rees¹² has used a single, high-resolution framework chosen by sequence homology. Individual loops are selected from a database of the available antibody crystal structures in the Brookhaven Protein Databank,¹⁵ first on length and, if two or more loops of equally suitable length exist, one loop is selected on sequence homology. These loops are then attached to the framework model.

The selected loop is fitted onto the model framework using interactive computer graphics and the sequence of the loop is then corrected using the MOP protocol. In both methods, if no loop of the required length is available, the loop of closest length to the unknown is used and insertions or deletions are made using interactive molecular graphics while attempting to maintain the overall shape of the loop and intraloop hydrogen bonding.

When this procedure has been repeated for each of the six loops and the

⁴⁹ S. Sheriff, personal communication.

⁵⁰ C. Chothia, A. M. Lesk, M. Levitt, A. G. Amit, R. A. Mariuzza, S. E. V. Phillips, and R. J. Poljak, *Science* 233, 755 (1986).

side chains of the framework the whole modeled Fv is strained²³ or unrestrained.

Key Residue Method. workers^{13,14,50} also exploit information is predicted on the. These are residues which Trp, Tyr, or Phe), can for Asn, Gln, Asp, Glu, Arg, tions (e.g., Gly or Pro). As manual although it is a pr

All Protein Database. restricted size of the anti only around a dozen str known immunoglobulin tions in CDRs when cont able among the known str procedures. Such changes computer graphics.^{12,23,24} (its repeatability is low. Th critical in determining th structures become availab CDR of identical length a eled and of identifying a crease. At the present time the Brookhaven Protein I two per year.

One way of expandin loops from all known pro database to antibodies. Th protein loops was propos developed a method for COMPOSER) which utili tures. They first define a st structural comparison of : body structures). Loop re; structured. If the loop is a ti

⁵¹ A. R. Rees and P. de la Paz, 7

⁵² T. A. Jones and S. Thirup, EA

⁵³ M. J. Sutcliffe, I. Haneef, D. C

ilarly poor.⁴⁹ There is, ing more correct. examined to identify ons should be replaced crystal structures with utomated framework- lersen and A. R. Rees,

cult task. By their very n some cases, extreme

sed the available CDRs ls for those in the anti-

the MOP method r endently by at least two 1 which they select loops Smith-Gill^{20,23} choose a variable loops, selecting uired loop length is not trast, the group of Rees¹² sosen by sequence homol- of the available antibody atabank,¹⁵ first on length length exist, one loop is are then attached to the

ramework using interactive is then corrected using the of the required length is own is used and insertions ar graphics while attempt-) and intraloop hydrogen ach of the six loops and the

side chains of the framework region have been resequenced using MOP, the whole modeled Fv is subjected to energy minimization, either restrained²³ or unrestrained,^{12,51} using the GROMOS²⁹ or CHARMM²⁸ potential.

Key Residue Method. The approach devised by Chothia and co-workers^{13,14,50} also exploits the antibody structure database. Loop conformation is predicted on the basis of the presence of critical or key residues. These are residues which affect loop packing (e.g., bulky residues such as Trp, Tyr, or Phe), can form hydrogen bonds or salt bridges (e.g., Ser, Thr, Asn, Gln, Asp, Glu, Arg, or Lys), or are able to adopt unusual conformations (e.g., Gly or Pro). As used by Chothia *et al.* the method is completely manual although it is a prime candidate for automation.

All Protein Database Method. The above methods are limited by the restricted size of the antibody crystal structure database, which contains only around a dozen structures. When the database is restricted to the known immunoglobulin crystal structures, making insertions and deletions in CDRs when conformations of the required lengths are not available among the known structures is one of the most unreliable parts of the procedures. Such changes are generally made by hand, using interactive computer graphics.^{12,23,24} Clearly, such a procedure is highly subjective and its repeatability is low. The site at which insertions or deletions are made is critical in determining the resulting conformation.⁴⁸ When more crystal structures become available, the probability of the database containing a CDR of identical length and high sequence homology to that being modeled and of identifying all key residues (as defined by Chothia) will increase. At the present time, however, the number of antibody structures in the Brookhaven Protein Databank is increasing only at the rate of one or two per year.

One way of expanding the database of loop structures is to examine loops from all known protein crystal structures rather than restricting the database to antibodies. The use of nonhomologous structures in modeling protein loops was proposed by Jones and Thirup.⁵² Sutcliffe *et al.*⁵³ have developed a method for modeling loops (implemented in the program COMPOSER) which utilizes a database of all high-resolution crystal structures. They first define a structurally conserved framework region (SCR) by structural comparison of all homologous known structures (here the antibody structures). Loop regions not defined within this SCR are then constructed. If the loop is a tight turn, it is built *ab initio*. Otherwise, geometri-

⁵¹ A. R. Rees and P. de la Paz, *Trends Biochem. Sci.* **11**, 144 (1986).

⁵² T. A. Jones and S. Thirup, *EMBO J.* **5**, 819 (1986).

⁵³ M. J. Sutcliffe, I. Haneef, D. Carney, and T. L. Blundell, *Protein Eng.* **1**, 377 (1987).

are determined and constraints in the framework from the database. of structurally determined the SCR, a sliding window onto which the

ologous proteins, they element in sheet, helix, torsion angles within related. Their distributions were calculated side-chain atom positioned and these form the the parent side-chain is for the required usually, the side chain is a set of rules has been split the rules into the that for structurally and ability that side-chain mutant side chain and al occupancy of atom application this supports mutated residues, the ns and 56–84% for δ 7% (γ) and 50–90% (δ) between proteins result hly accessible residues is important in surface

proposed a method for by analysis of backbone β -sheet proteins. Since each amino acid, 1000–. However, this number off points of the framework of the framework

ing. 1, 385 (1987).
ol. 196, 175 (1987).

also applies van der Waals constraints on the possible conformations of the loops. An iterative method is applied to adjust the torsion angles to fit the end-points of the loops onto the framework and the conformation with the best fit onto the framework is selected.

Random Search Method. Levinthal's group has suggested a method which combines molecular dynamics and generation of random sets of conformations for the CDRs.²⁵ These are screened to eliminate structures with van der Waals overlap either within the loop or with the rest of the molecule. Energy minimization, or molecular dynamics followed by energy minimization, is then applied to the system. By generating a large number of random structures, they attempt to saturate conformational space, selecting possible structures for the loop to be modeled by refining a set of the randomly generated structures and selecting loops of low energy. The results obtained by applying this method to modeling the anti-phosphorylcholine antibody McPC603³⁷ are mixed. The method works well for the shortest loops studied (CDR-H1, five residues; CDR-L2, seven residues) especially in the presence of the other CDRs from the crystal structure. When constructions were made in the absence of the other CDRs the prediction for CDR-H1 was still good, but the prediction for CDR-L2 in the region where it interacts with the other CDRs was poor. For these shorter loops, they find that the length of the loop to be modeled is the primary factor in determining the structure of the loop. For the longer loops studied (CDR-L3, 9 residues; CDR-H3, 11 residues) multiple energy minima were found, some significantly lower than the energy of predicted loops showing minimal root mean square (RMS) deviation from the crystal structure. No actual RMS deviations are cited as the nature of the method means that clusters of similar possible conformations are selected. At the time of publication of these results, the computer power available to them was insufficient to tackle the two longest loops (CDR-L1 and CDR-H2, 17 and 19 residues, respectively) in McPC603.

Conformational Search Methods. An alternative method to saturate conformational space is to generate all possible loop conformations by a tree search.^{27,46} A large number of alternative loop conformations is generated by rotation about the dihedral angles, ϕ and ψ , of the peptide main chain together with side-chain χ dihedral angles. One of the major problems with this type of approach is that the size of the tree search grows exponentially as the number of degrees of freedom included in the search increases. Thus, the conformational search procedure becomes impractical for longer loops. Bruccoleri *et al.*²⁶ have employed methods such as "real

²⁷ S. Rudikoff, Y. Satow, E. A. Padlan, D. R. Davies, and M. Potter, *Mol. Immunol.* **18**, 705 (1981).

space renormalization⁵⁸⁻⁶⁰ to reduce the search time for longer loops. Here, loops are built from both points at which they leave the framework simultaneously. The bottom residues are constructed first and the conformations generated are written into a conformations file or "CG" file. Around 10 low-energy conformations for this pair of residues are read back from the CG file and used on which to build the next pair of residues. The 10 lowest energy conformations of these four residues are, in turn, used to construct the next pair of residues until the size of the loop remaining for construction is small enough for a full conformational search. Alternatively, rather than constructing one residue at a time, fragments of two or three residues may be built. Numerous construction protocols thus become possible in order to construct a long loop. In Bruccoleri and Karplus's procedure, CONGEN,⁴⁶ final "chain closure" of the loop over three residues is performed analytically using a modification⁶¹ of Gō and Scheraga's chain closure algorithm.⁶² Additionally, cycles of CHARMM energy minimization²⁸ can be incorporated within the conformational search; this frees the resulting model from being restricted to the step size (generally 30°) used for the conformational search, thus permitting structures which might otherwise be rejected on the basis of van der Waals overlap.

The ϕ and ψ angles for an amino acid are combined into a single degree of freedom and the amount of conformational space explored is restricted by looking up ϕ/ψ combinations in a tabulated form of the Ramachandran plot⁶³ and examining only those combinations with energies below a specified cutoff. Three such Ramachandran maps are used: one for glycine, one for proline and one for alanine, which is taken as being representative of all the other amino acids. Three residues of the section being constructed are generated using the analytical chain closure algorithm of Gō and Scheraga⁶² (as modified by Bruccoleri and Karplus⁶¹). Thus, for a five-residue section, two residues are constructed by full conformational search while three are constructed by chain closure. Residues constructed by full conformational search will be termed *backbone* while those constructed by chain closure will be termed *chain*.

In order to improve the efficiency of the algorithm, after each combination of *backbone* torsion angles is generated, the distance between the endpoints is checked to ensure that it is possible to span the remaining gap with the number of residues left to construct (when these are in all-trans

⁵⁸ H. A. Scheraga, *Biopolymers* 22, 1 (1983).

