

MODELLER
A Program for Protein Structure Modeling
Release 8v1

Andrej Šali

with help from

Ben Webb, M.S. Madhusudhan, Min-Yi Shen, Marc A. Martí-Renom,
Narayanan Eswar, Frank Alber, Maya Topf, Baldomero Oliva, Andrés Fiser,
Roberto Sánchez, Bozidar Yerkovich, Azat Badretdinov,
Francisco Melo, John P. Overington, and Eric Feyfant

email sali@salilab.org

URL <http://salilab.org/modeller/>

20 June, 2005

Contents

Copyright notice	xv
Acknowledgments	xvii
1 Introduction	1
1.1 What is MODELLER?	1
1.2 MODELLER bibliography	2
1.3 Distribution	3
1.4 Installation	4
1.5 Bug reports	7
1.6 Method for comparative protein structure modeling by MODELLER	8
1.7 Tutorial on using MODELLER for comparative modeling	10
1.7.1 Preparing input files	10
1.7.2 Running MODELLER	11
1.8 Frequently asked questions (FAQ) and examples	13
1.9 MODELLER updates	21
1.9.1 Changes since release 8v0	21
1.9.2 Changes since release 7v7	21
1.9.3 Changes since release 6v2	22
1.9.4 Changes between releases 4 and 6v2	23
2 Automated comparative modeling with automodel	25
2.1 Simple usage	25
2.2 More advanced usage	26
2.2.1 Including water molecules, HETATM residues, and hydrogen atoms	26
2.2.2 Getting a very fast and approximate model	26
2.2.3 Building an all hydrogen model	27
2.2.4 Refining only part of the model	27
2.2.5 Providing your own restraints file	28
2.2.6 Using your own initial model	29
2.2.7 Adding additional restraints to the defaults	29
2.2.8 Fully automated alignment and modeling	30
2.2.9 Loop optimization	31
3 Comparative modeling class reference	35

3.1	automodel reference	35
3.1.1	automodel() — prepare to build one or more comparative models	35
3.1.2	automodel.md_level — control the model refinement level	36
3.1.3	automodel.outputs — all output data for generated models	36
3.1.4	automodel.rand_method — control initial model randomization	36
3.1.5	automodel.generate_method — control initial model generation	36
3.1.6	automodel.max_var_iterations — select length of optimizations	36
3.1.7	automodel.repeat_optimization — number of times to repeat optimization	36
3.1.8	automodel.max_molpdf — objective function cutoff	37
3.1.9	automodel.initial_malign3d — initial template alignment	37
3.1.10	automodel.starting_model — first model to build	37
3.1.11	automodel.ending_model — last model to build	37
3.1.12	automodel.final_malign3d — final template-model alignment	37
3.1.13	automodel.write_intermediates — write intermediate files during optimization	37
3.1.14	automodel.select_atoms() — select region for optimization	37
3.1.15	automodel.auto_align() — generate an automatic initial alignment	37
3.1.16	automodel.very_fast() — request rapid optimization	37
3.1.17	automodel.make() — build all models	37
3.1.18	automodel.cluster() — cluster all built models	38
3.1.19	automodel.special_restraints() — add additional restraints	38
3.1.20	automodel.special_patches() — add additional patches to the topology	38
3.1.21	automodel.user_after_single_model() — analyze or refine each model	38
3.2	allhmodel reference	38
3.2.1	allhmodel() — prepare to build all-hydrogen models	38
3.3	loopmodel reference	39
3.3.1	loopmodel() — prepare to build models with loop refinement	39
3.3.2	loopmodel.loop.md_level — control the loop model refinement level	39
3.3.3	loopmodel.loop.starting_model — first loop model to build	40
3.3.4	loopmodel.loop.ending_model — last loop model to build	40
3.3.5	loopmodel.loop.outputs — all output data for generated loop models	40
3.3.6	loopmodel.select_loop_atoms() — select region for loop optimization	40
4	MODELLER command reference	41
4.1	Miscellaneous rules and features of MODELLER	41
4.1.1	MODELLER system	41
4.1.2	Controlling breakpoints and the amount of output	41
4.1.3	File naming	41
4.1.4	File types	42
4.1.5	Key for command descriptions	43
4.2	The environ class: MODELLER environment	45
4.2.1	environ() — create a new MODELLER environment	45
4.2.2	environ.io — default input parameters	45
4.2.3	environ.edat — default objective function parameters	45
4.2.4	environ.libs — MODELLER libraries	46

4.2.5	<code>environ.dendrogram()</code> — clustering	46
4.2.6	<code>environ.principal_components()</code> — clustering	46
4.2.7	<code>environ.system()</code> — execute system command	46
4.3	The <code>energy_data</code> class: objective function parameters	47
4.3.1	<code>energy_data()</code> — create a new set of objective function parameters	47
4.3.2	<code>energy_data.contact_shell</code> — nonbond distance cutoff	47
4.3.3	<code>energy_data.update_dynamic</code> — nonbond recalculation threshold	47
4.3.4	<code>energy_data.dynamic_sphere</code> — calculate soft-sphere overlap restraints	47
4.3.5	<code>energy_data.dynamic_lennard</code> — calculate Lennard-Jones restraints	48
4.3.6	<code>energy_data.dynamic_coulomb</code> — calculate Coulomb restraints	48
4.3.7	<code>energy_data.dynamic_modeller</code> — calculate non-bonded spline restraints	48
4.3.8	<code>energy_data.excl_local</code> — exclude certain local pairs of atoms	48
4.3.9	<code>energy_data.radii_factor</code> — scale atomic radii	48
4.3.10	<code>energy_data.lennard_jones_switch</code> — Lennard-Jones switching parameters	48
4.3.11	<code>energy_data.coulomb_switch</code> — Coulomb switching parameters	48
4.3.12	<code>energy_data.relative_dielectric</code> — relative dielectric	48
4.3.13	<code>energy_data.covalent_cys</code> — use disulfide bridges in residue distance	48
4.3.14	<code>energy_data.nonbonded_sel_atoms</code> — control interaction with picked atoms	48
4.3.15	<code>energy_data.nlogn_use</code> — select non-bond list generation algorithm	49
4.4	The <code>io_data</code> class: coordinate file input parameters	50
4.4.1	<code>io_data()</code> — create a new input parameters object	50
4.4.2	<code>io_data.hetatm</code> — whether to read HETATM records	50
4.4.3	<code>io_data.hydrogen</code> — whether to read hydrogen atoms	50
4.4.4	<code>io_data.water</code> — whether to read water molecules	50
4.4.5	<code>io_data.atom_files_directory</code> — search path for coordinate files	50
4.5	The <code>libraries</code> class: stereochemical parameters and molecular topology	51
4.5.1	Modeling residues with non-existing or incomplete entries in the topology and parameter libraries	51
4.5.2	<code>libraries.topology</code> — topology library information	52
4.5.3	<code>libraries.parameters</code> — parameter library information	52
4.5.4	<code>topology.append()</code> — append residue topology library	52
4.5.5	<code>topology.clear()</code> — clear residue topology library	52
4.5.6	<code>topology.read()</code> — read residue topology library	52
4.5.7	<code>parameters.append()</code> — append parameters library	53
4.5.8	<code>parameters.clear()</code> — clear parameters library	53
4.5.9	<code>parameters.read()</code> — read parameters library	53
4.5.10	<code>topology.make()</code> — make a subset topology library	53
4.5.11	<code>topology.submodel</code> — select topology model type	54
4.5.12	<code>topology.write()</code> — write residue topology library	54
4.6	The <code>model</code> class: handling of atomic coordinates, and model building	55
4.6.1	<code>model()</code> — create a new 3D model	55
4.6.2	<code>model.seq_id</code> — sequence identity between the model and templates	55
4.6.3	<code>model.restraints</code> — all restraints which act on the model	55
4.6.4	<code>model.read()</code> — read coordinates for MODEL	55

4.6.5	<code>model.write()</code> — write MODEL	56
4.6.6	<code>model.generate_topology()</code> — generate MODEL topology	56
4.6.7	<code>model.patch()</code> — patch MODEL topology	57
4.6.8	<code>model.patch_ss_templates()</code> — guess MODEL disulfides from templates	58
4.6.9	<code>model.patch_ss()</code> — guess MODEL disulfides from model structure	59
4.6.10	<code>model.mutate()</code> — mutate selected MODEL residues	60
4.6.11	<code>model.build()</code> — build MODEL coordinates from topology	61
4.6.12	<code>model.unbuild()</code> — undefine MODEL coordinates	63
4.6.13	<code>model.transfer_xyz()</code> — copy templates' coordinates to MODEL	63
4.6.14	<code>model.res_num_from()</code> — residue numbers from MODEL2 to MODEL	64
4.6.15	<code>model.rename_segments()</code> — rename MODEL segments	65
4.6.16	<code>model.pick_atoms()</code> — select atoms in MODEL	65
4.6.17	<code>model.pick_hot_atoms()</code> — pick atoms violating restraints	68
4.6.18	<code>model.randomize_xyz()</code> — randomize MODEL coordinates	70
4.6.19	<code>model.to_iupac()</code> — standardize certain dihedral angles	70
4.6.20	<code>model.reorder_atoms()</code> — standardize order of MODEL atoms	71
4.6.21	<code>model.rotate_dihedrals()</code> — change dihedral angles	71
4.6.22	<code>model.orient()</code> — center and orient MODEL	72
4.6.23	<code>model.rotate()</code> — rotate and translate MODEL	73
4.6.24	<code>model.write_data()</code> — write derivative MODEL data	74
4.6.25	<code>model.write_pdb_xref()</code> — write residue number/index correspondence	76
4.6.26	<code>model.make_region()</code> — define a random surface patch of atoms	76
4.6.27	<code>model.color()</code> — color MODEL according to alignment	77
4.6.28	<code>model.superpose()</code> — superpose MODEL2 on MODEL given alignment	78
4.6.29	<code>model.make_chains()</code> — Fetch sequences from PDB file	81
4.6.30	<code>model.schedule.make()</code> — create optimization schedule	82
4.6.31	<code>model.schedule.step</code> — current optimization schedule step	83
4.6.32	<code>model.schedule.read()</code> — read optimization schedule	83
4.6.33	<code>model.schedule.write()</code> — write optimization schedule	83
4.6.34	<code>model.energy()</code> — evaluate MODEL given restraints	84
4.6.35	<code>model.optimize()</code> — optimize MODEL given restraints	86
4.6.36	<code>model.switch_trace()</code> — open new optimization trace file	88
4.6.37	<code>model.debug_function()</code> — test code self-consistency	89
4.6.38	<code>model.assess_ga341()</code> — assess a model with the GA341 method	90
4.6.39	<code>model.assess_dope()</code> — assess a model with the DOPE method	91
4.7	Calculation of spatial restraints	92
4.7.1	Specification of restraints	92
4.7.2	Specification of pseudo atoms	93
4.7.3	<code>model.restraints.make()</code> — make restraints	98
4.7.4	<code>model.symmetry.define()</code> — define similar segments	102
4.7.5	<code>model.restraints.unpick_all()</code> — unselect all restraints	104
4.7.6	<code>model.restraints.clear()</code> — delete all restraints	104
4.7.7	<code>model.restraints.pick()</code> — pick restraints for selected atoms	104

4.7.8	<code>model.restraints.condense()</code> — remove unselected restraints	105
4.7.9	<code>model.restraints.add()</code> — add restraint	106
4.7.10	<code>model.restraints.unpick()</code> — unselect restraint	107
4.7.11	<code>model.restraints.reindex()</code> — renumber MODEL2 restraints for MODEL	107
4.7.12	<code>model.restraints.spline()</code> — approximate restraints by splines	108
4.7.13	<code>model.restraints.append()</code> — read spatial restraints	109
4.7.14	<code>model.restraints.write()</code> — write spatial restraints	109
4.8	The <code>group_restraints</code> class: restraints on atom groups	110
4.8.1	<code>group_restraints()</code> — create a new set of group restraints	110
4.8.2	<code>group_restraints.append()</code> — read group restraint parameters	110
4.9	The <code>alignment</code> class: comparison of sequences and structures	111
4.9.1	Alignment file format	111
4.9.2	<code>alignment()</code> — create a new alignment	113
4.9.3	<code>alignment.append()</code> — read sequences and/or their alignment	114
4.9.4	<code>alignment.clear()</code> — delete all sequences from the alignment	115
4.9.5	<code>alignment.check()</code> — check alignment for modeling	115
4.9.6	<code>alignment.compare_with()</code> — compare two alignments	116
4.9.7	<code>alignment.append_model()</code> — copy model sequence and coordinates to alignment	116
4.9.8	<code>alignment.append_profile()</code> — add profile sequences to the alignment	117
4.9.9	<code>alignment.write()</code> — write sequences and/or their alignment	117
4.9.10	<code>alignment.edit()</code> — edit overhangs in alignment	118
4.9.11	<code>alignment.describe()</code> — describe proteins	119
4.9.12	<code>alignment.id_table()</code> — calculate percentage sequence identities	119
4.9.13	<code>alignment.compare_sequences()</code> — compare sequences in alignment	120
4.9.14	<code>alignment.align()</code> — align two (blocks of) sequences	121
4.9.15	<code>alignment.align2d()</code> — align sequences with structures	122
4.9.16	<code>alignment.malign()</code> — align two or more sequences	125
4.9.17	<code>alignment.consensus()</code> — consensus sequence alignment	126
4.9.18	<code>alignment.compare_structures()</code> — compare 3D structures given alignment	127
4.9.19	<code>alignment.align3d()</code> — align two structures	128
4.9.20	<code>alignment.malign3d()</code> — align two or more structures	129
4.9.21	<code>alignment.salign()</code> — align two or more sequences/structures of proteins	131
4.9.22	<code>alignment.to_profile()</code> — convert alignment to profile format	141
4.9.23	<code>alignment.expand()</code> — put all models into alignment	141
4.9.24	<code>alignment.segment_matching()</code> — align segments	142
4.10	The <code>profile</code> class: using sequence profiles	144
4.10.1	<code>profile()</code> — create a new profile	144
4.10.2	<code>profile.read()</code> — read a profile of a sequence	144
4.10.3	<code>profile.write()</code> — write a profile	145
4.10.4	<code>profile.to_alignment()</code> — profile to alignment	146
4.10.5	<code>profile.scan()</code> — Compare a target profile against a database of profiles	146
4.10.6	<code>profile.build()</code> — Build a profile for a given sequence or alignment	148
4.11	The <code>sequence_db</code> class: using sequence databases	151