⁵⁹ M. R. Pincus and R. D. Klausner, *Proc. Natl. Acad. Sci.* 79, 3413 (1982).

⁶⁰ M. R. Pincus, R. D. Klausner, and H. A. Scheraga, *Proc. Natl. Acad. Sci.* 79, 5107 (1982).

⁶¹ R. E. Bruccoleri and M. Karplus, *Macromolecules* 18, 2767 (1985).

⁶² N. Gō and H. A. Scheraga, *Macromolecules* 3, 178 (1970).

⁶³ G. N. Ramachandran, C. Ramakrishnan, and V. Sasisekharan, *J. Mol. Biol.* 7, 95 (1963).

configuration with the
Only if this condition
followed by the *chain*
conformations are rejec
cutoff (typically 20 kcal/
kcal/mol for side chains

Two major problem
apart from the extreme
very good at reproducin
chain orientations²⁶—t
achieved for each loop, i
structures is generated f
the conformation obser
lowest potential energy s
of solvent accessibility a
in energy by less than 2
ing both conformations
conformation with the
However, this is frequ
generated by CONGEN.
employed, the order in wh
come. In a real modeling
put together will be relat

Results of Antibody Data

So far, all these meth
structures have become
large deviations betwe
obtained, especially in
the modeled structure²⁰
tions of between 1.1 Å (c
seen for C α s and betwe
atoms. Although it is m
quacies in the modeling
crystal packing leads to
of the antibody D1.3⁵⁰
(backbone—N, C α , C,
RMS of 2.07 Å. (All at

⁶⁴ B. K. Lee and F. M. Richards

⁶⁵ S. W. Suh, T. N. Bhat, M. Davies, *Proteins: Struct. Fu*

h time for longer loops. They leave the framework fixed first and the conformations file or "CG" file.

of residues are read back next pair of residues. The residues are, in turn, used to of the loop remaining for conformational search. Alternatively, fragments of two or more protocols thus become Bruccoleri and Karplus's the loop over three residues⁶¹ of Gō and Scheraga's CHARMM energy minimization search; this frees step size (generally 30°) of structures which might overlap.

combined into a single degree of freedom explored is restricted form of the Ramachandran plot energies below a specified: one for glycine, one being representative of all other amino acids being constructed are the algorithm of Gō and Scheraga. Thus, for a five-residue loop, a conformational search while the loop is constructed by full conformational search while those constructed by

them, after each combination, the distance between the two atoms span the remaining gap when these are in all-trans

1, 3413 (1982).

Ann. Acad. Sci. 79, 5107 (1982).
7 (1985).

Iran, *J. Mol. Biol.* 7, 95 (1963).

configuration with the bond angles stretched by 5° from the optimum). Only if this condition is met are further *backbone* residues constructed followed by the *chain* residues and, finally, the side chains. In addition, conformations are rejected with van der Waals energy above a specified cutoff (typically 20 kcal/mol for *backbone*, 100 kcal/mol for *chain*, and 20 kcal/mol for side chains).

Two major problems exist with the conformational search algorithms apart from the extreme expense in CPU time. First, although the method is very good at reproducing crystal structures, including quite precise side-chain orientations²⁶—typically, RMS deviations below 2 Å may be achieved for each loop, including side-chains, a large number of low-energy structures is generated from which only one must be selected. Generally the conformation observed in the crystal structure does not represent the lowest potential energy structure. Bruccoleri *et al.*²⁶ have suggested the use of solvent accessibility as a filter. If the bottom two energy structures differ in energy by less than 2 kcal (three times the Boltzmann factor kT , indicating both conformations are well populated at room temperature) then the conformation with the smaller solvent accessible surface⁶⁴ is chosen.²⁶ However, this is frequently still unable to select the best conformations generated by CONGEN. Second, where real space renormalization is employed, the order in which the segments are constructed affects the outcome. In a real modeling situation any such decisions about how a loop is put together will be relatively arbitrary.

Results of Antibody Database Methods

So far, all these methods have had only limited success; where crystal structures have become available after modeling has been performed, quite large deviations between predicted and observed structures have been obtained, especially in side-chain conformations.^{50,65} In a comparison of the modeled structure²⁰ of J539 with the crystal structure,⁶⁵ RMS deviations of between 1.1 Å (CDR-H1 and CDR-L2) and 4.0 Å (CDR-H3) are seen for C_αs and between 2.0 Å (CDR-L1) and 6.5 Å (CDR-H3) for all atoms. Although it is most likely that these differences result from inaccuracies in the modeling, the authors cannot rule out the possibility that crystal packing leads to distortions in loop conformation. Chothia's model of the antibody D1.3⁵⁰ is better, with RMS deviations of 0.50–0.97 Å (backbone—N, C_α, C, C_β) for five of the six loops, CDR-H1 having an RMS of 2.07 Å. (All atom RMS deviations have not been published; a

⁶⁴ B. K. Lee and F. M. Richards, *J. Mol. Biol.* 55, 379 (1971).

⁶⁵ S. W. Suh, T. N. Bhat, M. A. Navia, G. H. Cohen, D. N. Rao, S. Rudikoff, and D. R. Davies, *Proteins: Struct. Funct. Genet.* 1, 74 (1986).

higher resolution map for D1.3, the map used in the comparison is at 2.8-Å resolution,⁶⁶ is required to make a more critical evaluation of the accuracy of side-chain predictions. However, examination of Fig. 1 in Chothia *et al.*⁵⁰ suggests that the side-chain predictions in some of the loops may be poor as the C_β positions inferred from peptide plane orientations appear to be badly misplaced. However, it should be noted that, at 2.8-Å resolution, the exact orientation of peptide planes outside regions of defined secondary structure may not be well defined in the electron density map.

In general, when the key residues defined by Chothia are present and loops of identical length are available in the database, the predictions are good, at least on backbone, as indicated above. However, for loops of lengths not represented in the database or lacking the critical residues required for this type of analysis, the predictions tend to be poor. In addition, relatively unconserved residues can have a profound effect on loop packing.⁶⁷ A variable residue at the base of CDR-H2 in Gloop2 (an antibody against the loop region of lysozyme), when changed from Glu to Ser, resulted in abolition of binding. Roberts *et al.*⁶⁷ argued that, since this residue is relatively buried in the model of the antibody,¹² it is most likely that the Glu50H → Ser mutation disrupts the packing of CDRs in the ACS, thus drastically altering the topology of the combining site surface and hence abolishing antigen binding. Chothia's approach would not identify this residue as being critical in loop packing. The recent crystal structure of Gloop2^{68,69} supports this prediction, as Glu-50H is not accessible on the surface of the ACS and is involved in inter-CDR contacts.

A Combined Algorithm

Clearly, none of the methods presented above is wholly satisfactory, although all show a limited amount of success. The ideal modeling procedure may thus involve a combination of the best features of a number of these approaches. Database methods can be used to reduce the search times of conformational search procedures for long loops, thus eliminating the need for real-space renormalization.⁶⁰ A number of conformations can be extracted from a database, either by using methods such as that of Sutcliffe *et al.*⁵³ (see All Protein Database Method, above) or, alternatively, by using distance constraints within the loops themselves defined from the known antibody structures. A conformational search program such as

⁶⁶ A. G. Amit, R. A. Mariuzza, S. E. V. Phillips, and R. J. Poljak, *Science* **233**, 747 (1986).

⁶⁷ S. Roberts, J. C. Cheetham, and A. R. Rees, *Nature (London)* **328**, 731 (1987).

⁶⁸ P. D. Jeffrey, R. E. Griest, G. L. Taylor, and A. R. Rees, manuscript in preparation (1991).

⁶⁹ P. D. Jeffrey, D.Phil. Thesis University of Oxford, Oxford (1989).

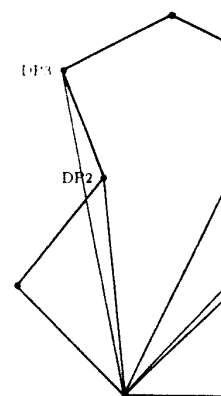


FIG. 1. C_α distances antibody and Bence-Jon then used to search the d mation.

CONGEN⁴⁶ may the best backbone confo ing only the top of searching. By remov from the knowledge renormalization from minimally dependen Such a method h

Database Searching

In order to ident seen in antibodies an antibody being mod loops are used as sh from an analysis of k CDRs, both antibody body structures were search the database is virtually 100% of a n was chosen in order distribution and to p

⁷⁰ A. C. R. Martin, J. C. C

comparison is at 2.8-Å resolution of the accuracy of the predictions. Fig. 1 in Chothia *et al.* shows that the orientations of the loops may be different at 2.8-Å resolution, but the defined secondary structure is of defined secondary structure. In the case of defined secondary structure, the predictions are generally good, but for loops of length greater than 10 residues, the predictions tend to be poor. In the case of a profound effect on the structure of DR-H2 in Gloop2 (an antibody), it is argued that, since this is a profound effect, it is most likely due to the packing of CDRs in the combining site surface. The approach would not identify the recent crystal structure of DR-H2 is not accessible on the surface contacts.

The method is wholly satisfactory, and the ideal modeling procedure features of a number of methods used to reduce the search space of conformations can be used, such as that of Chothia (1, above) or, alternatively, a search program such as

Chothia, *Science* 233, 747 (1986).
 Chothia, *Nature* 328, 731 (1987).
 Chothia, manuscript in preparation (1991).
 Chothia (1989).

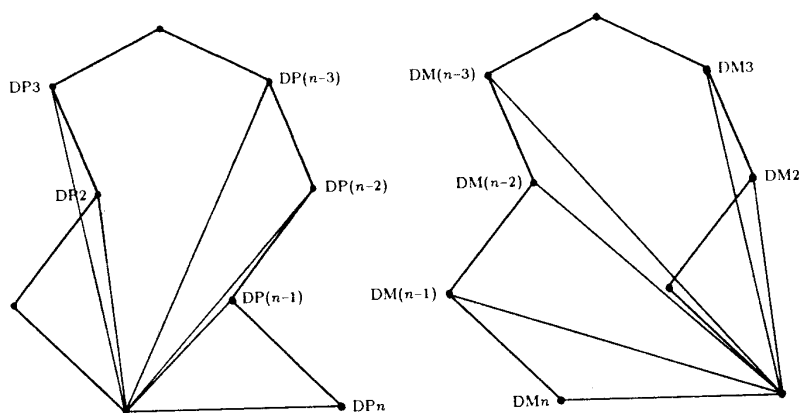


FIG. 1. C_{α} distances are measured from the N and C termini of the loops of known antibody and Bence-Jones protein structures as shown in the diagram. These distances are then used to search the database of all known protein structures for loops of similar conformation.

CONGEN⁴⁶ may then be used to select, from this set of database loops, the best backbone conformations for the bases of the loops, while reconstructing only the top of the loops and the side chains by conformational searching. By removing the requirement for manual insertions or deletions from the knowledge-based methods and the requirement for real space renormalization from the *ab initio* method, a protocol results which is minimally dependent on user choice.

Such a method has been developed^{48,70} and is described below.

Database Searching

In order to identify protein loops of the same general shape as loops seen in antibodies and, most importantly, of the required length for the antibody being modeled, sets of inter- C_{α} distance constraints within the loops are used as shown in Fig. 1. These distance constraints are derived from an analysis of known antibody crystal structures. For the light-chain CDRs, both antibody and Bence-Jones structures were used, while antibody structures were used for the heavy chain CDRs. The range used to search the database is the mean $\pm 3.5\sigma$ (σ = standard deviation). 3σ covers virtually 100% of a normal distribution; since the sample size is small, 3.5σ was chosen in order not to exclude values falling just outside the current distribution and to prevent the search from being overrestrictive.

⁷⁰ A. C. R. Martin, J. C. Cheetham, and A. R. Rees, *Proc. Natl. Acad. Sci.* 86, 9268 (1989).

A database has been created containing inter- C_α distances for every pair of residues in the range $C_{\alpha_i} \rightarrow C_{\alpha_{i+n}}$ ($-20 \leq n \leq +20$) within each protein. The database is searched by indexing a column from the database and performing a binary search on the sorted column.

Loop Processing

Having identified a set of loops matching the distance constraints applied at both the N and C termini, redundancies resulting from updated entries in the protein databank and from crystal structures containing multiple copies of a protein fragment are removed. The loops are extracted from the protein databank and are fitted onto the framework (the known crystal structure in the case of single-loop replacements, or a modeled framework where a complete antibody is being modeled). In order to orientate the loops correctly, they are overlapped onto the original loops present on the framework. Overlapping is performed in four stages:

1. The N terminus of the database loop (N_D) is moved onto the N terminus of the framework loop being replaced (N_F).
2. The C terminus of the database loop (C_D) is rotated onto the vector between N_F and the C terminus of the framework loop (C_F).
3. The database loop is translated along the vector between N_F and C_F such that the distance $(N_F - N_D) = (C_F - C_D)$.
4. The database loop is rotated about the $N_D - C_D$ vector such that the plane formed by the center of gravity of its backbone atoms, N_D and C_D , is coplanar with that formed by the center of gravity of the framework loop backbone atoms, N_F and C_F .

This procedure relies on the assumption that the loop takeoff angles defined by the planes through the backbone center of gravity, and the N and C terminal takeoff points, do not vary greatly between different antibodies with loops of different lengths. This has indeed been shown to be the case.⁴⁸

The sequences of the database loops are then corrected to match the loop being modeled by placing side chains using template conformations which are least squares fitted onto the backbone atoms (N, C_α , C, O) of the residue being replaced. Since the side chains are later repositioned using the conformational search program, CONGEN, the actual positioning of side-chain atoms with the exception of the C_β is not critical. CONGEN treats the C_β as a backbone atom, as its position is defined by the orientation of the backbone atoms, N, C_α , C, and O. Thus a side-chain conformational search which first rotates the χ_1 torsion angle ($C_\alpha - C_\beta$) is unable to position the C_β . The position of the C_β is taken from the parent amino acid except when this is a Gly in which case it is taken from the template

conformation. More sophisticated methods have been proposed^{54,55}; however, they were not considered.

As with the majority of other programs, CONGEN reduces the computational time implicitly by using "extended" coordinates for the "heavy" atom to be modeled to account for steric hindrance. Side chains potentially involve large conformational changes. Since coordinates for hydrogen atoms are not available, standard techniques, such as "polar coordinates," are used to define standard orientations.

The atom order is controlled by a list of residues which residues will be considered in the search. CONGEN conformations file.

In some cases, the number of conformations is too large to process—the conformational search. In such cases, the number of conformations is made manageable by using the algorithm of Sutcliffe⁷² which assigns a value to each conformation as being structurally favorable. The conformation is then moved relative to the position of the backbone and the structurally determined conformation being modeled using the algorithm of Sutcliffe⁷².

Conformational Search

Conformational search is performed by CONGEN,⁴⁶ which saturates the conformational space of a fragment by rotation about sidechain χ angles.

Side chains are considered in a search which starts with an energetically favorable conformation. All possible conformations are considered and the most favorable conformation with the

⁷¹ J. A. McCammon, P. G. Wo

^{71a} J. A. McCammon and S. D. Bywater, *J. Mol. Biol.*, **182**, Cambridge Univ. Press, 1985.

⁷² M. J. Sutcliffe, Ph.D. Thesis

⁷³ W. C. Barker and M. O. Dayhoff, ed.), Vol. 5. National Academy Press, 1972.

C_α distances for every pair (1-20) within each protein. C_α from the database and

the distance constraints result from updated crystal structures containing side chains. The loops are extracted from the framework (the known residues, or a modeled loop modeled). In order to add loops onto the original loops modeled in four stages:

1. N_D is moved onto the N_F plane.
2. N_D is rotated onto the vector $N_F - C_F$.
3. C_D is moved onto the vector between N_F and C_F .
4. C_D is moved onto the vector $N_F - C_F$.

at the loop takeoff angles center of gravity, and the N_D is rotated onto the vector $N_F - C_F$. It has already been shown to be the

are then corrected to match the template conformations. Atoms (N, C_α, C, O) of the side chain are later repositioned using the actual positioning of the side chain. CONGEN is defined by the orientation of the side chain. The angle (C_α - C_β) is unable to be taken from the parent amino acid taken from the template

conformation. More sophisticated side-chain replacement techniques have been proposed^{54,55}; however, since only the C_β position is important here, they were not considered worthwhile.

As with the majority of molecular mechanics-based programs, CONGEN reduces the computational burden by treating most hydrogen atoms implicitly by using "extended" (or "united") atoms⁷¹ (i.e., parameters for the "heavy" atom to which one or more hydrogens are attached are modified to account for the presence of the hydrogens). Only those hydrogens potentially involved in hydrogen bonding are treated explicitly.^{71a} Since coordinates for hydrogens are not provided by X-ray crystallographic techniques, such "polar" hydrogens are added to the heavy atoms in standard orientations.

The atom order is corrected to a standard order and, having specified which residues will be constructed by conformational searching, a CONGEN conformations file is generated for further processing by CONGEN.

In some cases, the number of loops extracted from the database is too large to process—the computer time required would become excessive. In such cases, the number of database fragments may be reduced to a more manageable number by using a structurally determining residue filter, using the algorithm of Sutcliffe.⁷² Each residue in each database fragment is assigned as being structurally determining if it causes the next C_α to be moved relative to the position of that C_α in any of the other database loops and the structurally determining residues are scored against the sequence being modeled using the Dayhoff mutation matrix.⁷³ (For details, see Sutcliffe⁷²)

Conformational Searching

Conformational searching is performed using the program CONGEN,⁴⁶ which saturates the conformational space available to a peptide fragment by rotation about the backbone ϕ and ψ torsion angles and the sidechain χ angles.

Side chains are constructed using the ITERATIVE algorithm,⁴⁶ which starts with an energetically acceptable position for all the side chains. All possible conformations for the first side chain are then regenerated and the conformation with the lowest energy is selected and its energy recorded.

⁷¹ J. A. McCammon, P. G. Wolynes, and M. Karplus, *Biochemistry* **18**, 927 (1979).

^{71a} J. A. McCammon and S. C. Harvey, in "Dynamics of Proteins and Nucleic Acids," p. 182. Cambridge Univ. Press, Cambridge, 1987.

⁷² M. J. Sutcliffe, Ph.D. Thesis University of London Birkbeck College, London (1988).

⁷³ W. C. Barker and M. O. Dayhoff, in "Atlas of Protein Sequence and Structure" (M. O. Dayhoff, ed.), Vol. 5. National Biomedical Research Foundation, Silver Spring, Maryland, 1972.

The process is repeated for each side chain in turn and, when all side chains have been searched, the process starts again from the first side chain until the energy does not change, or an iteration limit is reached. This procedure generates only one energetically feasible side-chain conformation per backbone conformation and has been shown to be the most effective⁴⁶ in generating an accurate conformation in a minimum amount of computer time.

In practice, it is rarely practical to construct more than five residues (occasionally six or seven if the conformational flexibility of the region is restricted). As discussed previously, longer regions may be constructed by "real space renormalization"⁶⁰ CONGEN treats the extracted set of database loops which have been processed into the form of a CONGEN conformations file as it would the base residues of a peptide being constructed by real space renormalization (see Conformational Search Methods, above).

A five-residue fragment in the middle of each database loop is reconstructed using a 30° search of the torsion angles of the outer two *backbone* residues and the Gō and Scheraga chain closure on the middle three *chain* residues. If this fails to find more than 100 conformations, the torsional search is reduced to 15° or, if necessary, 5°. If the search still fails, additional *backbone* residues are searched until conformations are generated.

It has been shown⁴⁸ that when building a complete combining site, the constructions are best performed in the absence of the other loops. If a loop is built between two other previously built loops, cumulative errors can cause the current loop to be built badly. In the case of short loops, however, where the FILTER algorithm (see below) cannot be used, energy is the only criterion on which the screening is performed and these loops must thus be built last in the presence of the other loops.