4.11.1	<code>sequence_db()</code> — create a new sequence database	151
4.11.2	<code>sequence_db.read()</code> — read a database of sequences	151
4.11.3	<code>sequence_db.write()</code> — write a database of sequences	152
4.11.4	<code>sequence_db.search()</code> — search for similar sequences	152
4.11.5	<code>sequence_db.filter()</code> — cluster sequences by sequence-identity	154
4.12	The <code>density</code> class: handling electron microscopy density data	156
4.12.1	<code>density()</code> — create a new density map	156
4.12.2	<code>density.read()</code> — read an EM (electron microscopy) density map file	156
4.12.3	<code>density.grid_search()</code> — dock a structure into an EM (electron microscopy) density map	156
4.13	The <code>info</code> object: obtaining information about the MODELLER build	159
4.13.1	<code>info.version</code> — return the full MODELLER version number	159
4.13.2	<code>info.version_info</code> — return the version number, as a tuple	159
4.13.3	<code>info.build_date</code> — return the date this binary was built	159
4.13.4	<code>info.exe_type</code> — return the executable type of this binary	159
4.13.5	<code>info.debug</code> — returns this binary’s debug flag	159
4.13.6	<code>info.time_mark()</code> — print current date, time, and CPU time	159
4.14	The <code>log</code> object: controlling the amount of output	160
4.14.1	<code>log.level()</code> — Set all log output levels	160
4.14.2	<code>log.minimal()</code> — display minimal log output	160
4.14.3	<code>log.verbose()</code> — display verbose log output	160
4.14.4	<code>log.very_verbose()</code> — display verbose log output, and dynamic memory information	160
4.15	The <code>modfile</code> module: handling of files	161
4.15.1	<code>modfile.default()</code> — generate an ‘automatic’ filename	161
4.15.2	<code>modfile.delete()</code> — delete a file	161
4.15.3	<code>modfile.inquire()</code> — check if file exists	161
5	Converting TOP scripts from old MODELLER versions	163
5.1	Running old scripts unchanged	163
5.2	Converting TOP scripts to Python	163
5.2.1	TOP commands and variables	163
5.2.2	TOP models and alignments	164
5.2.3	TOP to Python correspondence	165
6	Methods	169
6.1	Dynamic programming for sequence and structure comparison and searching	169
6.1.1	Pairwise comparison	169
6.1.2	Variable gap penalty	170
6.1.3	Local <i>versus</i> global alignment	170
6.1.4	Similarity <i>versus</i> distance scores	171
6.1.5	Multiple comparisons	171
6.2	Optimization of the objective function by MODELLER	171
6.2.1	Function	171
6.2.2	Optimizers	172
6.3	Equations used in the derivation of the molecular pdf	173

6.3.1	Features and their derivatives	173
6.3.2	Restraints and their derivatives	175
6.4	Flowchart of comparative modeling by MODELLER	179
6.5	Loop modeling method	181

List of Figures

1.1	<i>Comparative protein modeling by satisfaction of spatial restraints.</i>	8
1.2	<i>Sample spatial restraint.</i>	9
1.3	<i>Optimization of the objective function.</i>	9

List of Tables

4.1	<i>List of file types.</i>	44
4.2	<i>Columns in an optimization trace file.</i>	89
4.3	<i>List of mathematical forms of restraints.</i>	95
4.4	<i>List of feature types that can be restrained.</i>	96
4.5	<i>List of “physical” restraint types.</i>	97
5.1	<i>Correspondence between TOP and Python commands.</i>	167
5.2	<i>Correspondence between TOP and Python variables.</i>	167

Copyright notice

MODELLER, a protein structure modeling program.

Copyright © 1989–2005 Andrej Šali.

This program is distributed in the hope that it will be useful, but without any warranty; without even the implied warranty of merchantability or fitness for any purpose. The entire risk as to the quality and performance of the program is with the user.

Distribution of the program is allowed only with the author's written consent.

Python 2.3 is incorporated into MODELLER. This is Copyright © 2001–2003 Python Software Foundation; All Rights Reserved.

gzip is included in Windows builds, which is

Copyright © 1999, 2001 Free Software Foundation, Inc.

Copyright © 1992, 1993 Jean-loup Gailly

gzip is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2, or (at your option) any later version.

gzip is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

Under the terms of the GPL, the source code for gzip can be provided on request by the MODELLER developers. Alternatively, it can be downloaded from <http://gnuwin32.sourceforge.net/>.

bzip2 is also included in Windows builds, which is

Copyright (C) 1996–2002 by Julian Seward.

This program is free software; you can redistribute it and/or modify it under the terms set out in the LICENSE file, which is included in the bzip2-1.0 source distribution.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the LICENSE file for more details.

The source code for bzip2 can be downloaded from <http://gnuwin32.sourceforge.net/>.

Acknowledgments

I am grateful to my PhD supervisor Professor Tom L. Blundell in whose laboratory at Birkbeck College the program was initiated.

I would also like to thank Professor Martin Karplus who allowed some of the data in the CHARMM topology and library files to be used with MODELLER.

I am in debt to the MODELLER users for their constructive criticisms and suggestions.

MODELLER was written when at

1989–1990: Department of Crystallography, Birkbeck College
University of London, Malet St, London WC1E 7HX, UK.

1990–1991: ICRF Unit of Structural Molecular Biology, Birkbeck College
Malet St, London WC1E 7HX, UK.

1991–1994: Department of Chemistry, Harvard University
12 Oxford St, Cambridge, MA 02138, USA.

1995–2003: The Rockefeller University,
1230 York Ave, New York, NY 10021, USA.

2003–: University of California, San Francisco,
1700 4th Street, San Francisco, CA 94143, USA.

Chapter 1

Introduction

1.1 What is MODELLER?

MODELLER is a computer program that models three-dimensional structures of proteins and their assemblies by satisfaction of spatial restraints.

MODELLER is most frequently used for homology or comparative protein structure modeling: The user provides an alignment of a sequence to be modeled with known related structures and MODELLER will automatically calculate a model with all non-hydrogen atoms.

More generally, the input to the program are restraints on the spatial structure of the amino acid sequence(s) and ligands to be modeled. The output is a 3D structure that satisfies these restraints as well as possible. Restraints can in principle be derived from a number of different sources. These include related protein structures (comparative modeling), NMR experiments (NMR refinement), rules of secondary structure packing (combinatorial modeling), cross-linking experiments, fluorescence spectroscopy, image reconstruction in electron microscopy, site-directed mutagenesis, intuition, residue-residue and atom-atom potentials of mean force, *etc.* The restraints can operate on distances, angles, dihedral angles, pairs of dihedral angles and some other spatial features defined by atoms or pseudo atoms. Presently, MODELLER automatically derives the restraints only from the known related structures and their alignment with the target sequence.

A 3D model is obtained by optimization of a molecular probability density function (pdf). The molecular pdf for comparative modeling is optimized with the variable target function procedure in Cartesian space that employs methods of conjugate gradients and molecular dynamics with simulated annealing.

MODELLER can also perform multiple comparison of protein sequences and/or structures, clustering of proteins, and searching of sequence databases. The program is used with a scripting language and does not include any graphics. It is written in standard FORTRAN 90 and will run on UNIX, Windows, or Mac computers.

1.2 MODELLER bibliography

In your publications using MODELLER, please quote

A. Šali and T. L. Blundell. Comparative protein modelling by satisfaction of spatial restraints. *J. Mol. Biol.* **234**, 779–815, 1993.

More information about the methods implemented in MODELLER, their use, applications, and limitations can be found in the papers listed on our web site at <http://salilab.org/publications/>. Here is a subset of these publications:

1. A. Šali and T. L. Blundell. Comparative protein modelling by satisfaction of spatial restraints. *J. Mol. Biol.* **234**, 779–815, 1993.
2. A. Šali, A. and J. P. Overington. Derivation of rules for comparative protein modeling from a database of protein structure alignments. *Protein Science* **3**, 1582–1596, 1994.
3. R. Sánchez and A. Šali. Comparative protein structure modeling: Introduction and practical examples with MODELLER. In *Protein Structure Prediction: Methods and Protocols*, D.M. Webster, editor, 97–129. Humana Press, 2000.
4. M. A. Martí-Renom, A. Stuart, A. Fiser, R. Sánchez, F. Melo and A. Šali. Comparative protein structure modeling of genes and genomes. *Ann. Rev. Biophys. Biomolec. Struct.* **29**, 291–325, 2000.
5. A. Fiser, R. K. G. Do and A. Šali. Modeling of loops in protein structures. *Protein Science* **9**, 1753–1773, 2000.
6. F. Melo, R. Sánchez, A. Sali. Statistical potentials for fold assessment. *Protein Science* **11**, 430–448, 2002.
7. M. A. Martí-Renom, B. Yerkovich, and A. Sali. Comparative protein structure prediction. John Wiley & Sons, Inc. *Current Protocols in Protein Science* 1, 2.9.1 – 2.9.22, 2002.
8. U. Pieper, N. Eswar, A. C. Stuart, V. A. Ilyin and A. Sali. MODBASE, a database of annotated comparative protein structure models. *Nucleic Acids Research* **30**, 255–259, 2002.
9. A. Fiser and A. Sali. MODELLER: generation and refinement of homology-based protein structure models. In *Methods in Enzymology*, C.W. Carter and R.M. Sweet, eds. Academic Press, San Diego, **374**, 463–493, 2003.
10. N. Eswar, B. John, N. Mirkovic, A. Fiser, V. A. Ilyin, U. Pieper, A. C. Stuart, M. A. Martí-Renom, M. S. Madhusudhan, B. Yerkovich and A. Sali. Tools for comparative protein structure modeling and analysis. *Nucleic Acids Research* **31**, 3375–3380, 2003.

1.3 Distribution

MODELLER is available free of charge to academic non-profit institutions.

First, please download the MODELLER distribution file from the MODELLER home page at <http://salilab.org/modeller/>. Next, please obtain the MODELLER key from the academic license server, also accessible from the MODELLER home page. This key is required to run MODELLER, and you should provide it to the Windows or Unix install program when prompted to do so. (If installing manually, the key needs to be assigned to the environment variable `KEY_MODELLER8v1` in your login or MODELLER startup script.) See file `INSTALLATION` for further installation instructions.

There is a MODELLER users email list. You can access it from the MODELLER home page or subscribe to it directly by sending an email message with the word “subscribe” in its body to modeller_usage-request@salilab.org.

A graphical interface to MODELLER is available as part of QUANTA, INSIGHTII, and DISCOVERY STUDIO, interactive molecular modeling programs from Accelrys Inc., San Diego, with many tools for protein modeling and structural analysis. These programs facilitate preparation of input files for MODELLER (*e.g.*, an alignment file) as well as an analysis of results (*e.g.*, an evaluation of the models). If you are interested in these programs, please contact

Mrs. Dana Haley-Vicente, Ph.D.

Accelrys Inc.

9685 Scranton Road

San Diego, CA 92121-3752, USA

tel +1-858-799-5322; fax +1-858-799-5100

email dhv@accelrys.com; URL <http://www.accelrys.com/>

1.4 Installation

The following installation instructions are from the INSTALLATION file in the root directory of the MODELLER distribution. See Section 1.3 for how to obtain MODELLER.

Please note that while MODELLER includes a Python interpreter, it does not include any of the standard Python modules. If you wish to write scripts which use standard modules, you will need to install these yourself, unless your system already has the correct version of Python installed:

For Windows Use the Windows installer for Python 2.3.5, available at <http://www.python.org/2.3.5/>.

For Mac OS X 10.3 (Panther) MODELLER will use the standard system Python (2.3); no additional configuration is necessary.

For Mac OS X 10.2 (Jaguar) MODELLER should work with the MacPython 2.3.3 distribution for OS X 10.2, available from <http://homepages.cwi.nl/~jack/macpython/download.html>.

For Linux with Python 2.3 (e.g. Fedora Core 3, Mandrake 10.1) MODELLER will use the standard system Python (2.3.4) under Fedora Core 3 or Mandrake 10.1; no additional configuration is necessary.

For Solaris MODELLER works with the Python 2.3.3 package available at <http://www.sunfreeware.com/>.

Other Linux/Unix systems You will need to build Python from the source code, available at the Python website. If you install it in a non-standard location (e.g. not `/usr/local/python2.3` or `/usr/lib/python2.3`) then you will need to set the PYTHONHOME environment variable before starting MODELLER, so that it can find the modules.

INSTALLATION

M O D E L L E R 8v1

Copyright(c) 1989-2005 Andrej Sali
All Rights Reserved

** PLATFORMS

MODELLER is written in Fortran 90 and runs on Pentium PC's (Linux and Windows XP), Apple Macintosh (OS 10.3) and workstations from Silicon Graphics (IRIX), Sun (Solaris), IBM (AIX), and DEC Alpha (OSF/1).

** INSTALLATION

See file README for information about how to get MODELLER. The source code is not generally available. Hence, most users are limited to the compiled versions of MODELLER. The program is distributed as a single install file that contains scripts, libraries, examples, documentation (in PDF and HTML formats) and executables for the supported platforms and operating systems. Please refer to the relevant section below for your platform:

** WINDOWS INSTALLATION

- 1) If using Windows NT/2000/XP, log on as a Computer Administrator. (For older systems, e.g. Windows 98, just log on as normal.)
- 2) Download the Windows installer modeller8v1.exe and save it to your Desktop.
- 3) Double-click on the modeller8v1 file to start the installer.
- 4) Tell the installer where to install Modeller, and enter your Modeller license key when prompted.
- 5) Once the install is complete, use the Modeller link from the Start Menu to start a Command Prompt from where you can run Modeller scripts. You can then delete the original installer file from your Desktop.

** MAC OS X INSTALLATION

- 1) Download the Modeller8v1.dmg file to your Desktop.
- 2) Double-click on the Modeller8v1.dmg file to open the disk image.
- 3) Double-click on the Install.command file within this image. Tell the installer where to install Modeller, and enter your Modeller license key when prompted.
- 4) Once the install is complete, you can run the Modeller script from a Terminal window. You can then drag both the 'Modeller 8v1' disk image and the Modeller8v1.dmg file to your trash.

** LINUX INSTALLATION (USING RPM)

(Note: replace 'i386' with 'ia64' or 'x86_64' if using an Itanium or AMD64/EM64T machine, respectively.)