Energy Screening

CONGEN is built around the CHARMM²⁸ energy minimization and dynamics package. The work of Brucoleri in using CONGEN has selected a conformation from those generated by selecting the lowest energy conformation as evaluated using the CHARMM potential *in vacuo*. The *in vacuo* CHARMM potential, however, does not appear to rank the conformations very well—conformations of low RMS deviation compared with the crystal structure frequently do not fall among the lowest energy structures, as shown in Fig. 2. An examination of the low-energy conformations of CDR-H2 of HyHEL-5 using molecular graphics shows the low energy to result from an optimization of the van der Waals packing. In solvent, this attractive effect would be counterbalanced by an attraction of the loop

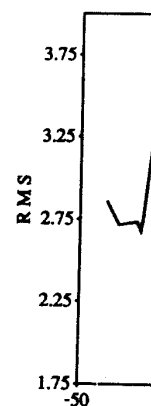


FIG. 2. Root mean square energy for each conformational state calculated using the CHARMM potential. Conformations of low energy are not the

most stable. The simulation of T1 ribonuclease collapse, while in the presence of solvent, is retained.

Use of a static representation, by chance, have not. In order to simulate the waters of hydration, a large number of water molecules are necessary when it is necessary to simulate the waters of hydration.

As an alternative, the use of a static representation of the water molecules, which occurs in the simulation of the waters of hydration, is not recommended. The Lennard-Jones potential, which is used to describe the van der Waals interactions, is not appropriate for the simulation of the waters of hydration. The distance between the water molecules and the protein atoms drops very rapidly with distance, and the simulation of the waters of hydration is not dynamically (whether or not) attempted to minimize the energy. Since the simulation is simply being used to simulate the waters of hydration, it is not necessary to minimize the energy. In addition, the use of a high dielectric constant is not necessary.

²⁴ A. D. Mackerell, Jr., L. J.

n and, when all side chains in the first side chain until is reached. This procedure -chain conformation per be the most effective⁴⁶ in num amount of computer

ct more than five residues l flexibility of the region is ons may be constructed by s the extracted set of data- the form of a CONGEN es of a peptide being con-formational Search Meth-

ch database loop is recon- s of the outer two *backbone* e on the middle three *chain* nformations, the torsional f the search still fails, addi- formations are generated. mplete combining site, the of the other loops. If a loop ops, cumulative errors can ase of short loops, however, t be used, energy is the only nd these loops must thus be

²⁸ energy minimization and using CONGEN has selected ting the lowest energy con- potential *in vacuo*. The *in* t appear to rank the confor- S deviation compared with ong the lowest energy struc- e low-energy conformations hics shows the low energy to als packing. In solvent, this y an attraction of the loop

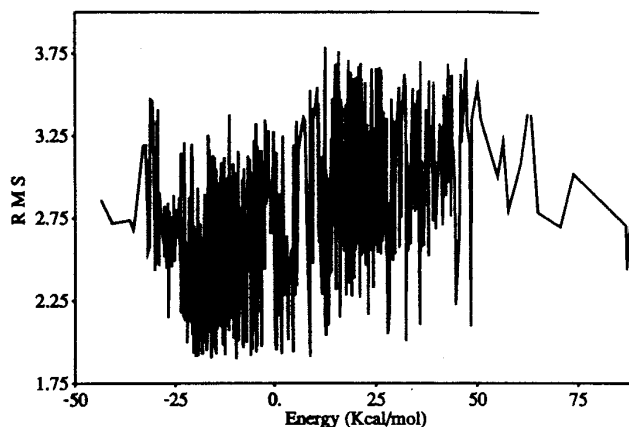


FIG. 2. Root mean square (RMS) deviation from the crystal structure is plotted against energy for each conformation generated for CDR-H2 of HyHEL-5. The energies are calculated using the CHARMM potential *in vacuo*. It is clear that those conformations which rank low in energy are not those of most similar structure to that observed in the crystal structure.

toward the solvent. This is analogous to the situation seen in a dynamics simulation of T1 ribonuclease,⁷⁴ where *in vacuo* the active site is seen to collapse, while in the presence of solvent, the native conformation is retained.

Use of a static random water box is not practical since some conformations, by chance, have favorable interactions with the solvent while others do not. In order to overcome this problem, it would be necessary to simulate the waters dynamically and computer time becomes unacceptably large when it is necessary to screen many hundreds of conformations.

As an alternative, in order to prevent the optimization of van der Waals packing which occurs *in vacuo*, the attractive ($-B/r_{ij}^6$) part of the 12-6 Lennard-Jones potential is removed. By doing this, the concept of a van der Waals radius is lost, since all atoms repel each other to some extent, whatever the distance between them [even though this repulsive effect drops very rapidly with distance (A/r_{ij}^{12})]. If the atoms were being simulated dynamically (whether by molecular dynamics, or simply by movements in attempting to minimize the potential energy), this would have to be taken into account. Since atoms are not being allowed to move, but the potential is simply being used to calculate comparative energies, this step is not necessary. In addition, the presence of solvent around the loops will lead to a high dielectric constant (around 50), reducing the importance of electro-

⁷⁴ A. D. Mackerell, Jr., L. Nilsson, and R. Rigler, *Biochemistry* 27, 4547 (1988).

statics. Thus the electrostatic potential is also removed. Similar "solvent-modified" potentials have been used in molecular dynamics of L-arabinose-binding protein,⁷⁵ in simulations of carbohydrates⁷⁶ and in analysis of misfolded proteins.⁷⁷ In the latter case, Novotný *et al.* omit the attractive part of the Lennard-Jones potential for all sulfur atoms and carbon atoms other than those in carbonyl groups, where either of the atoms in the pair being considered has a solvent accessibility greater than 0.1 Å². Since the loops of interest in antibody combining sites are on the protein surface and thus the majority of their atoms are solvent exposed, removing the attractive part of the potential for *all* atoms is effective and saves the computation time required to calculate the solvent accessibility. Novotný *et al.* calculate the solvent-modified electrostatic energy by multiplying the atomic charges by a factor dependent on the ratio of the distance of the atom from the center of the closest surface atom. The resulting "dielectric screening factor," first described by Northrup *et al.*⁷⁸ decreases the effective charge of an atom linearly from 1.0 in the center of the protein to 0.3 on the surface. For our purposes of modeling surface loops, we have found that removal of the electrostatic interaction altogether results in a greatly improved simulation and the added complexity of accounting more accurately for the center of the protein is unnecessary. This solvent-modified potential has been implemented in a modified version of GROMOS.

The rankings obtained by this solvent-modified potential are subsequently better than those from the *in vacuo* CHARMM potential, with low RMS conformations falling among the low-energy conformations (Fig. 3).

Filtering

In general, however, the lowest energy conformation is still not the lowest RMS conformation observed amongst the bottom five energies. Selection from this group is performed using one of two filtering procedures described below.

Solvent Accessibility. The solvent-modified potential used to rank the conformations considers only the enthalpic term of the free energy—entropy is ignored. In considering the local folding of small fragments of the structure of a protein, the major entropic component is that relating to the hydrophobic effect⁷⁹—the tendency for hydrophobic atoms to pack away from the solvent. By calculating the solvent-accessible surface area of

⁷⁵ B. Mao, M. R. Pear, J. A. McCammon, and F. A. Quicho, *J. Biol. Chem.* **257**, 1131 (1982).

⁷⁶ K. Bock, M. Meldal, D. R. Bundle, T. Iversen, P. J. Garegg, T. Norberg, A. A. Lindberg, and S. B. Svenson, *Carbohydr. Res.* **130**, 23 (1984).

⁷⁷ J. Novotný, A. A. Rashin, and R. E. Bruccoleri, *Proteins: Struct. Funct. Genet.* **4**, 19 (1988).

⁷⁸ S. H. Northrup, M. R. Pear, J. D. Morgan, J. A. McCammon, and M. Karplus, *J. Mol. Biol.* **153**, 1087 (1981).

⁷⁹ C. Tanford, "The Hydrophobic Effect," 2nd Ed. Wiley, New York, 1980.

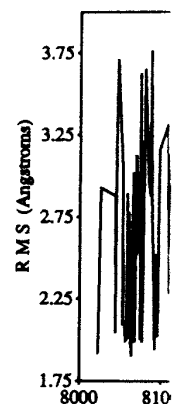


FIG. 3. RMS deviat conformation generated 1 solvent-modified GROM much improved, with two energies.

hydrophobic atoms, phobics to the solver for this purpose are number of possible c all side-chain carbon ized groups or adja

However, screeni of atoms which :

ATOMS DEFIN CONFORMATI

Ala-C_α,C_β
Cys-C_α,C_β,S_γ
Asp-C_α,C_β
Glu-C_α,C_β,C_γ
Phe-C_α,C_β,C_γ,C_{δ1},C_{δ2}
His-C_α,C_β,C_γ,C_{δ2},C_ε
Ile-C_α,C_β,C_{γ1},C_{γ2},C_δ
Lys-C_α,C_β,C_γ,C_δ
Leu-C_α,C_β,C_γ,C_{δ1},C_{δ2}

ed. Similar "solvent-dynamics of L-arabin-⁷⁶ and in analysis of l. omit the attractive ns and carbon atoms the atoms in the pair an 0.1 Å². Since the e protein surface and removing the attrac- d saves the computa- bility. Novotný *et al.* by multiplying the of the distance of the e resulting "dielectric decreases the effective of the protein to 0.3 on loops, we have found er results in a greatly accounting more accu- This solvent-modified on of GROMOS. d potential are subse- AM potential, with low conformations (Fig. 3).

mation is still not the bottom five energies. of two filtering proce-

ential used to rank the of the free energy— g of small fragments of onent is that relating to phobic atoms to pack accessible surface area of

Biol. Chem. 257, 1131 (1982).
T. Norberg, A. A. Lindberg,

uct. Funct. Genet. 4, 19 (1988).
non, and M. Karplus, *J. Mol.*

v York, 1980.

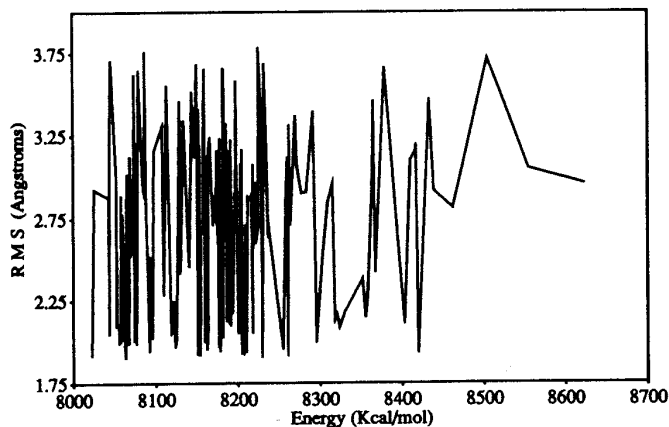


FIG. 3. RMS deviation from the crystal structure is plotted against energy for each conformation generated for CDR-H2 of HyHEL-5. The energies are calculated using the solvent-modified GROMOS potential described in the text. The ranking of conformations is much improved, with two of the lowest RMS conformations appearing in the bottom five energies.

hydrophobic atoms, that conformation with the lowest exposure of hydrophobics to the solvent can be selected. The atoms defined as hydrophobic for this purpose are listed in Table I. This list was compiled by testing a number of possible combinations to find the optimum set and represents all side-chain carbons (from C_α out), excluding those in charged, delocalized groups or adjacent to hydrophilic groups.

However, screening by solvent accessibility is oversensitive to positioning of atoms which are not part of the protein fragment being modeled.

TABLE I
ATOMS DEFINED AS HYDROPHOBIC FOR THE PURPOSE OF SCREENING
CONFORMATIONS ON THE BASIS OF THE SOLVENT ACCESSIBILITY OF
HYDROPHOBIC ATOMS

Ala-C _α C _β	Met-C _α C _β C _γ S _δ C _ε
Cys-C _α C _β S _γ	Asn-C _α C _β
Asp-C _α C _β	Pro-C _α C _β C _γ C _δ
Glu-C _α C _β C _γ	Glu-C _α C _β C _γ
Phe-C _α C _β C _γ C _{δ1} C _{δ2} C _{ε1} C _{ε2} C _ε	Arg-C _α C _β C _γ
His-C _α C _β C _γ C _{δ2} C _{ε1}	Thr-C _α C _γ
Ile-C _α C _β C _{γ1} C _{γ2} C _{δ1}	Val-C _α C _β C _γ C _{δ1} C _{δ2}
Lys-C _α C _β C _γ C _δ	Trp-C _α C _β C _γ C _{δ1} C _{δ2} C _{ε2} C _{ε3} C _{ε2} C _{ε3} C _{η2}
Leu-C _α C _β C _γ C _{δ1} C _{δ2}	Tyr-C _α C _β C _γ C _{δ1} C _{δ2} C _{ε1} C _{ε2}

Thus, while it performs well in testing the modeling procedure by reconstructing individual CDRs of a known crystal structure with the other CDRs present⁷⁰ and should perform equally well in modeling CDR replacements or mutations, it performs poorly when framework residues have also been modeled and the other CDRs are absent.

In addition, solvent accessibility cannot be used as a screen when loops are being constructed in a combining site with no prior knowledge of the structure of the other loops, since the solvent accessibility is meaningless with no other loops present. These problems led to the development of an alternative filtering method.

Structurally Determining Residues—FILTER. As an alternative to solvent accessibility, a structurally determining residue algorithm has been developed. This is applicable only to loops of six residues or longer, as it requires information from the database search (see later). For each residue of the loop in turn, the database loops are searched for the required residue type and the ϕ and ψ torsion angles are recorded. If the required residue type is not identified in any of the database loops, then the most similar residue types are identified from the Dayhoff mutation matrix and the confidence for the prediction at this residue is reduced by 5%. This procedure is repeated until residues are identified among the database loops.

The ϕ and ψ angles for each of the conformations being screened are then calculated and scored against the torsion angles observed in the database. The terminal residues start with a confidence of 50% since one of the two torsion angles is undefined. Two scoring schemes are implemented. In the first scheme, torsion angles are scored as $1/\theta$, where θ is the difference between model and database torsion angles (in radians) and θ is less than a specified cutoff. This value is then multiplied by the confidence. The second scheme simply scores a 1 for a residue if it has ϕ/ψ angles within the cutoff of any ϕ/ψ angles observed in the database conformations and thus represents the distribution of scores across the loops. Both scoring schemes are examined. If the first scoring scheme is unable to distinguish between the model conformations, the distribution of the scores is also examined—conformations which have their score distributed across the amino acids are chosen in preference to ones which score highly on only a few residues. In addition, if the conformation scoring highest in scheme 1 does not have a high distribution, this selection is rejected and a new angular cutoff is used.

The angular cutoff is initially chosen at 15°. If this is unable to select a conformation, the cutoff is increased to 30° and if this is still unable to distinguish between the models, the cutoff is further increased to 45°. This protocol is based on the ability of the different angular cutoffs to distinguish correctly between the low-energy conformations. Figure 4 shows the

ability of different angular algorithms still fails to select that all the low-energy energy conformations

Variations to the Procedure

For loops of six or more residues, the conformations constructed from the data extracted from the database backbone of a loop appear in crystal structures. Confidence is used only to select the conformation of the side chains generated by the MOP protocol.

For loops of five residues, the additional constraints means that the conformation is selected from the database is not searched and thus the database is not searched. In some cases, it has been shown

While the combination of the two methods is computer intensive. Ch

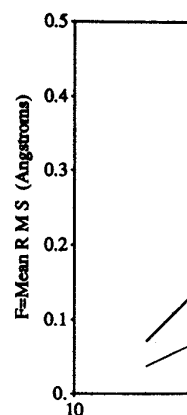


FIG. 4. Ability of the FILTER algorithm to select a conformation based on different angular cutoffs. $F = \sum(R_{sel} - R_{min})/n$, the mean RMS of the conformation selected by the FILTER algorithm and the conformation selected by the MOP protocol.

procedure by reconstruction with the other modeling CDR re-framework residues at.

a screen when loops for knowledge of the ability is meaningless development of an

As an alternative to the algorithm has been residues or longer, as it (ter). For each residue the required residue

If the required residues, then the most simulation matrix and reduced by 5%. This among the database

ns being screened are angles observed in the ce of 50% since one of 3 schemes are implemented as $1/\theta$, where θ is the es (in radians) and θ is lied by the confidence. e if it has ϕ/ψ angles atabase conformations the loops. Both scoring s unable to distinguish n of the scores is also distributed across the score highly on only a ng highest in scheme 1 is rejected and a new

this is unable to select a if this is still unable to r increased to 45°. This ngular cutoffs to distinctions. Figure 4 shows the

ability of different angular cutoffs to select the correct conformation. If the algorithm still fails to select a unique conformation, it is a fair indication that all the low-energy conformations are extremely similar and the lowest energy conformation may be selected.

Variations to the Procedure

For loops of six or seven residues, only one or two residues will be constructed from the database. In these cases, several hundred loops are extracted from the database and the conformational space available to the backbone of a loop appears to be well saturated in the current database of crystal structures. Conformational search of the backbone is thus unnecessary and is used only to construct the side chains. Conformational searching of the side chains gives better fits to the crystal structure than does use of the MOP protocol.

For loops of five residues or shorter, the limited number of distance constraints means that an enormous number of conformations ($> 10,000$) is selected from the database. Processing this number of loops becomes impractical and thus CONGEN is used alone in these cases. Because the database is not searched, the FILTER algorithm cannot be used. In these cases, it has been shown that the lowest energy structure is acceptable.

While the combined algorithm performs very well, it is extremely computer intensive. Chothia's approach of defining canonical ensembles

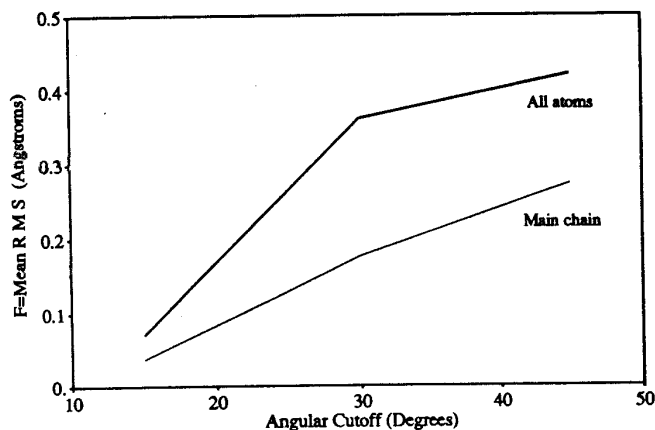


FIG. 4. Ability of the FILTER algorithm to distinguish correctly from the low-energy conformations of Gloop2 and HyHEL-5 CDRs at different angular cutoffs. For each cutoff, $F = \sum (R_{sel} - R_{min})/n$, the mean difference in RMS deviation between the selected conformation and the conformation with the minimum RMS deviation, is plotted. Cases where FILTER is unable to select a conformation are excluded, as are those where consideration of the distribution selects against the highest scoring conformation.

into which the CDRs fit requires very little computer time, but is limited to those CDRs which fit the currently defined set of canonicals. Where these structures do exist, the method works well (especially on backbone conformation) and the two approaches may successfully be combined. Thus a model may be produced by constructing those loops which fit canonicals using Chothia's approach, the remaining loops being built by the combined algorithm.

Summary

The overall protocol may be summarized as follows:

1. Construct the framework on the basis of sequence homology, rejecting sections with unusually high temperature factors and selecting them from other structures.
2. For loops of five residues or shorter, construct the loop using CONGEN. Go to step 8.
3. For loops of six residues or longer, search the distance database for backbone conformations.
4. Overlap the database loops onto the framework.
5. Correct the sequences.
6. Add explicit hydrogens.
7. For loops of eight or more residues, delete the midsection of the loop and reconstruct with CONGEN.
8. Reconstruct side chains with CONGEN.
9. Calculate the energy for each conformation with the solvent-modified potential.
10. Screen the bottom five energy conformations with the FILTER algorithm.

The procedures described for modeling by the combined algorithm (independent of both CHARMM and GROMOS) will be made available in an integrated modeling package from Oxford Molecular Limited (Oxford, England).