- 1) Download the modeller-8v1-1.i386.rpm file.
- 2) Install the RPM file with the following command, replacing XXXX with your Modeller license key:

```
env KEY_MODELLER8v1=XXXX rpm -ivh modeller-8v1-1.i386.rpm
```

** GENERIC UNIX INSTALLATION

- 1) Download the modeller8v1.tar.gz file into a temporary directory on your computer.
- 2) Open a console or terminal (e.g. xterm, Konsole, GNOME terminal) and change to the directory where you downloaded the .tar.gz file. Unpack the file with the following commands:

```
gunzip modeller8v1.tar.gz
tar -xvf modeller8v1.tar
```

The result of unpacking will be the directory `./modeller8v1`, containing the following uncompressed files and directories:

<code>doc/</code>	MODELLER documentation directory
<code>examples/</code>	directory containing examples and tutorials
<code>Install</code>	installation script
<code>INSTALLATION</code>	this file
<code>README</code>	file describing distribution and registration
<code>modlib/</code>	libraries and data files for the program
<code>bin/</code>	.top script files and MODELLER executables

3) Go to the `./modeller8v1` directory and run the installation script:

```
./Install
```

Answer several questions as prompted. If you make a mistake, you can re-run the script.

For additional information visit our web site:

<http://salilab.org/modeller/>

Sincerely,

MODELLER Team
June 2005

1.5 Bug reports

Please report MODELLER bugs by e-mail to the MODELLER developers (for contact information, see <http://salilab.org/modeller/contact.html>), or, if you suspect your inputs may be faulty (and the files are not confidential) to the users' mailing list (see the same web page).

In order to be able to reproduce the bug, we will need all of your original input files (*e.g.*, script file, alignment, PDBs). In most cases, the full output demonstrating the error(s) you receive is useful too.

Please do not paste your input files directly into your email, but instead put them in a `.zip` or `.tar.gz` file. That way, we can see the exact same files you're using.

1.6 Method for comparative protein structure modeling by MODELLER

MODELLER implements an automated approach to comparative protein structure modeling by satisfaction of spatial restraints (Figure 1.1) [Šali & Blundell, 1993]. The method and its applications to biological problems are described in detail in references listed in Section 1.2. Briefly, the core modeling procedure begins with an alignment of the sequence to be modeled (target) with related known 3D structures (templates). This alignment is usually the input to the program. The output is a 3D model for the target sequence containing all mainchain and sidechain non-hydrogen atoms. Given an alignment, the model is obtained without any user intervention. First, many distance and dihedral angle restraints on the target sequence are calculated from its alignment with template 3D structures (Figure 1.2). The form of these restraints was obtained from a statistical analysis of the relationships between many pairs of homologous structures. This analysis relied on a database of 105 family alignments that included 416 proteins with known 3D structure [Šali & Overington, 1994]. By scanning the database, tables quantifying various correlations were obtained, such as the correlations between two equivalent $C_{\alpha} - C_{\alpha}$ distances, or between equivalent mainchain dihedral angles from two related proteins. These relationships were expressed as conditional probability density functions (pdf's) and can be used directly as spatial restraints. For example, probabilities for different values of the mainchain dihedral angles are calculated from the type of a residue considered, from mainchain conformation of an equivalent residue, and from sequence similarity between the two proteins. Another example is the pdf for a certain $C_{\alpha} - C_{\alpha}$ distance given equivalent distances in two related protein structures (Figure 1.2). An important feature of the method is that the spatial restraints are obtained empirically, from a database of protein structure alignments. Next, the spatial restraints and CHARMM energy terms enforcing proper stereochemistry [MacKerell *et al.*, 1998] are combined into an objective function. Finally, the model is obtained by optimizing the objective function in Cartesian space. The optimization is carried out by the use of the variable target function method [Braun & Gö, 1985] employing methods of conjugate gradients and molecular dynamics with simulated annealing (Figure 1.3). Several slightly different models can be calculated by varying the initial structure. The variability among these models can be used to estimate the errors in the corresponding regions of the fold.

There are additional specialized modeling protocols, such as that for the modeling of loops (Section 2.2.9).

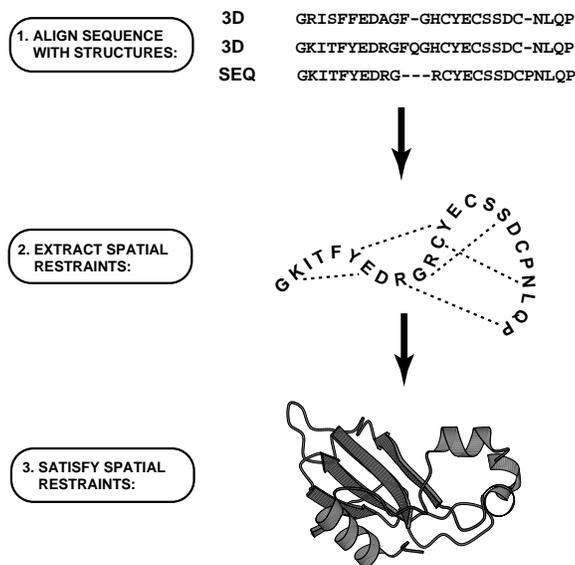


Figure 1.1: *Comparative protein modeling by satisfaction of spatial restraints.* First, the known, template 3D structures ('3D') are aligned with the target sequence to be modeled ('SEQ') Second, spatial features, such as $C_{\alpha} - C_{\alpha}$ distances, hydrogen bonds, and mainchain and sidechain dihedral angles, are transferred from the templates to the target. Thus, a number of spatial restraints on its structure are obtained. Third, the 3D model is obtained by satisfying all the restraints as well as possible.

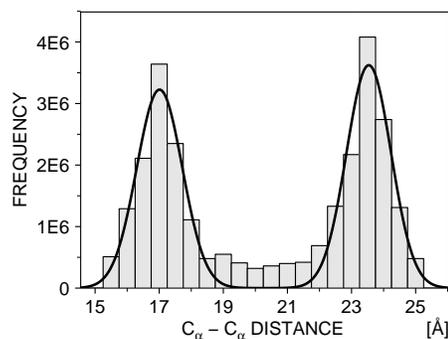


Figure 1.2: *Sample spatial restraint.* A restraint on a given C_{α} - C_{α} distance, d , is expressed as a conditional probability density function that depends on two other equivalent distances ($d' = 17.0$ and $d'' = 23.5$): $p(d/d', d'')$. The restraint (continuous line) is obtained by least-squares fitting a sum of two Gaussian functions to the histogram, which in turn is derived from the database of alignments of protein structures. In practice, more complicated restraints are used that depend on additional information, such as similarity between the proteins, solvent accessibility, and distance from a gap in the alignment [Šali & Blundell, 1993].

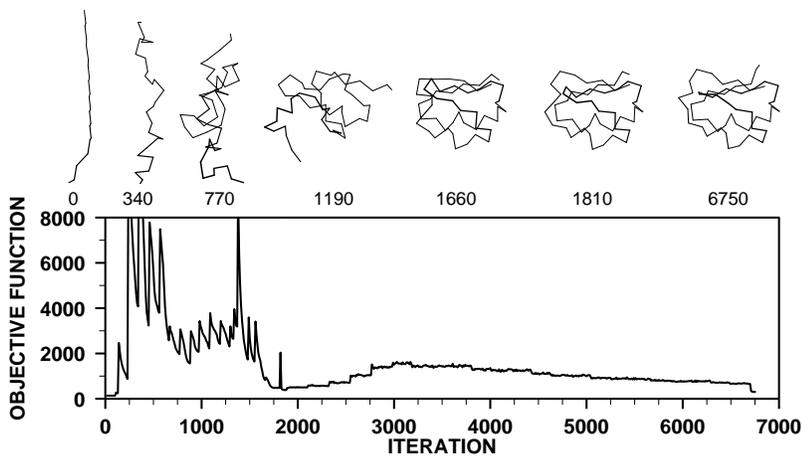


Figure 1.3: *Optimization of the objective function.* Optimization of the objective function (curve) starts with a distorted average of template structures (not with an extended structure as shown here). The iteration number is indicated below each sample structure. In this run, the first $\sim 2,000$ iterations correspond to the variable target function method relying on the conjugate gradients technique. This approach first satisfies sequentially local restraints and slowly introduces longer range restraints until the complete objective function is optimized. In the last 4,750 iterations for this model, molecular dynamics with simulated annealing is used to refine the model. Typically, a model is calculated in the order of minutes on a PC workstation.

1.7 Tutorial on using MODELLER for comparative modeling

Simple demonstrations of MODELLER in all steps of comparative protein structure modeling, including fold assignment, sequence-structure alignment, model building, and model assessment, can be found in references listed at <http://salilab.org/modeller/documentation.html>. A number of additional tools useful in comparative modeling are listed at http://salilab.org/bioinformatics_resources.shtml. Specifically, users have access to MODBASE, a comprehensive database of comparative models for all known protein sequences detectably related to at least one known protein structure; MODWEB, a web server for automated comparative protein structure modeling; and MODLOOP, a web server for automated modeling of loops in protein structures. For “frequently-asked-questions” (FAQ), see Section 1.8.

The rest of this section is a ‘hands on’ description of the most basic use of MODELLER in comparative modeling, in which the input are Protein Data Bank (PDB) atom files of known protein structures, their alignment with the target sequence to be modeled, and the output is a model for the target that includes all non-hydrogen atoms. Although MODELLER can find template structures as well as calculate sequence and structure alignments, it is better in the difficult cases to identify the templates and prepare the alignment carefully by other means. The alignment can also contain very short segments such as loops, secondary structure motifs, *etc.*

This tutorial assumes that MODELLER is already installed on your computer. See Section 1.4 for more details on installation (also in the INSTALLATION file in the MODELLER distribution directory).

1.7.1 Preparing input files

The sample input files in this tutorial can be found in the `examples/automodel` directory of the MODELLER distribution.

There are three kinds of input files: Protein Data Bank atom files with coordinates for the template structures, the alignment file with the alignment of the template structures with the target sequence, and MODELLER commands in a script file that instruct MODELLER what to do.

Atom files

Each atom file is named `code.atm` where `code` is a short protein code, preferably the PDB code; for example, *Peptococcus aerogenes* ferredoxin would be in a file `1fdx.atm`. If you wish, you can also use file extensions `.pdb` and `.ent` instead of `.atm`. The code must be used as that protein’s identifier throughout the modeling. The atom sets do not have to be superposed by the user before comparative modeling is done.

Alignment file

One of the formats for the alignment file is related to the PIR database format; this is the preferred format for comparative modeling:

C; A sample alignment in the PIR format; used in tutorial

```
>P1;5fd1
structureX:5fd1:1      : :106  : :ferredoxin:Azotobacter vinelandii: 1.90: 0.19
AFVVTDNCKIKKYTDCVEVCPVDCFYEGPNFLVIHPDECIDCALCEPECPAQAI FSEDEVPEDMQEFIQ LNAELA
EVWPNITEKKDPLPDAEDWDGKGLQH LER*

>P1;1fdx
sequence:1fdx:1      : :54   : :ferredoxin:Peptococcus aerogenes: 2.00:-1.00
AYVINDSC--IACGACKPECPVNI IQGS--IYAIDADSCIDCGSCASVCPVGAPNPED-----
-----*
```

See Section 4.9.1 for a detailed description of the alignment file format. Influence of the alignment on the quality of the model cannot be overemphasized. To obtain the best possible model, it is important to understand how the alignment is used by MODELLER [Šali & Blundell, 1993]. In outline, for the aligned regions, MODELLER tries to derive a 3D model for the target sequence that is as close to one or the other of the template structures as possible while also satisfying stereochemical restraints (*e.g.*, bond lengths, angles, non-bonded atom contacts, ...);

the inserted regions, which do not have any equivalent segments in any of the templates, are modeled in the context of the whole molecule, but using their sequence alone. This way of deriving a model means that whenever a user aligns a target residue with a template residue, he tells MODELLER to treat the aligned residues as **structurally** equivalent. Command `alignment.check()` can be used to find some trivial alignment mistakes.

Script file

MODELLER is a command-line only tool, and has no graphical user interface; instead, you must provide it with a script file containing MODELLER commands. This is an ordinary Python 2.3 script, although it must be run by MODELLER rather than the Python interpreter (since additional MODELLER-specific commands are provided).

If you are not familiar with Python, you can simply adapt one of the many examples in the `examples` directory, or look at the code for the classes used by MODELLER itself, in the `modlib/modeller` directory. Finally, there are many resources for learning Python itself, such as a comprehensive tutorial at <http://www.python.org/doc/2.3.5/tut/>.

A sample script file `model-default.py` to produce one model of sequence `1fdx` from the known structure of `5fd1` and from the alignment between the two sequences is

```
# Homology modelling by the automodel class

from modeller.automodel import *      # Load the automodel class

log.verbose()      # request verbose output
env = environ()   # create a new MODELLER environment to build this model in

# directories for input atom files
env.io.atom_files_directory = './../atom_files'

a = automodel(env,
              alnfile = 'alignment.ali',      # alignment filename
              knowns   = '5fd1',             # codes of the templates
              sequence = '1fdx')             # code of the target
a.starting_model= 1      # index of the first model
a.ending_model  = 1      # index of the last model
                  # (determines how many models to calculate)
a.make()              # do the actual homology modelling
```

See Chapter 2 for more information about the `automodel` class, and a more detailed explanation of what this script does.

1.7.2 Running MODELLER

To run MODELLER with the script file `model-default.py` above, do the following:

1. Open a command line prompt:
 - **On Linux/Unix:** `ssh` to the machine, or open an `xterm` or `GNOME Terminal`.
 - **On Windows:** Click on the 'Modeller' link on your Start Menu. This will give you a Windows Command Prompt, set up for you to run MODELLER.
 - **On Mac OS X:** Open the Terminal application.
2. Change to the directory containing the script and alignment files using the `'cd'` command.
3. Run MODELLER itself by typing the following at the command prompt:


```
mod8v1 model-default.py
```

A number of intermediary files are created as the program proceeds. After about 30 seconds on a Pentium IV workstation, the final `1fdx` model is written to file `1fdx.B99990001.pdb`. Examine the `model-default.log` file for information about the run. In particular, one should always check the output of the `alignment.check()`

command, which you can find by searching for '`check_a`'. Also, check for warning and error messages by searching for '`W>`' and '`E>`', respectively. There should be no error messages; most often, there are some warning messages that can usually be ignored.

1.8 Frequently asked questions (FAQ) and examples

Please also check the archive of the Users Mail List at http://salilab.org/archives/modeller_usage/.