Gloop2—A Case Study

Gloop2⁸⁰ is a monoclonal antibody raised against the loop peptide region of lysozyme and selected to cross-react with the native protein. This antibody has been the target of numerous site-directed mutagenesis experiments^{67,81,82} (S. Roberts, *et al.* unpublished, 1991) and modeling exercises^{12,48,70} (C. Chothia *et al.*, unpublished). The crystal structure has be-

⁸⁰ M. J. Darsley and A. R. Rees, *EMBO J.* 4, 393 (1985).

⁸¹ S. Roberts, D.Phil. Thesis University of Oxford, Oxford (1988).

⁸² K. Hilyard, D.Phil. Thesis University of Oxford, Oxford (1991).

	SE
LFR1	D I Q M T Q
L1	R A S Q E I
LFR2	W L Q Q K P
L2	A A S T L D
LFR3	G V P K R F
L3	L Q Y L S Y
LFR4	F G A G T K
HFR1	Q V Q L Q Q
H1	T F G I T
HFR2	W V K Q R T
H2	E I F P G N
HFR3	K A T L T A
H3	E I R Y
HFR4	W G Q G T T

^a Light-chain CDRs are the light-chain frame those from the heavy

^b The section of this C modeling exercise.

come available,^{68,69} all be tested.

The Gloop2 structure described above and used. The Gloop2 V_L model each loop is shown loops are all built on the two short CDRs, the other four loops. The region encloses the region used for chain come from database are constructed by C

Framework

The V_L and V_H domains built separately from

⁸³ O. Epp, P. Colman, H. I *J. Biochem.* 45, 513 (197

⁸⁴ S. Sheriff, E. W. Silvert *D. R. Davies, Proc. Natl*

ter time, but is limited to canonicals. Where these ally on backbone conforly be combined. Thus a ops which fit canonicals being built by the com-

ollows:

f sequence homology, re-erature factors and select-

construct the loop using

the distance database for

etwork.

ete the midsection of the

on with the solvent-modi-

ations with the FILTER

the combined algorithm S) will be made available Molecular Limited (Ox-

against the loop peptide th the native protein. This ected mutagenesis experi-991) and modeling exer- crystal structure has be-

(1988).
1991).

TABLE II
SEQUENCE OF THE V_L REGION OF GLOOP2^a

LFR1	D I Q M T Q S P S S L S A S L G E R V S L T C
L1	R A S Q E I S G Y L S
LFR2	W L Q Q K P D G T I K R L I Y
L2	A A S T L D S
LFR3	G V P K R F S G R R S G S D Y S L T I S S L E S E D F A D Y Y C
L3	L Q Y L S Y P L T
LFR4	F G A G T K L E
HFR1	Q V Q L Q Q S G T E L A R P G A S V R L S C K A S G Y T F T
H1	T F G I T
HFR2	W V K Q R T G Q G L E W I G
H2	E I F P G N S K T Y (Y A E R F K G) ^b
HFR3	K A T L T A D K S S T T A Y M Q L S S L T S E D S A V Y F C A R
H3	E I R Y
HFR4	W G Q G T T L T V

^a Light-chain CDRs are named L1, L2, and L3; heavy-chain CDRs are H1, H2, and H3; the light-chain framework regions are named LFR1, LFR2, LFR3, and LFR4 while those from the heavy chain are named HFR1, HFR2, HFR3, and HFR4.

^b The section of this CDR enclosed in parentheses was considered as framework for the modeling exercise.

come available,^{68,69} allowing some of the predictions from the modeling to be tested.

The Gloop2 structure has been modeled using the combined algorithm described above and will be used as a case study to describe the approach used. The Gloop2 V_L sequence is shown in Table II. The protocol used to model each loop is shown in Table III and will be described in turn. The loops are all built onto an "empty" combining site with the exception of the two short CDRs, H1 and H3, which are built last in the presence of the other four loops. The notation used in Table III to describe the constructions encloses the region built by CONGEN in square brackets and the region used for chain closure in parentheses. Residues not thus enclosed come from database structures. In all cases, the side chains of all residues are constructed by CONGEN.

Framework

The V_L and V_H domains of the framework for the Gloop2 model were built separately from REI⁸³ and HyHEL-5,⁸⁴ respectively. The β strands of

⁸³ O. Epp, P. Colman, H. Fehlhammer, W. Bode, M. Schiffer, R. Huber, and W. Palm, *Eur. J. Biochem.* **45**, 513 (1974).

⁸⁴ S. Sheriff, E. W. Silvertown, E. A. Padlan, G. H. Cohen, S. J. Smith-Gill, B. C. Finzel, and D. R. Davies, *Proc. Natl. Acad. Sci.* **84**, 8075 (1987).

TABLE III
CONSTRUCTION OF GLOOP2 USING THE COMBINED ALGORITHM

Loop	Sequence	Number of conformations	CPU time (hr)	Global fit RMS Deviation			
				$P2_1$		$P1$	
				All (Å)	C_α (Å)	All (Å)	C_α (Å)
L1	RAS[Q(EIS)G]YLS	1368	59 + 25	1.90	0.92	2.09	0.86
L2	AASTLDS	157	4 + 1	1.25	0.68	1.10	0.66
L3	LQ[Y(LSY)P]LT	721	24 + 14	2.08	0.76	2.00	0.75
H1	[T(FGI)T]	375	18 + 8	2.06	0.79	2.04	1.03
H2	EI[F(PGN)S]KTY	3066	123 + 57	2.51	1.35	2.23	1.20
H3	[R(EIR)Y]	144	7 + 3	2.17	1.00	1.76	0.76
Overall			235 + 108	2.11	1.23	1.96	1.25

^a For each of the six CDRs of Gloop2, the construction protocol is indicated together with the number of conformations generated, the computer time required for the conformational search and the energy screening (on a VAX 3200), and global RMS deviations compared with the $P2_1$ and $P1$ crystal structures. These RMS deviations were calculated by fitting the structures on the framework regions and calculating the deviation over the CDRs on all atoms and C_α s. The all-atom RMS deviations drop between 0.02 and 0.55 Å while the C_α deviations drop between 0.03 and 0.44 Å when a local RMS fit is calculated.

REI were fitted to the β strands of the HyHEL-5 V_L , using the following residue ranges (numbered consecutively):

REI	HyHEL-5
34-39	33-38
43-47	42-46
84-90	83-89
100-104	98-102

The V_L domain of HyHEL-5 was then deleted to leave a hybrid of REI V_L and HyHEL-5 V_H . Temperature factors of the main-chain atoms (N, C_α , C, O) of the REI and HyHEL-5 domains were examined as described above. The N-terminal residues of HyHEL-5 V_H were identified as having high-temperature factors (greater than the mean plus three standard deviations): PCA-1H, Val-2H, and Gln-3H. These residues were replaced with the equivalent conformation from J539.⁶⁵ Finally, side-chain replacements were performed with the REPLACE and REFI functions of the interactive molecular graphics program FRODO.⁴²

CDR-L1

CDR-L1 is 11 residues long. The middle five residues were constructed using CONGEN, the remaining residues coming from database loops. The

FILTER algorithm (Temperatures at any angular conformations are very simulated).

CDR-L2

CDR-L2 is seven residues long. This length, the main-chain atoms were ranked by the available side chains were built and ranked with the FILTER algorithm was used to select conformation 9 scores (distribution = 1); similar but ranks lowest on distribution both score and distribution.

CDR-L3

CDR-L3 is nine residues long. It was constructed using CONGEN loops. Conformation 8 was the seventh position was in cis conformation.

CDR-H1

CDR-H1 is only five residues long. The number of constraints is thus too small. Thousands of conformations were generated. The loop was constructed using the FILTER algorithm. This means the FILTER algorithm was used only available criteria were used that the other loops were used. CDR-H1 (and CDR-H2) in the presence of, CDR-H3 construction was performed. As expected, a worse conformation was selected (1.42 Å for N, C_α , C, O), the presence of the other loops (C_α , C).

ALGORITHM

Global fit RMS Deviation

Å)	$P2_1$		$P1$	
	C_α (Å)	All (Å)	C_α (Å)	
0	0.92	2.09	0.86	
5	0.68	1.10	0.66	
8	0.76	2.00	0.75	
6	0.79	2.04	1.03	
1	1.35	2.23	1.20	
17	1.00	1.76	0.76	
11	1.23	1.96	1.25	

indicated together with the number of conformational search and the is compared with the $P2_1$ and $P1$ using the structures on the framework atoms and C_α s. The all-atom RMS drops drop between 0.03 and 0.44 Å

V_L , using the following

ave a hybrid of REI V_L main-chain atoms (N, C_α), examined as described were identified as having plus three standard deviations were replaced with side-chain replacements actions of the interactive

residues were constructed from database loops. The

FILTER algorithm (Table IV) is unable to distinguish between the conformations at any angular cutoff, suggesting that all the low-energy conformations are very similar. Thus the lowest energy conformation was selected.

CDR-L2

CDR-L2 is seven residues long. It has been shown⁴⁸ that, for loops of this length, the main-chain conformational space is sufficiently well saturated by the available database of protein crystal structures. Thus only the side chains were built using the conformational search; the conformations were ranked with the solvent-modified potential and the FILTER algorithm was used to select a final conformation (Table IV). At the 15° cutoff, conformation 9 scores highest, but this score is achieved on a single residue (distribution = 1); similarly at 30°, conformation 9 has the highest score, but ranks lowest on distribution; at 45°, conformation 93 ranks highest on both score and distribution and is thus the conformation selected.

CDR-L3

CDR-L3 is nine residues long. The middle five residues were constructed using CONGEN, the remaining residues coming from database loops. Conformation 82 was selected using FILTER at 45°. The proline at the seventh position within the loop is correctly predicted in the unusual cis conformation.

CDR-H1

CDR-H1 is only five residues long. The number of available distance constraints is thus too small to perform a database search (tens of thousands of conformations are extracted from the database) and the whole loop was constructed using CONGEN. The absence of the database search means the FILTER algorithm cannot be used and, since energy is thus the only available criterion on which to select a conformation, it is important that the other loops are present in order to calculate the energy. Thus CDR-H1 (and CDR-H3, which is also short) was constructed after, and in the presence of, CDR-L1, CDR-L2, CDR-L3, and CDR-H2. A repeat construction was performed in the absence of the other loops and, as expected, a worse conformation was selected (RMS: 3.21 Å for all atoms, 1.42 Å for N, C_α , C) compared with the conformation selected in the presence of the other loops (RMS: 2.18 Å for all atoms, 1.24 Å for N, C_α , C).

TABLE IV
LOW-ENERGY CONFORMATIONS FOR THE SIX CDRs OF GLOOP2 CONSTRUCTED
ONTO THE EMPTY COMBINING SITE^a

CDR	Conformation	Energy (kcal/mol)	RMS (Å)		FILTER score			Identifier
			All	B/B	15°	30°	45°	
L1	247*	5975.5	2.94	2.04	1.36(3)	2.60(6)	2.56(6)	IREI-A-0024
	319	5975.7	2.97	2.11	1.36(3)	2.22(5)	2.62(7)	IREI-A-0024
	320	5978.2	2.96	2.09	1.36(3)	2.51(6)	2.62(7)	IREI-A-0024
	241	5985.9	2.64	1.88	1.36(3)	2.60(6)	2.56(6)	IREI-A-0024
	253	5986.0	2.99	2.11	1.36(3)	2.60(6)	2.56(6)	IREI-A-0024
L2	43	5273.0	4.22	2.56	1.92(2)	2.97(2)	9.48(4)	IPYP-0-0109
	93*	5297.9	1.51	1.08	0.27(1)	9.25(5)	26.23(7)	2RHE-0-0051
	11	5300.6	1.86	1.33	3.52(5)	5.81(5)	21.58(7)	1FB4-L-0049
	109	5358.7	2.26	1.33	2.22(3)	5.08(4)	7.81(5)	2YHX-0-0383
	9	5366.7	2.06	1.33	6.64(1)	12.45(2)	26.06(6)	1CN1-A-0143
L3	98	5816.3	3.67	1.95	1.61(3)	2.82(4)	3.40(4)	IREI-A-0089
	84	5816.7	3.68	1.88	1.61(3)	2.82(4)	3.40(4)	IREI-A-0089
	107	5816.9	3.23	2.04	1.61(3)	2.82(4)	3.71(5)	IREI-A-0089
	82*	5818.4	2.82	1.99	1.61(3)	2.82(4)	3.87(6)	IREI-A-0089
	96	5818.9	3.70	1.94	1.61(3)	2.82(4)	3.40(4)	IREI-A-0089
H1	14*	6228.4	2.18	1.24	—	—	—	—
	12	6233.0	2.30	1.40	—	—	—	—
	16	6238.7	2.21	1.54	—	—	—	—
	4	6239.4	2.37	1.40	—	—	—	—
	11	6239.8	2.28	1.53	—	—	—	—
H2	240	5559.3	3.81	1.87	0.63(2)	2.28(5)	2.72(5)	1FB4-H-0050
	463*	5559.7	3.00	1.06	2.20(5)	3.42(6)	4.51(7)	1FBJ-H-0050
	1657	5566.7	2.50	1.10	0.96(3)	2.71(5)	4.10(6)	2HFL-H-0050
	471	5567.9	3.04	1.07	1.41(4)	2.84(6)	4.26(7)	1FBJ-H-0050
	722	5568.8	3.78	1.63	0.43(1)	1.47(3)	2.66(5)	1IG2-H-0050
H3	128*	7495.3	2.36	1.45	—	—	—	—
	125	7191.4	2.38	1.59	—	—	—	—
	41	7495.9	4.01	3.03	—	—	—	—
	120	7500.6	2.26	1.52	—	—	—	—
	138	7509.7	2.85	1.49	—	—	—	—

^a FILTER scores show the normal score with the distribution (i.e., the number of residues over which the score is achieved) in parentheses. RMS deviations are calculated against the P₂ crystal structure by least squares fitting over the framework and are quoted over the regions constructed. The identifier refers to the parent database loop from which the conformation was built. The first four characters are the PDB code, followed by the chain identifier (0 if no identifier) and the residue number of the first residue of the loop.

CDR-H2

CDR-H2, as defined by the EIFPGNSKTYAER definition¹³ of just four residues, is too restrictive. The latter definition does not affect antigen binding because it includes the P1 sheet and part of the P2 sheet here, and recommends a combined algorithm, consisting of a minimal residue defined by the loop. For Gloop2, CDR-H2 (EIFPGNSKTY). Confirmed at 15°. The all-atom RMS of the Phe at position 3 is 120° when compared with the P1 structure, the all-atom RMS drops to 2.

This concerted side effect of crystallization is due to the surface side chain being high. In the Gloop2 P structure for the Fv domain, residues of H2 show no Phe at position 3 and they are relatively imm

CDR-H3

CDR-H3 is only four residues long, which precludes the construction of five residues. It is possible to construct five residues, but one residue was built at the P1 position where the algorithm could not be

Summary

Each of the six CDRs is shown in Fig. 5.

GLOOP2 CONSTRUCTED

score		
	45°	Identifier
(6)	2.56(6)	1REI-A-0024
(5)	2.62(7)	1REI-A-0024
(6)	2.62(7)	1REI-A-0024
(6)	2.56(6)	1REI-A-0024
(6)	2.56(6)	1REI-A-0024
7(2)	9.48(4)	1PYP-0-0109
5(5)	26.23(7)	2RHE-0-0051
1(5)	21.58(7)	1FB4-L-0049
8(4)	7.81(5)	2YHX-0-0383
5(2)	26.06(6)	1CN1-A-0143
2(4)	3.40(4)	1REI-A-0089
2(4)	3.40(4)	1REI-A-0089
2(4)	3.71(5)	1REI-A-0089
2(4)	3.87(6)	1REI-A-0089
2(4)	3.40(4)	1REI-A-0089
—	—	—
—	—	—
—	—	—
—	—	—
—	—	—
28(5)	2.72(5)	1FB4-H-0050
42(6)	4.51(7)	1FBJ-H-0050
71(5)	4.10(6)	2HFL-H-0050
84(6)	4.26(7)	1FBJ-H-0050
17(3)	2.66(5)	1IG2-H-0050
—	—	—
—	—	—
—	—	—
—	—	—
—	—	—

, the number of residues over which tested against the $P2_1$ crystal structure over the regions constructed. The information was built. The first four (0 if no identifier) and the residue

CDR-H2

CDR-H2, as defined by Wu and Kabat,¹¹ is 17 residues long (sequence: EIFPGNSKTYAERFKG). This contrasts with the Chothia and Lesk definition¹³ of just four residues at the top of the loop (sequence: PGNS). The latter definition would ignore residues such as Glu-50H, known to affect antigen binding (see above). The Wu and Kabat definition seems excessive since it includes, on the C-terminal side, a complete strand of β sheet and part of the loop at the back of the β barrel. The definition used here, and recommended for all future modeling exercises using the combined algorithm, consists of the Kabat and Wu loop, but with the C-terminal residue defined as the strand partner of the N-terminal residue of the loop. For Gloop2, CDR-H2 is thus defined as 10 residues (sequence: EIFPGNSKTY). Conformation 463 was selected using a FILTER cutoff of 15°. The all-atom RMS deviation is poor (3.00 Å), owing to rotations of the Phe at position 3 in the loop and Tyr at position 10 by approximately 120° when compared with the $P2_1$ crystal structure. Gloop2 has been solved in two separate crystal forms, $P2_1$ and $P1$.^{68,69} When compared with the $P1$ structure, the side chains are placed almost perfectly and the all-atom RMS drops to 2.23 Å (global fit).

This concerted side-chain motion between crystal forms illustrates the effects of crystallization conditions on surface side-chain placement. Even though surface side chains may show relatively low temperature factors as a result of crystal packing interactions, their mobility in solution may be high. In the Gloop2 $P1$ structure, the mean side-chain temperature factor for the Fv domain is 13.46 ($\sigma = 8.20$) while the side chains of these two residues of H2 show mean temperature factors of 5.56 ($\sigma = 0.68$) for the Phe at position 3 and 7.10 ($\sigma = 1.73$) for the Tyr at position 10, suggesting they are relatively immobile.

CDR-H3

CDR-H3 is only four residues long. As was the case with CDR-H1, this short length precludes the use of the database. It was actually necessary to construct five residues rather than four, using CONGEN, as it was not possible to construct four residues by conformational search. Thus an extra residue was built at the N terminus of the loop. Once again, the FILTER algorithm could not be used and conformation 128 was selected on energy.

Summary

Each of the six CDRs of Gloop2 is shown with the modeled structure in Fig. 5.