1. **I want to build a model of a chimeric protein based on two known structures. Alternatively, I want to build a multi-domain protein model using templates corresponding only to the individual domains.**

This can be accomplished using the standard modeling routine. The alignment should be as follows when the chimera is a combination of proteins A and B:

```
proteinA  aaaaaaaaaaaaaaaaaaaaaaaaaa-----
proteinB  -----bbbbbbbbbbbbbbbbbbbbbbbbbbbbbb
chimera   aaaaaaaaaaaaaaaaaaaaaaabbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbb
```

In the PIR format the alignment file is:

```
>P1;proteinA
structureX:proteinA
aaaaaaaaaaaaaaaaaaaaaaaaa-----*
>P1;proteinB
structureX:proteinB
-----bbbbbbbbbbbbbbbbbbbbbbbbbbbbbb*
>P1;chimera
sequence:chimera
aaaaaaaaaaaaaaaaaaaaaabbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbb*
```

If no additional information is available about the relative orientation of the two domains the resulting model will probably have an incorrect relative orientation of the two domains when the overlap between A and B is non-existing or short. To obtain satisfactory relative orientation of modeled domains in such cases, orient the two template structures appropriately before the modeling.

2. **I don't want to use one region of a template for construction of my model.**

The easiest way to achieve this is to not align that region of the template with the target sequence. If region 'bbbbbbb' of the template should not be used as a template for region 'eeee' of the target sequence the alignment should be like this:

```
template  aaaaaaaaaaaaaaaaaaaaaa-----bbbbbbcccccccccccccccccccccccccc
target    ddddddddddddddddddeeeee-----ffffffffffffffffffffffffffff
```

The effect of this alignment is that no homology-derived restraints will be produced for region 'eeee'.

3. **I want to define (additional) disulfide bonds in the target sequence because no equivalent disulfide bonds exist in any of the templates (in which case `model.patch_ss_templates()` cannot define them automatically).**

MODELLER can restrain disulfides in two ways: automatically (`model.patch_ss_templates()` or `model.patch_ss()`) and manually (`model.patch()`).

If there is an equivalent disulfide bridge in any of the templates aligned with the target, the `model.patch_ss_templates()` command will generate appropriate disulfide bond restraints without any other input from the user. This command is run automatically by the 'model' script used for comparative modeling. The restraints include bond, angle and dihedral angle restraints. The SG — SG atom pair also becomes an excluded atom pair that is not checked for an atom-atom overlap. The χ_i dihedral angle restraints will depend on the conformation of the equivalent disulfides in the template structure, as described in [Šali & Overington, 1994]. The command `model.patch_ss()` is similar, except that the current structure of MODEL, not templates, is used to guess the disulfide bonded CYS – CYS pairs.

Explicit manual restraints can be added by the `model.patch()` command relying on the PRES DISU patching residue in the CHARMM topology file. When building comparative models with the `automodel` class, the ‘manual’ disulfides should be defined in the `automodel.special_patches()` routine. The `model.patch()` command will establish the correct stereochemistry by relying on the CHARMM topology file and parameters to restrain the disulfide bond.

It is better to use `model.patch_ss_templates()` than `model.patch()` where possible because the dihedral angles are restrained more precisely by using the templates than the general rules of stereochemistry.

Some CHARMM parameter files have a multiple dihedral entry for the disulfide dihedral angle χ_3 that consists of three individual entries with periodicities of 1, 2 and 3. This is why you see three feature restraints for a single disulfide in the output of the `model.energy()` command.

```
from modeller.automodel import *

# Redefine the special_patches routine to include the additional disulfides
# (this routine is empty by default):
class mymodel(automodel):
    def special_patches(self):
        # A disulfide between residues 1 and 85 in chain A:
        self.patch(residue_type='DISU', residue_ids=('1:A', '85:A'))
        # A disulfide between residues 41 and 45 in chain B:
        self.patch(residue_type='DISU', residue_ids=('41:B', '45:B'))

# This is as usual:
log.verbose()
env = environ()

a = mymodel(env, alnfile='align1.ali', knowns='templ1', sequence='targ1')
a.make()
```

4. I want to explicitly force certain Pro residues to the *cis* ω conformation.

MODELLER should usually be allowed to handle this automatically *via* the omega dihedral angle restraints, which are calculated by default.

```
from modeller.automodel import *
from modeller.scripts import cispeptide

# Redefine the special_restraints routine to force Pro to cis conformation:
# (this routine is empty by default):
class mymodel(automodel):
    def special_restraints(self, aln):
        cispeptide(self.restraints, atom_ids1=('O:4', 'C:4', 'N:5', 'CA:5'),
                  atom_ids2=('CA:4', 'C:4', 'N:5', 'CA:5'))

# This is as usual:
log.verbose()
env = environ()

a = mymodel(env, alnfile='align1.ali', knowns='templ1', sequence='targ1')
a.make()
```

5. How can I select/remove/add a set of restraints?

Restraints can be read from a file by `model.restraints.append()`, calculated by `model.restraints.make()`, or added “manually” by `model.restraints.add()`. `model.restraints.pick()` picks those restraints for objective function calculation that restrain the selected atoms only, as specified in the selected

atoms set 1. Initially, all atoms are selected; this can be changed by the `model.pick_atoms()` command. `model.restraints.make()` command for some restraint types (*e.g.*, distance) constructs restraints of the selected type between the atoms in the selected atoms sets 2 and 3. The `'automodel.homcsr()'` routine contains examples of the `model.pick_atoms()` command when generating restraints by `model.restraints.make()`. There are also commands for adding and unselecting single restraints, `model.restraints.add()` and `model.restraints.unpick()`, respectively. If you do `model.restraints.condense()`, the unselected restraints will be deleted. This is useful for getting rid of the unwanted restraints completely.

6. What are the different refinement levels really doing?

There are two different optimization approaches available within MODELLER: variable target function method (VTFM) with conjugate gradients (CG) [Šali & Blundell, 1993] and molecular dynamics (MD) with simulated annealing (SA) [Šali & Blundell, 1993]. They can both be done to a different degree (with more or less cycles of CG and MD, faster or slower schedule for VTFM and SA). The exact details are best obtained from the scripts themselves because a detailed description would probably be longer than the scripts. For example, the QUANTA and INSIGHTII implementations of MODELLER have these three levels of optimization: no optimization (only copying coordinates from templates and building the undefined atoms using internal geometry information from the RTF entries); only VTFM with CG; also MD with SA. Most of the time (70%) is spent on the MD&SA part. Our experience is that when MD&SA are used, if there are violations in the best of the 10 models, they probably come from an alignment error, not an optimizer failure (if there are no insertions longer than approximately 15 residues).

7. I want to change the default optimization schedule.

See chapter 2 for the variables that could be changed and for their possible values.

```
from modeller.automodel import *

log.verbose()
env = environ()

a = automodel(env, alnfile='align1.ali', knowns='templ1', sequence='targ1')
# Very thorough VTFM optimization:
a.library_schedule = 1
a.max_var_iterations = 300
# Very thorough MD optimization:
a.md_level = refine.very_slow
# Repeat the whole cycle 3-times and do not stop unless obj.func. > 1E6
a.repeat_optimization = 3
a.max_molpdf = 1e6

a.make()
```

8. I want to build an all hydrogen atom model with water molecules and other non-protein atoms (atoms in the HETATM records in the PDB file).

```
from modeller.automodel import *

log.verbose()
env = environ()
env.io.hydrogen = env.io.hetatm = env.io.water = True

a = allhmodel(env, alnfile='align1.ali', knowns='templ1', sequence='targ1')
a.make()
```

9. How do I build a model with water molecules or residues that do not have an entry in the topology and/or parameter files?

Water molecules are indicated by 'w' in the alignment file and the special block residue ('BLK') that does not have entries in the residue topology and parameter libraries is indicated by '.'.

See Section 4.5.1 for information about block residues.

```
from modeller.automodel import *

log.verbose()
env = environ()
env.io.hetatm = env.io.water = True

a = automodel(env, alnfile='align1.ali', knowns='templ1', sequence='targ1')
a.make()
```

The alignment file:

```
>P1;templ1
structureX:templ1:1::10::
FAYVI/.wwww*

>P1;targ1
sequence:targ1:1::8::
-GWIV/.ww-w*
```

10. How do I define my own residue types, such as D-amino acids, special ligands, and unnatural amino-acids?

This is a painful area in all molecular modeling programs. However, CHARMM and X-PLOR provide a reasonably straightforward solution *via* the residue topology and parameter libraries. MODELLER uses CHARMM topology and parameter library format and also extends the options by allowing for a generic "BLK" residue type (Section 4.5.1). This BLK residue type circumvents the need for editing any library files, but it is not always possible to use it. Due to its conformational rigidity, it is also not as accurate as a normal residue type. In order to define a new residue type in the MODELLER libraries, you have to follow the series of steps described below. As an example, we will define the ALA residue without any hydrogen atoms. You can add an entry to the MODELLER topology or parameter file; you can also use your own topology or parameter files. For more information, please see the CHARMM manual.

(a) Define the new residue entry in the residue topology file (RTF), say 'top_heav.lib'.

```
RESI ALA      0.00000
ATOM N      NH1    -0.29792
ATOM CA     CT1     0.09563
ATOM CB     CT3    -0.17115
ATOM C      C       0.69672
ATOM O      O      -0.32328
BOND CB CA      N CA      O C      C CA      C +N
IMPR C CA +N O      CA N C CB
IC -C      N      CA      C          1.3551  126.4900  180.0000  114.4400  1.5390
IC N       CA      C      +N         1.4592  114.4400  180.0000  116.8400  1.3558
IC +N      CA      *C      O          1.3558  116.8400  180.0000  122.5200  1.2297
IC CA      C      +N      +CA        1.5390  116.8400  180.0000  126.7700  1.4613
IC N       C      *CA     CB          1.4592  114.4400  123.2300  111.0900  1.5461
IC N       CA      C      O          1.4300  107.0000   0.0000  122.5200  1.2297
PATC FIRS  NTER  LAST  CTER
```

You can obtain an initial approximation to this entry by defining the new residue type using the residue type editor in QUANTA and then writing it to a file.

The RESI record specifies the CHARMM residue name, which can be up to four characters long and is usually the same as the PDB residue name (exceptions are the potentially charged residues where the different charge states correspond to different CHARMM residue types). The number gives the total residue charge.

The ATOM records specify the IUPAC (*i.e.*, PDB) atom names and the CHARMM atom types for all the atoms in the residue. The number at the end of each ATOM record gives the partial atomic charge. The BOND records specify all the covalent bonds between the atoms in the residue (*e.g.*, there are bonds CB-CA, N-CA, O-C, etc.). In addition, symbol '+' is used to indicate the bonds to the subsequent residue in the chain (*e.g.*, C +N). The covalent angles and dihedral angles are calculated automatically from the list of chemical bonds.

The IMPR records specify the improper dihedral angles, generally used to restrain the planarity of various groups (*e.g.*, peptide bonds and sidechain rings). See also below.

The IC (internal coordinate) records are used for constructing the initial Cartesian coordinates of a residue. An entry

$$IC \quad a \quad b \quad c \quad d \quad d_{ab} \quad \alpha_{abc} \quad \Theta_{abcd} \quad \alpha_{bcd} \quad d_{cd}$$

specifies distances d , angles α , and either dihedral angles or improper dihedral angles Θ between atoms a , b , c and d , given by their IUPAC names. The improper dihedral angle is specified when the third atom, c , is preceded by a star, '*'. As before, the '-' and '+' pre-fixes for the atom names select the corresponding atom from the preceding and subsequent residues, respectively. The distances are in angstroms, angles in degrees. The distinction between the dihedral angles and improper dihedral angles is unfortunate since they are the same mathematically, except that by convention when using the equations, the order of the atoms for a dihedral angle is $abcd$ and for an improper dihedral angle it is $acbd$.

The PATC record specifies the default patching residue type when the current residue type is the first or the last residue in a chain.

- (b) You have to make sure that all the CHARMM atom types of the new residue type occur in the MASS records at the beginning of the topology library: Add your entry at the end of the MASS list if necessary. If you added any new CHARMM atom types, you also have to add them to the radii libraries, 'modlib/radii.lib' and 'modlib/radii14.lib'. These libraries list the atomic radii for the different topology models, for the long range and 1-4 non-bonded soft-sphere terms, respectively. The full names of the files that are used during calculation are given by the environment variables \$RADII.LIB and \$RADII14.LIB.
- (c) Optionally, you can add the residue entry to the library of MODELLER topology models, 'modlib/models.lib'. The runtime version of this library is specified by the environment variable \$MODELS.LIB. This library specifies which subsets of atoms in the residue are used for each of the possible topologies. Currently, there are 10 topologies selected by topology.submodel (3 is default):

1	ALLH	all atoms
2	POL	polar hydrogens only
3	HEAV	non-hydrogen atoms only
4	MCCB	non-hydrogen mainchain (N, C, CA, O) and CB atoms
5	MNCH	non-hydrogen mainchain atoms only
6	MCWO	non-hydrogen mainchain atoms without carbonyl O
7	CA	CA atoms only
8	MNSS	non-hydrogen mainchain atoms and disulfide bonds
9	CA3H	reduced model with a small number of sidechain interaction centers
10	CACB	CA and CB atoms only

The Ala entry is:

```
#
      ALLH POLH HEAV MCCB MNCH MCWO CA  MNSS CA3H CACB
*
RESI ALA
ATOM   NH1  NH1  NH1  NH1  NH1  NH1  ##### NH1  ##### #####
```

```

ATOM      H      HN      #### #### #### #### #### #### ####
ATOM      CT1    CT1    CT1    CT1    CT1    CT1    CT1    CT1    CAH    CT1
ATOM      HB      #### #### #### #### #### #### #### #### CH3E ####
ATOM      CT3    CT3    CT3    CT3    #### #### #### #### #### CT2
ATOM      HA      #### #### #### #### #### #### #### #### ####
ATOM      HA      #### #### #### #### #### #### #### #### ####
ATOM      HA      #### #### #### #### #### #### #### #### ####
ATOM      C      C      C      C      C      C      #### C      ####
ATOM      O      O      O      O      O      #### #### O      ####

```

The residue entries in this library are separated by stars. The '####' string indicates a missing atom. The atom names for the present atoms are arbitrary. The order of the atoms must be the same as in the CHARMM residue topology library. If a residue type does not have an entry in this library, all atoms are used for all topologies.