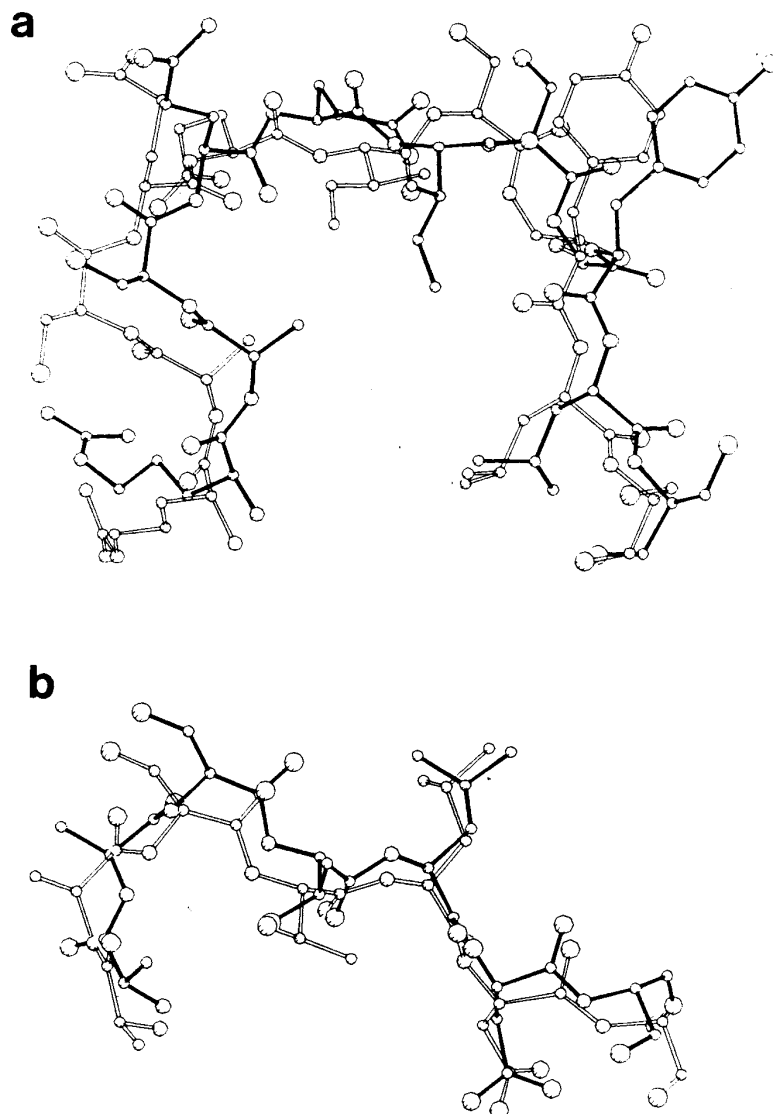
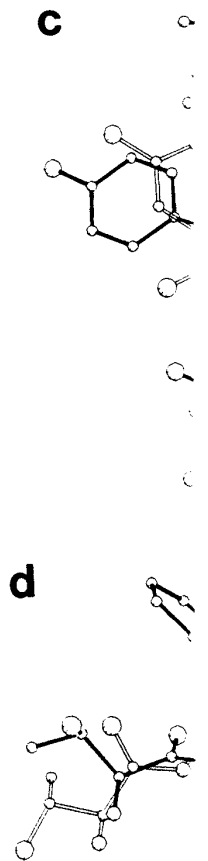
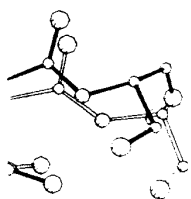
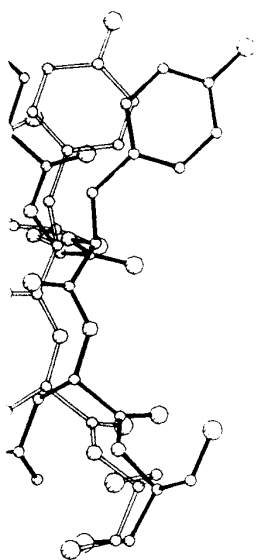


FIG. 5. The six CDRs of the Gloop2 *P2*₁ crystal structure are shown in filled lines with the modeled loops generated by the combined algorithm shown in open lines. All the fits are global (i.e., the framework regions are fitted rather than a local fit being performed on the CDRs themselves). (a) CDR-L1; (b) CDR-L2; (c) CDR-L3; (d) CDR-H1; (e) CDR-H2; (f) CDR-H3.



Overall, the results in accuracy to those as However, the canonic residues identified by structures currently a canonical ensembles. remaining CDRs do the protein engineer r the use of Chothia's n



re shown in filled lines with the
in open lines. All the fits are
cal fit being performed on the
(d) CDR-H1; (e) CDR-H2; (f)

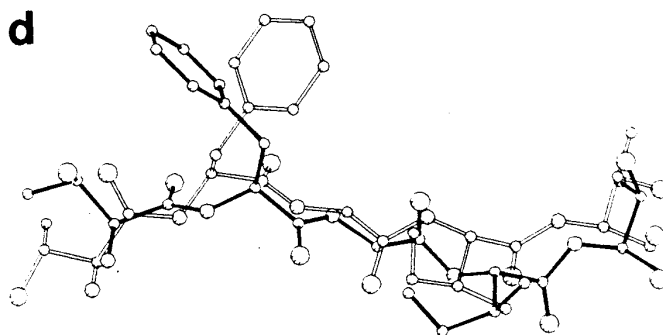
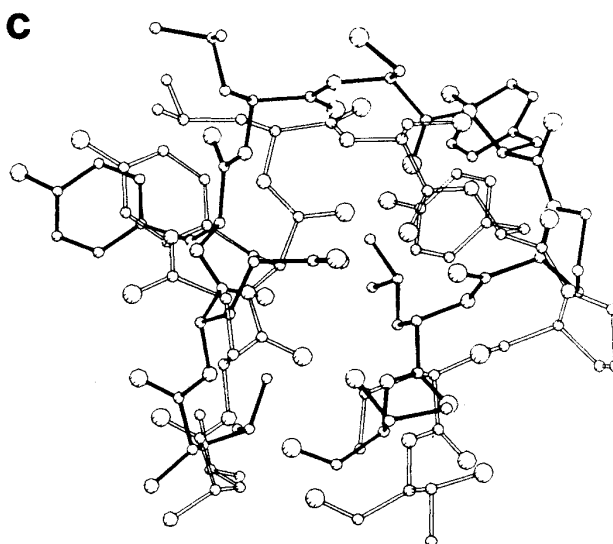


FIG. 5c, d

Overall, the results obtained using the combined algorithm are similar in accuracy to those achieved using the canonical method of Chothia *et al.* However, the canonical method is limited to those loops where the key residues identified by Chothia are present. With the number of antibody structures currently available, it is not possible to classify CDR-H3 into canonical ensembles. Additionally, a small percentage of examples in the remaining CDRs do not match the current canonical classifications and the protein engineer may well wish to mutate the key residues, precluding the use of Chothia's method for modeling the resulting conformation.

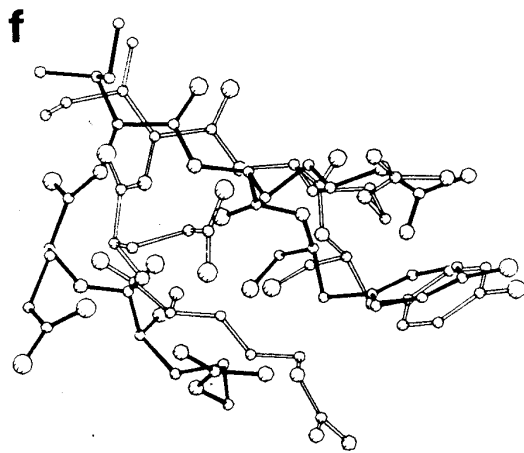
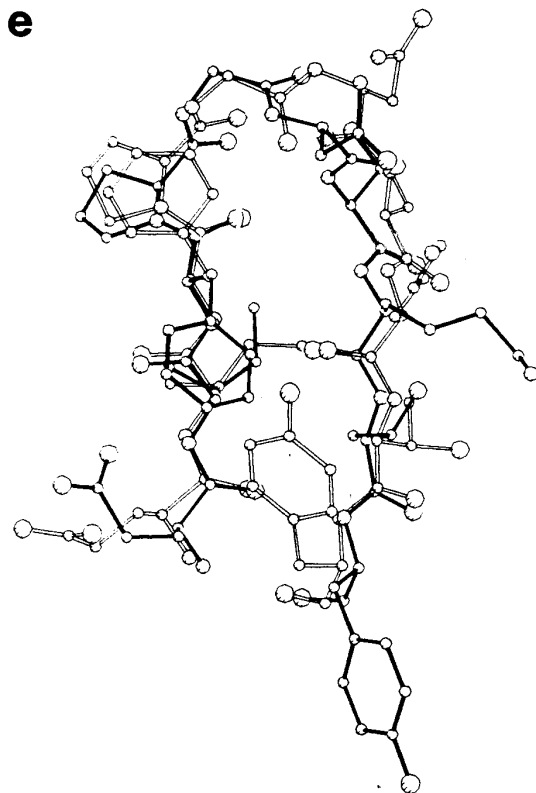


FIG. 5e, f

Thus the best approach to model the backbone represented in the data unrepresented among where mutations have eled by the combined a

Acknowledgments

The authors would like to mention the assistance of Phil J and Gloop2, and Bob Bruce SERC, UK, for financial sup

[7] X-Ray Cr Antigen-Comp Structura

By IAN A. WIL
GAIL

Introduction

The development of antibodies to produce unlimited numbers of a given antigen. This opportunity for detailed structural information and has driven the progress of antigen recognition. The structure has been solved both as free Fab determinations, combined with multiple myeloma cells, and edge of antibody-antigen interaction. The rapid increase

¹ G. Kohler and C. Milstein
² Throughout this review, the terms Fab or Fab' fragments are used to denote distinct from Fab will be t

Thus the best approach appears to be to use Chothia's method (at least to model the backbone conformation) when the loop to be modeled is represented in the database of canonical structures. Any other loops, either unrepresented among the known canonicals (including CDR-H3), or where mutations have been made to the key residues, may then be modeled by the combined algorithm presented here.

Acknowledgments

The authors would like to thank Celltech for providing a V53000 on which the computations were performed, Phil Jeffrey and Steven Sheriff for providing coordinates of HyHEL-5 and Gloop2, and Bob Bruccoleri for providing CONGEN and its source code. We also thank SERC, UK, for financial support.

[7] X-Ray Crystallographic Analysis of Free and Antigen-Complexed Fab Fragments to Investigate Structural Basis of Immune Recognition

By IAN A. WILSON, JAMES M. RINI, DAVED H. FREMONT,
GAIL G. FIESER, and ENRICO A. STURA

Introduction

The development of hybridoma technology¹ made it possible to produce unlimited numbers of antibodies of known specificity against any given antigen. This opportunity led to an upsurge of interest in obtaining detailed structural information concerning antibody-antigen recognition and has driven the present investigations of the structural basis of immune recognition. The structures of many monoclonal antibodies have now been solved both as free Fabs^{1a} and as Fab-antigen complexes. These structure determinations, combined with four crystal structures of Fabs derived from multiple myeloma antibodies, have markedly increased our knowledge of antibody-antigen interactions.

The rapid increase in the number of free and antigen-bound antibody

¹ G. Kohler and C. Milstein, *Nature (London)* **256**, 495 (1975).

^{1a} Throughout this review, the term Fab will frequently be used as the general terminology for Fab or Fab' fragments of antibodies. In specific instances, where appropriate, Fab' as distinct from Fab will be used.