- (d) You have to add the new residue type to the residue type library, 'modlib/restyp.lib'. The execution version of this file is specified by the environment variable \$RESTYP_LIB. For the ALA residue,

```

1 | ALA          | A | ALA | alanine

```

You would generally add the new residue type at the end of the file. There are 5 fields in each line, separated by the '|' characters. The first field is an integer index corresponding to the integer residue type. The standard residue types have their indices smaller than 24. These are also the indices corresponding to the residue-residue substitution matrices. The second field contains the list of equivalent PDB or IUPAC 3-character residue names, used in the PDB files. A list rather than a single name is allowed because PDB can unfortunately use different names for the same residue type (*e.g.*, water can be HOH, WAT, etc.). The third field gives a single character code for the residue type, which is used in the alignment file. This does not have to be unique, but if it is not unique you cannot use it in the alignment file. Any ASCII character is fine, it does not have to be a letter. If you run out of characters you can re-define the existing ones that you do not need. The fourth field gives the four-character CHARMM residue name, as specified in the RESI record of the topology library. The last field contains an optional comment.

Every residue in the CHARMM topology file has to have an entry in the \$RESTYP_LIB library, but not every residue entry in the \$RESTYP_LIB library needs an entry in the residue topology file. If you need the edit the \$RESTYP_LIB file, it is recommended that you change a copy of it, and provide that file to the `environ()` constructor.

- (e) In general, when you add a new residue type, you also add new chemical bonds, angles, dihedral angles, improper dihedral angles, and non-bonded interactions, new in the sense that a unique combination of CHARMM atoms types is involved whose interaction parameters are not yet specified in the parameter library (see also Section 4.5.1). In such a case, you will get a number of warning and/or error messages when you generate the stereochemical restraints by the `model.restraints.make()` command. These messages can sometimes be ignored because MODELLER will guess the values for the missing parameters from the current Cartesian coordinates of the model. When this is not accurate enough or if the necessary coordinates are undefined you have to specify the parameters explicitly in the parameter library. Search for BOND, ANGL, DIHE, and IMPR sections in the parameters library file and use the existing entries to guess your new entries. Note that you can use dummy atom types 'X' to create general dihedral (*i.e.*, X A A X) and improper dihedral angle (*i.e.*, A X X A) entries, where A stands for any of the real CHARMM atom types. For the dihedral angle cosine terms, the CHARMM convention for the phase is different for 180° from MODELLER's (Eq. 6.57). If you use non-bonded Lennard-Jones terms, you also have to add a NONB entry for each new atom type. If you use the default soft-sphere non-bonded restraints, you have already taken care of it by adding the new atom types to the \$RADII_LIB and \$RADII_LIB libraries.

11. How do I define my own patching residue types?

This is even messier than defining a new residue type. As an example, we will define the patching residue for establishing a disulfide bond between two CYS residues.

```

PRES DISU          -0.36 ! Patch for disulfides. Patch must be 1-CYS and 2-CYS.

```

```

ATOM 1:CB CT2 -0.10 !
ATOM 1:SG SM -0.08 ! 2:SG--2:CB--
ATOM 2:SG SM -0.08 ! /
ATOM 2:CB CT2 -0.10 ! -1:CB--1:SG
DELETE ATOM 1:HG
DELETE ATOM 2:HG
BOND 1:SG 2:SG
IC 1:CA 1:CB 1:SG 2:SG 0.0000 0.0000 180.0000 0.0000 0.0000
IC 1:CB 1:SG 2:SG 2:CB 0.0000 0.0000 90.0000 0.0000 0.0000
IC 1:SG 2:SG 2:CB 2:CA 0.0000 0.0000 180.0000 0.0000 0.0000

```

The PRES record specifies the CHARMM patching residue type (up to four characters). As for the normal RESI residue types, patching residue types also have to be defined in the residue type library, 'modlib/restyp.lib'.

The ATOM records have the same meaning as for the RESI residue types described above. The extension is that the IUPAC atom names (listed first) must be pre-fixed by the index of the residue that is patched. In this example, there are two CYS residues that are patched, thus the prefixes 1 and 2. When using the `model.patch()` command, the order of the patched residues specified by `residue_ids` must correspond to these indices (this is only important when the patch is not symmetric, unlike the 'DISU' patch in this example).

DELETE records specify the atoms to be deleted, the two hydrogens bonded to the two sulfurs in this case. The BOND and IC (internal coordinate) records are the same as those for the RESI residues, except that the atom names are prefixed with the patched residue indices.

12. Is it possible to restrain secondary structure in the target sequence?

Yes. There are 'ALPHA', 'STRAND' and 'SHEET' restraint types that the `model.restraints.make()` command can generate. One specifies the segment which is then restrained to the specified secondary structure conformation. For example,

```

from modeller.automodel import *

# Redefine the special_restraints routine to include the secondary
# structure restraints (this routine is empty by default):
def mymodel(automodel):
    def special_restraints(self, aln):
        rsr = self.restraints
        # An alpha-helix:
        rsr.make(aln, restraint_type='ALPHA', residue_ids=('20', '30'))
        # Two strands:
        rsr.make(aln, restraint_type='STRAND', residue_ids=('1', '6'))
        rsr.make(aln, restraint_type='STRAND', residue_ids=('9', '14'))
        # An anti-parallel sheet:
        rsr.make(aln, restraint_type='SHEET', atom_ids=('N:1', 'O:14'),
                sheet_h_bonds=-5)

# This is as usual:
log.verbose()
env = environ()

a = automodel(env, alnfile='align1.ali', knowns='templ1', sequence='targ1')
a.make()

```

13. I want to patch the N-terminal or (C-terminal) residue (e.g., to model acetylation properly), but the model.patch() command does not work.

This is probably because the N-terminus is patched by default with the NTER patching residue (corresponding to $-\text{NH}_3^+$) and a patched residue must not be patched again. The solution is to turn the default patching off by `env.patch_default = False` before the `model.generate_topology()` command is called.

14. **Is it possible to use templates with the coordinates for C_α atoms only?**

Yes. You do not have to do anything special.

15. **How do I analyze the output log file?**

First, check for the error messages by searching for string '`_E>`'. These messages can only rarely be ignored. Next, check for the warning messages by searching for string '`_W>`'. These messages can almost always be ignored. If everything is OK so far, the most important part of the log file is the output of the `model.energy()` command for each model. This is where the violations of restraints are listed. When there are too many too violated restraints, more optimization or a different alignment is needed. What is too many and too much? It depends on the restraint type and is best learned by doing `model.energy()` on an X-ray structure or a good model to get a feel for it. You may also want to look at the output of command `alignment.check()`, which should be self-explanatory. I usually ignore the other parts of the log file.

16. **How do I prevent “knots” in the final models?**

The best way to prevent knots is to start with a starting structure that is as close to the desired final model as possible. Other than that, the only solution at this point is to calculate independently many models and hope that in some runs there won't be knots. Knots usually occur when one or more neighboring long insertions (*i.e.*, longer than 15 residues) are modeled from scratch. The reason is that an insertion is build from a randomized distorted structure that is located approximately between the two anchoring regions. Under such conditions, it is easy for the optimizer to “fall” into a knot and then not be able to recover from it. Sometimes knots result from an incorrect alignment, especially when more than one template is used. When the alignment is correct, knots are a result of optimization not being good enough. However, making optimization more thorough by increasing the CPU time would not be worth it on the average as knots occur relatively infrequently. The excluded volume restraints are already included in the standard comparative modeling routine.

17. **What is considered to be the minimum length of a sequence motif necessary to derive meaningful constraints from the alignment to use in modeling.. one, two, three, or more?**

Usually more than that (dozens if you want just to detect reliable similarity, and even more if you want a real model). It is good to have at least 35-40% sequence identity to build a model. Sometimes even 30% is OK.

18. **Does Modeller have a graphical interface (GUI) ?**

No; Modeller is run from the command line, and uses a Python script to direct it. However, a graphical interface to Modeller is commercially available from Accelrys, as part of Discovery Studio Modeling 1.1, at http://www.accelrys.com/dstudio/ds_modeling/ds_modeler.html.

19. **What do the ‘Alignment sequence not found in PDB file’ or ‘Number of residues in the alignment and pdb files are different’ errors mean?**

When you give MODELLER an alignment, it also needs to read the structure of the known proteins (templates) from PDB files. In order to correctly match coordinates to the residues specified in the alignment, the sequences in the PDB file and the alignment file must be the same (although obviously you can add gap or chain break characters to your alignment). If they are not, you see this error. (Note that MODELLER takes the PDB sequence from the ATOM and HETATM PDB records, not the SEQRES records.) You should also check the header of your alignment file, to make sure that you are reading the correct chain and residue numbers from your PDB.

To see the sequence that MODELLER reads from the PDB file '`1BY8.pdb`', use this short script to produce a '`1BY8.seq`' sequence file:

```
env = environ()
code = '1BY8'
mdl = model(env, file=code)
aln = alignment(env)
aln.append_model(mdl, align_codes=code)
aln.write(file=code+'.seq')
```

1.9 MODELLER updates

1.9.1 Changes since release 8v0

- Removal of limits on the number of pseudo atoms, segments, the length of file names and most other strings, and the number of parameters in parameter files.
- Removal of the `read_restyp()` command. To read a user-defined residue type file, set `restyp_lib_file` when creating an environ object.
- Bug fix: the correct VDW radii (as selected by `topology.submodel`) are now used for PSA calculations (previously the all-hydrogen radii were used in all cases).
- Improved error handling: almost all errors should be recoverable now (using the standard Python exception handling mechanism) and incorrect input should trigger an error rather than causing Modeller to crash (please report any crashes as bugs to the developers).
- SALIGN bug fixes, improved binning of access classes.
- Modeller scripts are now more like 'real' Python scripts; the PYTHONPATH environment variable is honored, and additional command line arguments can be given to your scripts (accessible as `sys.argv`).
- Bug fix for `alignment.segment_matching()`: correct alignment scores are now written.
- Support for AMD64 and EM64T (Opteron/Xeon64) Linux machines.
- DOPE model assessment now returns only the DOPE energy itself by default, ignoring other energy terms.
- Bug fix for `model.restraints.pick()` when used with absolute position restraints; always allow 1-atom restraints to pass the `residue.span_range` check.

1.9.2 Changes since release 7v7

- TOP is no longer used as the scripting language. MODELLER now reads Python 2.3 scripts instead. (TOP scripts can still be used, but will not take advantage of all new features.) See Chapter 5 for more information.
- Incorporation of the GA341 (`model.assess_ga341()`) and DOPE (`model.assess_dope()`) methods for model assessment.
- 'Automatic' filename generation (when filenames were set to 'default') is no longer supported; use the `modfile.default()` function instead.
- Memory is used much more efficiently, and various limits evident in older releases (*e.g.*, errors concerning MAXRES, MAXCNS and MAXBND) have been removed. Consequently, you can no longer assign to the MAXRES variable; this is not necessary now.
- New `alignment.salign()` alignment command, using structural information, and based on the COMPARER program.
- Various fixes to `profile.build()`, and a new command `profile.scan()`, for building and using sequence profiles.
- Generated models can be docked to electron microscopy density maps, using the new `density.grid_search()` command.
- Superpositions are now improved by using local substructures. This local refinement can be disabled (thus getting the behavior of old MODELLER versions) by turning off the `refine_local` option to `model.superpose()`.
- MODELLER will now read compressed files in the .gz and .bz2 formats (as well as the existing .Z format) provided the gzip and bzip2 programs are available.
- Fixes for various problems reading alignment files with very long lines, or with Mac, Unix or DOS-style line endings.

- Bug fix for `alignment.compare_structures()`; disulfide bridges were not being properly considered in the comparison.
- Charges for all topologies are now assigned using a ‘nearest neighbor’ technique; see `topology.make()`.
- `sequence_db.search()` (formerly `SEQUENCE_SEARCH`) now takes a 1-sequence alignment as input, rather than reading in that alignment itself.

1.9.3 Changes since release 6v2

- The `CUT_OVERHANGS` argument to the `WRITE_ALIGNMENT` command has been removed, and replaced by a more powerful `EDIT_ALIGNMENT` command.
- The `MAX_LOOP_LENGTH` argument to the `PICK_ATOMS` command has been replaced by `MIN-MAX_LOOP_LENGTH`, such that both the minimum and maximum loop segment lengths can be selected for in ‘SEGMENT’ mode.
- The `SEARCH_CHAINS_LIST` and `SEARCH_CHAINS_FILE` arguments to `SEQUENCE_SEARCH` have been removed. Instead, a sequence database must now be read into memory prior to using `SEQUENCE_SEARCH` by using the new `READ_SEQUENCE_DB` command. (Such a database can also be written out with `WRITE_SEQUENCE_DB`.)
- New ‘all-hydrogens’ example directory, for building all-hydrogen models.
- CHARMM topology libraries are now PDB (IUPAC) compliant. Most obviously, this has resulted in the following name changes:
 - The HSD (neutral histidine) residue is now HIS.
 - The ILE CD atom is now CD1.
 - The LEU CD1 and CD2 atoms have been swapped.
 - PDB hydrogen atom naming conventions now apply.

See also the ‘top-charmmH.lib’ file for the old CHARMM-style naming, and the comments at the start of ‘top.lib’.

- `READ_ALIGNMENT` can now read ‘FASTA’ format alignments. Additionally, the `CLOSE_FILE`, `REWIND_FILE`, and `END_OF_FILE` variables can be used to read partial ‘PIR’ or ‘FASTA’ files.
- `MALIGN3D` allows the filenames of fitted atom files to be customized with the `EDIT_FILE_EXT` variable.
- `MAKE_RESTRAINTS` can now impose additional gap-distance weighting on distance restraints, using the `RESTRAINT_STDEV2` variable.
- New commands for dealing with profiles: `ALN_TO_PROF`, `PROF_TO_ALN`, `WRITE_PROFILE`, `READ_PROFILE`, `BUILD_PROFILE`.
- New options to `ALIGN` and `ALIGN2D` for dealing with profiles: `WEIGH_SEQUENCES`, `SMOOTH_PROF_WEIGHT`, `READ_PROFILE`, `INPUT_PROFILE_FILE`, `WRITE_PROFILE`, `OUTPUT_PROFILE_FILE`. `ALIGN_WHAT` can also now take the value ‘PROFILE’.
- Other new commands: `SEQFILTER`, `TIME_MARK`, `MAKE_CHAINS`, `VOLUME`, `VOLUME_CAVITY`.

1.9.4 Changes between releases 4 and 6v2

MODELLER 5 has not been generally released; the major changes in MODELLER 6v2 relative to version 4 include (this is a very incomplete list):

- New TOP loop modeling routine 'loop' significantly improves the accuracy of loop modeling (Section 2.2.9). It can also be used in modeling of sidechains or other parts of the structure.
- Several new statistical atomic distance-dependent potentials can typically be used for loop modeling and model evaluation.
- The ENERGY OUTPUT = 'ENERGY_PROFILE' allows construction of energy profiles for model evaluation, based on any combination of MODELLER energy terms. It replaces the **ENERGY_PROFILE** command.
- New functional form and parameters for the binormal Φ, Ψ restraints improve their numerical stability and accuracy.
- New TOP command **SEGMENT_MATCHING** facilitates exploring many different alignments in the difficult comparative modeling cases.
- New dynamic memory allocation makes MODELLER more memory efficient.
- Many bug fixes.
- Many new arguments, changed naming/meaning of arguments, and several new commands.
- Conversion from Fortran 77 to Fortran 90 was the main culprit for a very long delay before the latest release.

Chapter 2

Automated comparative modeling with automodel

The simplest way to build comparative models with MODELLER is to use the `automodel` class. This automates many of the steps required for simple modeling, and can be customized for more complex tasks. It also allows refinement of loop regions.

See section 6.4 for a flowchart of the comparative modeling procedure, and section 6.5 for an overview of the loop modeling algorithm.

This chapter gives an overview of simple applications of the `automodel` class. For more detailed information, see the comparative modeling class reference, in Chapter 3, or the comments in the Python file itself, `modlib/modeller/automodel/automodel.py`.

2.1 Simple usage

The simple example below constructs a single comparative model for the `1fdx` sequence from the known `5fd1` structure, using `alignment.ali`, a PIR format alignment of `5fd1` and `1fdx`. The final model is written into the PDB file `1fdx.B99990001.pdb`. See Section 1.7.2 for instructions on how to run this script.

Example: `examples/automodel/model-default.py`

```
# Homology modelling by the automodel class

from modeller.automodel import *    # Load the automodel class

log.verbose()    # request verbose output
env = environ() # create a new MODELLER environment to build this model in

# directories for input atom files
env.io.atom_files_directory = './:../atom_files'

a = automodel(env,
               alnfile = 'alignment.ali',    # alignment filename
               knowns   = '5fd1',           # codes of the templates
               sequence = '1fdx')           # code of the target
a.starting_model= 1    # index of the first model
a.ending_model  = 1    # index of the last model
                  # (determines how many models to calculate)
a.make()             # do the actual homology modelling
```

Example: examples/automodel/alignment.ali

```

C; A sample alignment in the PIR format; used in tutorial

>P1;5fd1
structureX:5fd1:1      : :106  : :ferredoxin:Azotobacter vinelandii: 1.90: 0.19
AFVVTDNCIKCKYTDCEVEVCPVDCFYEGPNFLVIHPDECIDCALCEPECPAQAI FSEDEVPEDMQEFIQLN AELA
EVWPNITEKKDPLPDAEDWDGKGLQHLER*

>P1;1fdx
sequence:1fdx:1      : :54   : :ferredoxin:Peptococcus aerogenes: 2.00:-1.00
AYVINDSC--IACGACKPECPVNI IQGS--IYAIDADSCIDCGSCASVCPVGAPNPED-----
-----*

```

Stepping through the script, first we load the `automodel` class, using standard Python syntax to load a module. Next, we request verbose output (see Section 4.14) so that we can more easily spot errors. We then create an `environ()` object (see Section 4.2) and call it `env`. This object holds the MODELLER ‘environment’, which comprises default values for many parameters, as well as the libraries used for comparative modeling (topology, parameters, dihedral classes, *etc*). An `environ` object is needed to create most other MODELLER objects, but you can call it whatever you like (it doesn’t have to be called `env`).

Once we have the `environ` object, we can set some global parameters. In this case, we set `io_data.atom_files_directory` to set the directories to look for PDB files in.

Next, we create an `automodel` object, tell it which PIR alignment file to use, and which sequences are templates and which one we want to build a model for, and call it `a`. This doesn’t actually build any models, but creates the object, ready to be tweaked for our purposes. In this case, we simply tell it to build a single model, by setting both `automodel.starting_model` and `automodel.ending_model` to 1. Finally, we actually build the model by running `automodel.make()`.

2.2 More advanced usage

2.2.1 Including water molecules, HETATM residues, and hydrogen atoms

This can be done quite easily by setting `env.io.water`, `env.io.hetatm` and `env.io.hydrogen` respectively to `True` (they are `False` by default). Note that `env.io.hydrogen` is automatically turned on for you if you want to build an all hydrogen model, using the `allhmodel` class; see section 2.2.3.

Note that simply turning on `env.io.hetatm` is **not** usually sufficient to build models containing ligands or unnatural amino acids; MODELLER also requires topology and parameter information. In such a case, you can either provide such information, or you can treat the residues as rigid bodies (‘BLK’ ligands). See FAQ entry 9 and section 4.5.1 for more information, and the example in the advanced modeling tutorial, at <http://salilab.org/modeller/tutorial/advanced.html>.

2.2.2 Getting a very fast and approximate model

To get an approximate model very quickly (to get a rough idea of what it looks like, or to confirm that the alignment is reasonable) call the `automodel.very_fast()` method before `automodel.make()`. This uses only a very limited amount of variable target function optimization with conjugate gradients, and thus is roughly 3 times faster than the default procedure.

Note that no randomization of the starting structure is done in this case, so only a single model can be produced.

This example also demonstrates the use of the `assess_methods` keyword, to request model assessment. In this case the GA341 method is requested. See section 3.1.1.

Example: examples/automodel/model-fast.py

```

# Very fast homology modelling by the automodel class

from modeller.automodel import *    # Load the automodel class

log.verbose()
env = environ()
# directories for input atom files
env.io.atom_files_directory = './:../atom_files'

a = automodel(env,
              alnfile='alignment.ali',    # alignment filename
              knowns='5fd1',             # codes of the templates
              sequence='1fdx',           # code of the target
              assess_methods=assess.GA341) # request GA341 model assessment

a.very_fast()                          # prepare for extremely fast optimization

a.starting_model = 2
a.ending_model = 2
a.final_malign3d = True

a.make()                                # make the homology model

```

2.2.3 Building an all hydrogen model

This is done by using the `allhmodel` class rather than `automodel`. Otherwise, operation is identical. Note that the `allhmodel` class automatically turns on `env.io.hydrogen` for you and selects the all-atom topology and radii files.

Example: `examples/automodel/model-default-allh.py`

```

from modeller.automodel import *

log.verbose()
env = environ()

env.io.atom_files_directory = './:../atom_files'

a = allhmodel(env, alnfile='alignment.ali', knowns='5fd1', sequence='1fdx')
a.starting_model = a.ending_model = 4

a.make()

```

2.2.4 Refining only part of the model

The `automodel` class contains a `'select_atoms'` method which selects the atoms to be moved during optimization. By default, the routine selects all atoms, but you can redefine it to select any subset of atoms and then only those atoms will be refined. (To redefine the routine, it is necessary to create a subclass of `automodel`, here called `mymodel`, which has the modified routine within it. We then use `mymodel` in place of `automodel`.)

By default, these atoms will “feel” the presence of other atoms via all the static and possibly dynamic restraints that include both selected and un-selected atoms. However, you can turn off dynamic interactions between the

selected and unselected regions by setting `env.edat.nonbonded_sel_atoms` to 2 (by default it is 1). For example, the script below would refine only atoms in residues 1 and 2. The difference between this script and the one for loop modeling is that here the selected regions are optimized with the default optimization protocol and the default restraints, which generally include template-derived restraints. In contrast, the loop modeling routine does not use template-dependent restraints, but does a much more thorough optimization.

Example: `examples/automodel/model-segment.py`

```
# Homology modelling by the automodel class
#
# Demonstrates how to refine only a part of the model.
#
# You may want to use the more exhaustive "loop" modeling routines instead.
#

from modeller.automodel import *      # Load the automodel class

log.verbose()

# Override the 'select_atoms' routine in the 'automodel' class
class mymodel(automodel):
    def select_atoms(self):
        self.pick_atoms(selection_segment=('1:', '2:'),
                        selection_search='segment', pick_atoms_set=1,
                        res_types='all', atom_types='all',
                        selection_from='all', selection_status='initialize')

env = environ()
# directories for input atom files
env.io.atom_files_directory = './:../atom_files'
# selected atoms do not feel the neighborhood
env.edat.nonbonded_sel_atoms = 2

a = mymodel(env,
             alnfile = 'alignment.ali',      # alignment filename
             knowns   = '5fd1',             # codes of the templates
             sequence = '1fdx')              # code of the target

a.starting_model= 3                # index of the first model
a.ending_model  = 3                # index of the last model
# (determines how many models to calculate)
a.make()                            # do homology modelling
```

2.2.5 Providing your own restraints file

To force `automodel` not to construct restraints at all, but to instead use your own restraints file, simply use the `csrfile` keyword when creating the `automodel` class, as in the example below. Note that MODELLER does only rudimentary checking on your restraints file, so you must be careful that it applies correctly to the generated model.

Example: `examples/automodel/model-mysr.py`

```
# Modeling using a provided restraints file (csrfile)
```

```

from modeller.automodel import *    # Load the automodel class

log.verbose()
env = environ()

# directories for input atom files
env.io.atom_files_directory = './:../atom_files'

a = automodel(env,
               alnfile = 'alignment.ali',    # alignment filename
               knowns   = '5fd1',           # codes of the templates
               sequence = '1fdx',           # code of the target
               csrfile  = 'my.rsr')         # use 'my' restraints file
a.starting_model= 1                    # index of the first model
a.ending_model  = 1                    # index of the last model
                                           # (determines how many models to calculate)
a.make()                                # do homology modelling

```

2.2.6 Using your own initial model

Normally, `automodel` generates an initial model by transferring coordinates from the templates. However, if you have a prepared PDB file containing an initial model, you can have `automodel` use this instead with the `inifile` keyword, as in the example below. (This automatically sets `automodel.generate_method` to `generate.read_xyz` for you, which is necessary for this to work.)

Example: `examples/automodel/model-myini.py`

```

# Modeling using a provided initial structure file (inifile)

from modeller.automodel import *    # Load the automodel class

log.verbose()
env = environ()

# directories for input atom files
env.io.atom_files_directory = './:../atom_files'

a = automodel(env,
               alnfile = 'alignment.ali',    # alignment filename
               knowns   = '5fd1',           # codes of the templates
               sequence = '1fdx',           # code of the target
               inifile  = 'my-initial.pdb')  # use 'my' initial structure
a.starting_model= 1                    # index of the first model
a.ending_model  = 1                    # index of the last model
                                           # (determines how many models to calculate)
a.make()                                # do homology modelling

```

2.2.7 Adding additional restraints to the defaults

You can add your own restraints to the restraints file, with the homology-derived restraints, by redefining the `'special_restraints'` routine (by default it does nothing). This can be used, for example, to add information from NMR experiments. The example below enforces an additional restraint on a single CA-CA distance.

Example: examples/automodel/model-addrsr.py

```

# Addition of restraints to the default ones

from modeller.automodel import *    # Load the automodel class

log.verbose()
env = environ()

# directories for input atom files
env.io.atom_files_directory = './:../atom_files'

class mymodel(automodel):
    def special_restraints(self, aln):
        rsr = self.restraints
#       Add some restraints from a file:
#       rsr.append(file='my_rsrs1.rsr')
#       Restrain the specified CA-CA distance to 10 angstroms (st. dev.=0.1)
#       Use a harmonic potential and X-Y distance group.
        rsr.add(atom_ids=('CA:35', 'CA:40'),
                restraint_parameters=(3, 1, 1, 27, 2, 2, 0, 10.0, 0.1))

a = mymodel(env,
            alnfile = 'alignment.ali',    # alignment filename
            knowns  = '5fd1',            # codes of the templates
            sequence = '1fdx')           # code of the target
a.starting_model= 1                    # index of the first model
a.ending_model  = 1                    # index of the last model
# (determines how many models to calculate)
a.make()                                # do homology modelling

```

2.2.8 Fully automated alignment and modeling

If you do not have an initial alignment between your templates and target sequence, MODELLER can derive one for you, fully automatically. All MODELLER requires is a PIR file containing the target sequence and the template PDB codes (their sequences are not required — just use a single '*' character — as MODELLER will read these from the PDBs). Use the `automodel` class as per usual, but call the `automodel.auto_align()` method before `automodel.make()`; see the example below. (MODELLER has a variety of other alignment methods which you can use instead for this purpose; see Section 4.9 for more details.)

Please be aware that the single most important factor that determines the quality of a model is the quality of the alignment. If the alignment is incorrect, the model will also be incorrect. **For this reason, automated alignment for comparative modeling should not be used unless the sequences are so similar that the calculated alignment is likely to be correct, which usually requires more than 50% sequence identity.** Instead, the alignment should be carefully inspected, optimized by hand, and checked by the `alignment.check()` command before used in modeling. Moreover, several iterations of alignment and modeling may be necessary in general.

Example: examples/automodel/model-full.py

```

# A sample script for fully automated comparative modeling

from modeller.automodel import *    # Load the automodel class

log.verbose()

```

```

env = environ()

# directories for input atom files
env.io.atom_files_directory = './:../atom_files'

a = automodel(env,
              # file with template codes and target sequence
              alnfile = 'alignment.seg',
              # PDB codes of the templates
              knowns = ('5fd1', '1fdn', '1fxd', '2fxb'),
              # code of the target
              sequence = '1fdx')
a.auto_align()          # get an automatic alignment
a.make()                # do homology modelling

```

Example: examples/automodel/alignment.seg

```

>P1;1fdx
structureX:1fdx:FIRST:@:54:@:ferredoxin:Peptococcus aerogenes: 2.00:-1.00
AYVINDSCIACGACKPECPVNIQGSIIYAIADADSCIDCGSCASVCPVGAPNPED*
>P1;1fdn
structureX:1fdn:FIRST:@:55:@:ferredoxin:Clostrodiium acidurici: 1.84:-1.0
*
>P1;5fd1
structureX:5fd1:FIRST:@:60:@:ferredoxin:Azotobacter vinelandii: 1.90:0.192
*
>P1;1fxd
structureX:1fxd:FIRST:@:58:@:ferredoxin:Desolfovibrio gigas: 1.70:-1.0
*
>P1;2fxb
structureX:2fxb:FIRST:@:60:@:ferredoxin:Bacillus thermoproteolyticus: 2.30:-1.0
*

```

2.2.9 Loop optimization

The new loop optimization method relies on a scoring function and optimization schedule adapted for loop modeling [Fiser *et al.*, 2000]. It is used automatically to refine comparative models if you use the `loopmodel` class rather than `automodel`; see the example below.

Example: examples/automodel/model-loop.py

```

# Homology modelling by the automodel class

from modeller.automodel import *    # Load the automodel class

log.verbose()
env = environ()

# directories for input atom files
env.io.atom_files_directory = './:../atom_files'

```

```

a = loopmodel(env,
               alnfile = 'alignment.ali',      # alignment filename
               knowns   = '5fd1',             # codes of the templates
               sequence = '1fdx')              # code of the target
a.starting_model= 1      # index of the first model
a.ending_model  = 1      # index of the last model
                    # (determines how many models to calculate)
a.md_level = None       # No refinement of model

a.loop.starting_model = 1      # First loop model
a.loop.ending_model   = 4      # Last loop model
a.loop.md_level       = refine.fast # Loop model refinement level

a.make()                # do homology modelling

```

After generating the standard model(s), a number of loop models are generated for each model, from `loopmodel.starting_model` to `loopmodel.loop.ending_model`. Each loop model is written out with the `.BL` extension. See section 6.5 for more information.

The `loopmodel` class can also be used to refine a region of an existing PDB file, without comparative modeling, as in the example below. Note that this example also redefines the `loopmodel.select_loop_atoms` routine. This is necessary in this case, as the default selection selects all gaps in the alignment for refinement, and in this case no alignment is available. You can also redefine this routine in the above case if you want to select a different region for optimization.

Example: `examples/automodel/loop.py`

```

# Loop refinement of an existing model

from modeller.automodel import *

log.verbose()
env = environ()

# directories for input atom files
env.io.atom_files_directory = './:~/atom_files'

# Create a new class based on 'loopmodel' so that we can redefine
# select_loop_atoms (necessary)
class myloop(loopmodel):
    # This routine picks the residues to be refined by loop modeling
    def select_loop_atoms(self):
        # 4 residue insertion (1st loop)
        self.pick_atoms(selection_segment=('19:', '28:'),
                        selection_status='INITIALIZE')

m = myloop(env,
            inimodel='1fdx.B99990001.pdb', # initial model of the target
            sequence='1fdx')              # code of the target

m.loop.starting_model= 20      # index of the first loop model
m.loop.ending_model  = 23      # index of the last loop model
m.loop.md_level = refine.very_fast # loop refinement method

m.make()

```


Chapter 3

Comparative modeling class reference

3.1 automodel reference

All of the functions and data members of the `automodel` are listed for reference below. Please note that the `automodel` class is derived from the `model`, so all properties and commands of the `model` class are also available. Please see section 4.6 for more information.

3.1.1 automodel() — prepare to build one or more comparative models

Options:

<code>env = <environ></code>	MODELLER environment
<code>alnfile = <str : 1></code>	file containing template-sequence alignment
<code>knowns = <str : 0></code>	template codes in alignment
<code>sequence = <str : 1></code>	sequence code in alignment
<code>deviation = <float : 1></code>	control amount of randomization
<code>library_schedule = <int : 1></code>	optimization schedule
<code>toplib = <str : 1></code>	topology library to use
<code>parlib = <str : 1></code>	parameter library to use
<code>topology_model = <int : 1></code>	topology submodel
<code>csrfile = <str : 1></code>	user-provided restraints file
<code>infile = <str : 1></code>	user-provided initial model file
<code>assess_methods = function(s)</code>	assessment functions

Description: `alnfile` is required, and specifies the PIR file which contains an alignment between `knowns` (the templates) and `sequence` (the target sequence).

`deviation` controls the amount of randomization done by `randomize.xyz` or `randomize.dihedrals`. (This can also be set after the object is created, by assigning to `'automodel.deviation'`.)

`library_schedule` selects the optimization schedule for model building. See `model.schedule.make()` for more information. (This can also be set after the object is created, by assigning to `'automodel.library_schedule'`.)

If `toplib` and/or `parlib` are specified, the topology and parameters used to build the model(s) are taken from these files; by default, heavy atom-only parameters and topology are used. (If necessary, `topology_model` can also be given to set `topology.submodel`, but in most cases this is read automatically from the topology library, and so does not need to be set.)

If `csrfile` is set, restraints are not constructed, but are instead read from the user-supplied file of the same name. See section 2.2.5 for an example.

If `infile` is set, an initial model is read from the user-supplied file of the same name. See section 2.2.6 for an example.

`assess_methods` allows you to request assessment of the generated models (by default, none is done). You can provide a function, or list of functions, for this purpose, or use one or more of the standard functions provided in the `assess` module — `assess.GA341`, which uses the GA341 method (see `model.assess_ga341()`) or `assess.DOPE`, which uses the DOPE method (see `model.assess_dope()`). (This can also be set after the object is created, by assigning to `'automodel.assess_methods'`.) See section 2.2.2 for an example. Note that only standard models are assessed in this way; if you are also building loop models, see `loop-model.loop.assess_methods`.

See section 2.1 for a general example of using this class.

3.1.2 `automodel.md_level` — control the model refinement level

Description: This allows the degree of refinement of the models to be adjusted. You can define your own function for this purpose, set it to the special Python value `None` (in which case no additional refinement is done) or use one of the predefined functions in the `refine` module — `refine.very_fast`, `refine.fast`, `refine.slow`, `refine.very_slow` or `refine.slow_large` — which perform different amounts of MD annealing.

3.1.3 `automodel.outputs` — all output data for generated models

Description: This is produced after `automodel.make()` is finished. It contains filenames and model scores for every generated model. This information is provided for your own post-processing (*e.g.*, ranking and further refinement of the models) although a summary of it is printed at the end of the model run.

3.1.4 `automodel.rand_method` — control initial model randomization

Description: This is used to randomize the initial model before producing each final model. (If set to `None` then no randomization is done, and every model will be the same.) You can provide your own function for this purpose, or use one of the functions in the `randomize` module — `randomize.xyz` to randomize all coordinates, or `randomize.dihedrals` to randomize dihedral angles.

3.1.5 `automodel.generate_method` — control initial model generation

Description: This is used to build the initial model. It is usually set to `generate.transfer_xyz`, which builds the model based on the templates, but you can also set it to `generate.generate_xyz` to build it purely from the internal coordinates, or to `generate.read_xyz` to read it from a file (see section 2.2.6 for the easiest way to do this).

3.1.6 `automodel.max_var_iterations` — select length of optimizations

Description: This is used to set `max_iterations` for every call to `model.optimize()`. Smaller numbers may lead to more rapid (but less accurate) model building.

3.1.7 `automodel.repeat_optimization` — number of times to repeat optimization

Description: The entire optimization cycle is repeated this many times; increase the value from 1 (the default) to request more thorough optimization.

3.1.8 automodel.max_molpdf — objective function cutoff

Description: Optimization of each model is automatically aborted (continuing with the next model, if any) if the objective function becomes larger than this value.

3.1.9 automodel.initial_malign3d — initial template alignment

Description: If set to True, then an initial structural alignment of all templates is done.

3.1.10 automodel.starting_model — first model to build

Description: This determines the number of the first model to build; models are built from `automodel.starting_model` through to `automodel.ending_model`.

3.1.11 automodel.ending_model — last model to build

Description: This determines the number of the last model to build; see `automodel.starting_model`.

3.1.12 automodel.final_malign3d — final template-model alignment

Description: If set to True, then all of the generated models (and loop models, if using `loopmodel`) are fit to the templates, and written out with the `_fit.pdb` extension.

3.1.13 automodel.write_intermediates — write intermediate files during optimization

Description: If set to True, then PDB files are written out during the optimization (containing intermediate, partially optimized coordinates) in addition to the final model file(s).

3.1.14 automodel.select_atoms() — select region for optimization

Description: By default, this selects all atoms in the system. Only the selected atoms are optimized in model building, so you can redefine this routine to select the region of interest, if so desired. See section 2.2.4 for an example.

3.1.15 automodel.auto_align() — generate an automatic initial alignment

Description: This generates an initial alignment between the templates and the target sequence. See section 2.2.8 for an example.

3.1.16 automodel.very_fast() — request rapid optimization

Description: This sets parameters to request very fast optimization of the model(s). See section 2.2.2 for an example.

3.1.17 automodel.make() — build all models

Options:

`exit_stage = <int : 1>`

early exit flag

Description: You should call this command after creating an `automodel` object and setting any desired parameters, to then go ahead and build all models.

If `exit_stage` is 2, then this routine exits after generating the initial model; no optimized models are built. If it is 1, then it also creates the restraints file before exiting. If it is 0 (the default) then the full comparative modeling procedure is followed.

3.1.18 `automodel.cluster()` — cluster all built models

Options:

`cluster_cut = <float : 1>` cluster cutoff distance

Description: This can be called after model building is complete. It clusters all of the output models, and outputs an averaged cluster structure, both optimized (in the file `cluster.opt`) and unoptimized (in `cluster.ini`). `cluster_cut` gives the cluster cutoff distance, as used in `model.transfer_xyz()`.

3.1.19 `automodel.special_restraints()` — add additional restraints

Options:

`aln = <alignment>` template-target alignment

Description: By default, this routine does nothing. However, you can redefine it to add additional user-defined restraints to those calculated by `automodel`. The routine is passed `aln`, which is the alignment between the templates and the target sequence. See section 2.2.7 for an example.

3.1.20 `automodel.special_patches()` — add additional patches to the topology

Description: This routine, which is usually empty, can be redefined by the user to add additional patches to the model topology (for example, user-defined disulfides).

3.1.21 `automodel.user_after_single_model()` — analyze or refine each model

Description: This routine is called after building each model, before the output PDB file is written. It usually does nothing, but can be redefined by the user to perform analysis or additional refinement on each model.

3.2 `allhmodel` reference

The `allhmodel` class is derived from `automodel`, so all properties and commands of both the `automodel` and `model` classes are available in addition to those listed below.

3.2.1 `allhmodel()` — prepare to build all-hydrogen models

Options:

<code>env = <environ></code>	MODELLER environment
<code>alnfile = <str : 1></code>	file containing template-sequence alignment
<code>knowns = <str : 0></code>	template codes in alignment
<code>sequence = <str : 1></code>	sequence code in alignment
<code>deviation = <float : 1></code>	control amount of randomization
<code>library_schedule = <int : 1></code>	optimization schedule

<code>parlib = <str : 1></code>	parameter library to use
<code>csrfile = <str : 1></code>	user-provided restraints file
<code>infile = <str : 1></code>	user-provided initial model file
<code>assess_methods = function(s)</code>	assessment functions

Description: This creates a new object for building all-hydrogen models. All of the arguments are the same as those for `automodel()`.

See section 2.2.3 for an example.

3.3 loopmodel reference

The `loopmodel` class is derived from `automodel`, so all properties and commands of both the `automodel` and `model` classes are available in addition to those listed below.

3.3.1 loopmodel() — prepare to build models with loop refinement

Options:

<code>env = <environ></code>	MODELLER environment
<code>alnfile = <str : 1></code>	file containing template-sequence alignment
<code>knowns = <str : 0></code>	template codes in alignment
<code>sequence = <str : 1></code>	sequence code in alignment
<code>inimodel = <str : 1></code>	initial model file
<code>deviation = <float : 1></code>	control amount of randomization
<code>library_schedule = <int : 1></code>	optimization schedule
<code>toplib = <str : 1></code>	topology library to use
<code>parlib = <str : 1></code>	parameter library to use
<code>topology_model = <int : 1></code>	topology submodel
<code>csrfile = <str : 1></code>	user-provided restraints file
<code>infile = <str : 1></code>	user-provided initial model file
<code>assess_methods = function(s)</code>	assessment functions
<code>loop_assess_methods = function(s)</code>	assessment functions

Description: This creates a new object for loop modeling. It can either build standard comparative models (in identical fashion to the `automodel` class) and then refine each of them, in which case you should set the `alnfile` and `knowns` arguments appropriately (see the `automodel()` documentation) or it can refine a given region of a PDB file, in which case you should set `inimodel` to the name of the PDB file instead. In both cases, `sequence` identifies the code of the target sequence.

All other arguments are the same as those for `automodel()`, with the exception of those below:

`loop_assess_methods` is the analog of `automodel.assess_methods` for loop modeling, and allows you to request assessment of the generated loop models. (This can also be set after the object is created, by assigning to `'loopmodel.loop.assess_methods'`.)

See section 2.2.9 for examples.

3.3.2 loopmodel.loop.md_level — control the loop model refinement level

Description: This is the analog of `automodel.md_level` for loop modeling, and allows the loop model refinement to be customized.

3.3.3 `loopmodel.loop.starting_model` — first loop model to build

Description: This is the analog of `automodel.starting_model` and determines the number of the first loop model to build for each regular model.

3.3.4 `loopmodel.loop.ending_model` — last loop model to build

Description: This is the analog of `automodel.ending_model` and determines the number of the last loop model to build for each regular model.

3.3.5 `loopmodel.loop.outputs` — all output data for generated loop models

Description: This is the analog of `automodel.outputs` for loop modeling; it contains filenames and model scores for every generated loop model.

3.3.6 `loopmodel.select_loop_atoms()` — select region for loop optimization

Description: By default, this selects all atoms near gaps in the alignment for loop optimization. You should redefine this routine if you do not have an alignment, or you wish to set a different region for loop optimization. See section 2.2.9 for an example.

Chapter 4

MODELLER command reference

Sections in this Chapter describe technical aspects of MODELLER. They include:

- miscellaneous rules and features of MODELLER (Section 4.1);
- dealing with stereochemical parameters and molecular topology (Section 4.5);
- handling of atomic coordinates (Section 4.6);
- comparing and searching of sequences and structures (Section 4.9);
- calculating spatial restraints (Section 4.7),

4.1 Miscellaneous rules and features of MODELLER

This Section describes several features of the program, including file naming conventions, various file types, and the control of the amount of output.

4.1.1 MODELLER system

One of the main aims of MODELLER is to allow for flexible exploration of various modeling protocols to facilitate the development of better modeling methods. MODELLER can be seen as an interpreted language that is specialized for modeling of protein 3D structure.

See section 1.7.2 for basic information on writing and running scripts.

4.1.2 Controlling breakpoints and the amount of output

Some errors are recoverable. For those errors, a Python `ModellerError` exception is raised. It is then up to your script to deal sensibly with the failure of the preceding command using a standard `'except'` clause. For example, this flexibility allows derivation of multiple models and searching for many sequences, even if some cases abort due to convergence problems.

There are five kinds of messages that MODELLER writes to the `log` file: long output from the MODELLER commands, short notes to do with the execution of the program (files opened, *etc.*), warnings identified by `'_W>'`, errors identified by `'_E>'`, and the messages about the status of dynamic memory allocation. To control how much of this output is displayed, use the `log` object; for more information, see Section 4.14.

4.1.3 File naming

There are several filename generating mechanisms that facilitate file handling. Not all of them apply to all file types.

Environment variables

There can be UNIX shell environment variables in any input or output filename. The environment variables have to be in the format `${VARIABLE}` or `$(VARIABLE)`. Also, two predefined macros are available for string variables:

- `${LIB}` is expanded into the `$LIB_MODELLE` variable defined in `modlib/libs.lib` (equal to `$MODINSTALL8v1/modlib`);
- `${JOB}` is expanded into the root of the script filename, or `'(stdin)'` if instructions are being read from standard input;

Output directory prefix

For all output filenames, except for those that start with `'/'`, the full output filename is obtained by prefixing the filename with `output_directory`.

Coordinate files and derivative data

When accessing an atom file, a specified filename is tried first. If this is unsuccessful, MODELLER automatically expands the original filename by adding extensions `'.Z'`, `'.gz'`, or `'.bz2'`. This allows it to detect atom files compressed with the UNIX `compress`, `gzip`, or `bzip2` commands. If the compressed file exists, MODELLER automatically uncompresses it, reads it, and puts it back into the original state after the reading is finished (note that the `gzip` and `bzip2` programs must be installed on your system in order for this to work). If the specified file is still not found, the extensions `'.atm'`, `'.pdb'`, `'.ent'`, and `'.crd'` are tried in this order, with and without the compressed extensions, then also with the `'pdb'` prefix. This search for the atom file is repeated through all the directories in `io_data.atom_files.directory` (directories are separated by `':'`), unless the input atom filename starts with `'/'`, in which case `io_data.atom_files.directory` is neglected. Finally, if still unsuccessful and the file specified by the environment variable `$PDBENT` exists, the coordinate filename (*e.g.*, the 4 character PDB code) is matched to the list of the full PDB filenames in `$PDBENT` (compressed and uncompressed). For example, `$PDBENT` file may be:

```
/disk2/pdb/pdb.pdb.bnl.gov/all_entries/uncompressed_files/pdb1ema.ent
/disk2/pdb/pdb.pdb.bnl.gov/all_entries/uncompressed_files/pdb1hbp.ent
/disk2/pdb/pdb.pdb.bnl.gov/all_entries/uncompressed_files/pdb1gpy.ent
/disk2/pdb/pdb.pdb.bnl.gov/all_entries/uncompressed_files/pdb6gpb.ent
/disk2/pdb/pdb.pdb.bnl.gov/all_entries/uncompressed_files/pdb1fia.ent
etc.
```

Any derivative data that MODELLER may need, including residue solvent accessibilities, hydrogen bonding information, dihedral angles, residue neighbors, *etc.*, are calculated on demand from the atomic coordinates. The most time consuming operation is calculating solvent accessibility, but even this calculation takes less than 1 sec for a 200 residue protein on a Pentium III workstation.

MODELLER stores the filenames of coordinate sets in the alignment arrays. These arrays are used by `alignment.compare_structures()`, `model.restraints.make()`, `alignment.malign3d()`, `alignment.align2d()`, and several other commands. If these filenames do not change when the structures are needed for the second time, the coordinate files are not re-read because they should already be in memory. This creates a problem only when the contents of a structure file changes since it was last read during the current job.

4.1.4 File types

MODELLER uses a number of standard filename extensions to indicate the type of data stored in a file (Table 4.1). The extensions are generally not mandatory, only very helpful.

4.1.5 Key for command descriptions

Most commands in MODELLER take one or more keyword arguments. The valid keywords are listed, one per line, after each command. For convenience, many of these keywords take a default value if you not specify them. These default values themselves are stored in the `environ` class, so that you can specify global defaults (see section 4.2).

MODELLER will report an error if you do not specify a keyword which does not have a default.

Each keyword argument must be of a specific type, either a MODELLER class object (*e.g.*, a `model` or `alignment`) or one of the standard types shown below:

Extension	Description
.top	TOP script with instructions for a MODELLER job
.log	log output produced by a MODELLER run
.ali	alignment or sequences in the PIR format
.pap	alignment or sequences in the PAP format
.aln	alignment or sequences in the QUANTA format
.aln	alignment or sequences in the INSIGHTII format
.seq, .chn	sequence(s) in the PIR alignment format
.cod	list of sequence codes
.grp	list of families in PDB
.atm, .pdb, .ent	atom coordinates in the PDB or GRASP format
.crd	atom coordinates in the CHARMM format
fit.pdb	fitted protein structures in the PDB format
.ini	initial MODELLER model
.B*	MODELLER model in the PDB format
.D*	the progress of optimization
.BL*	MODELLER model in the PDB format, in loop modeling
.DL*	the progress of optimization, in loop modeling
.IL*	initial MODELLER model, in loop modeling
.V*	violations profile
.E*	energy profile
.rsr	restraints in MODELLER or USER format
.sch	schedule file for the variable target function optimization
.mat	matrix of pairwise protein distances from an alignment
.mat	matrix of pairwise residue type-residue type distance scores
.sim.mat	matrix of pairwise residue type-residue type similarity scores
.lib	various MODELLER libraries
.psa	residue solvent accessibilities
.sol	atomic solvent accessibilities
.ngh	residue neighbors
.dih	mainchain and sidechain dihedral angles
.ssm	secondary structure assignment
.var	sequence variability profile from multiple alignment
.asgl	data for plotting by ASGL

Table 4.1: *List of file types.*

<code><int : 1></code>	an integer variable or constant
<code><float : 1></code>	a floating-point variable or constant
<code><str : 1></code>	a string variable or constant
<code><bool : 1></code>	a boolean variable or constant
<code><int : 0></code>	a tuple of any length with elements <code><int : 1></code>
<code><int : N></code>	a tuple of N elements <code><int : 1></code>
<i>etc.</i>	the same for float, str, and bool types

4.2 The environ class: MODELLER environment

The `environ` class contains most information about the MODELLER environment, in particular the default values for many functions (see section 4.1.5 for more information). Usually it is the first class to be used in a MODELLER script, as it provides methods to create the other main classes. In addition, some miscellaneous commands are also provided as methods of the `environ` class.

4.2.1 `environ()` — create a new MODELLER environment

Options:

<code>rand_seed = <int : 1></code>	random seed
<code>restyp_lib_file = <str : 1></code>	residue type library file

Description: This creates a new `environ` object. `rand_seed` is used to seed the random number generator used throughout MODELLER, and should be set to a negative integer between -2 and -50000 if you do not want to use the default value. `restyp_lib_file` specifies the file to read the residue type library from; if unspecified, the default (`'restyp.lib'`) file is used. This file contains the mapping between one-letter residue types and CHARMM and PDB names; see the FAQ Section 1.8, Question 10 for the format of this file.

You can assign the new `environ` object to the Python variable `'env'` with the following:

```
env = environ()
```

You can release the object from memory when you no longer need it in standard Python fashion, either by an explicit `del(env)` or by reassigning `env` to some other object.

You can set a MODELLER default value by assigning to this `env` object; *e.g.*, to set the default sequence/sequence alignment penalties, you can use:

```
env.gap_penalties_1d = (-900., -50.)
```

When you create new MODELLER objects (such as `model` or `alignment` objects) they require an `environ` object, which they use for their own default values. Note that each object gets a *copy* of the environment, so it is not affected by any changes you make to the global environment *after* its creation. You can, however, modify the object's own environment directly, by assigning to its `.env` member:

```
env = environ()
env.cluster_cut = 2.0 # New objects will use cluster_cut = 2.0 by default
mdl = model(env) # This 'mdl' object gets cluster_cut = 2.0
mdl.env.cluster_cut = 1.0 # cluster_cut is now 1.0, but only for 'mdl'
```

If in doubt, set anything you need to set within `environ` before you create any objects, or always specify parameters explicitly rather than relying on defaults. (You can also turn on verbose logging with `log.verbose()` to see which defaults are being used in each case.)

4.2.2 `environ.io` — default input parameters

Description: This is an `io_data` object, which is used as the default by all routines which take an `io` argument. See section 4.4.

4.2.3 `environ.edat` — default objective function parameters

Description: This is an `energy_data` object, which is used as the default by all routines which take an `edat` argument. See section 4.3.

4.2.4 `environ.libs` — MODELLER libraries

Description: This is a `libraries` object, which contains all of the MODELLER topology and parameter libraries. See section 4.5.

4.2.5 `environ.dendrogram()` — clustering

Options:

<code>cluster_cut = <float : 1></code>	1.0	definition of a cluster
<code>matrix_file = <str : 1></code>	'family.mat'	the filename of the pairwise distance matrix

Description: This command calculates a clustering tree from the input matrix of pairwise distances. This matrix must be in the PHYLIP format and can be produced by the `alignment.id_table()`, `alignment.compare_sequences()`, or `alignment.compare_structures()` commands. The weighted pair-group average clustering method is used.

The tree is written to the `log` file.

This command is useful for deciding about which known 3D structures are to be used as templates for comparative modeling.

Example: See `alignment.id_table()` command.

4.2.6 `environ.principal_components()` — clustering

Options:

<code>matrix_file = <str : 1></code>	'family.mat'	the filename of the pairwise distance matrix
<code>file = <str : 1></code>	'default'	output file

Description: This command calculates principal components clustering for the input matrix of pairwise distances. This matrix must be in the PHYLIP format and can be produced by the `alignment.id_table()`, `alignment.compare_sequences()`, or `alignment.compare_structures()` commands.

The projected coordinates p and q are written to file `file`. The output file can be used with ASGL to produce a principal components plot.

This command is useful for deciding about which known 3D structures are to be used as templates for comparative modeling.

Example: See `alignment.id_table()` command.

4.2.7 `environ.system()` — execute system command

Options:

<code>command = <str : 1></code>	'nothing'	Unix or DOS command
--	-----------	---------------------

Description: This command executes the specified operating system command, for example 'rm' or 'ls' on a Unix system, or 'dir' on a Windows machine. This should be avoided in portable scripts, precisely because the available commands differ between operating systems.

4.3 The energy_data class: objective function parameters

The `energy_data` class stores objective function parameters.

4.3.1 `energy_data()` — create a new set of objective function parameters

Description: This creates a new `energy_data` object. The object will have the default parameters. You can, however, specify any of these parameters when you create the object:

```
edat = energy_data(contact_shell=7.0)
```

Alternatively, you can set parameters in an existing object:

```
edat.contact_shell = 7.0
```

Many commands use `energy_data` objects. However, for convenience, the `environ` class also contains an `energy_data` object, as `environ.edat`. This is used as the default if you do not give an `energy_data` parameter, so you can set this to change the objective function used by all functions:

```
env = environ()
env.edat.contact_shell = 7.0
```

4.3.2 `energy_data.contact_shell` — nonbond distance cutoff

Description: This defines the maximal distance between atoms that flags a non-bonded atom pair. Such pairs are stored in the list of non-bonded atom pairs. Only those non-bonded pairs that are sufficiently close to each other will result in an actual non-bonded restraint.

If undefined (`-999`), the default value is the maximum of the three possibilities: twice the radius of the largest atom multiplied by `energy_data.radii_factor` (in the case of the all non-hydrogen atoms model, this is 3.2 Å); `energy_data.lennard_jones_switch[1]`; or `energy_data.coulomb_switch[1]`. Only those values of the three possibilities are compared that have the corresponding `energy_data.dynamic_sphere`, `energy_data.dynamic_lennard`, or `energy_data.dynamic_coulomb` set to `True`.

The best value for `energy_data.contact_shell` must be found in combination with `energy_data.update_dynamic` (see also below). Good values are 4Å for `energy_data.contact_shell` and 0.39Å for `energy_data.update_dynamic` when no Lennard-Jones and Coulomb terms are used; if `energy_data.contact_shell` is larger, there would be many pairs in the non-bonded pairs list which would slow down the evaluation of the molecular pdf. If it is too small, however, the increased frequency of the pair list recalculation may slow down the optimization.

This distance is also used for the calculation of atomic density; see section 6.3.1.

4.3.3 `energy_data.update_dynamic` — nonbond recalculation threshold

Description: This sets the cumulative maximal atomic shift during optimization that triggers recalculation of the list of atom–atom non-bonded pairs. It should be set in combination with `energy_data.contact_shell`.

For soft-sphere overlap, to be absolutely sure that no unaccounted contacts occur, `energy_data.update_dynamic` has to be equal to $(\text{energy_data.contact_shell} - \text{maximal_overlap_distance}) / 2$. `maximal_overlap_distance` is equal to the diameter of the largest atom in the model; it is 3.2 Å in the case of the all non-hydrogen atoms model. This distance is the `energy_data.contact_shell` value if a default is requested. Factor 2 comes from the fact that the moves of both atoms can reduce the distance between them.

4.3.4 `energy_data.dynamic_sphere` — calculate soft-sphere overlap restraints

Description: If set to `True`, the dynamic soft-sphere overlap restraints are calculated. Note that they are only calculated if the scaled standard deviation of the soft-sphere overlap restraints is greater than zero. It is simpler not to pre-calculate any soft-sphere overlap restraints and to use the dynamically generated restraints alone, although this may be slower.

4.3.5 `energy_data.dynamic_lennard` — calculate Lennard-Jones restraints

Description: If set to `True`, dynamic Lennard-Jones restraints are calculated, using equation 6.63.

4.3.6 `energy_data.dynamic_coulomb` — calculate Coulomb restraints

Description: If set to `True`, dynamic Coulomb (electrostatic) restraints are calculated, using equation 6.60.

4.3.7 `energy_data.dynamic_modeller` — calculate non-bonded spline restraints

Description: If set to `True`, dynamic MODELLER non-bonded spline restraints are calculated.

4.3.8 `energy_data.excl_local` — exclude certain local pairs of atoms

Description: This specifies whether or not the atoms in a chemical bond, chemical angle, dihedral/improper angle, and in the excluded pairs list respectively are considered in the construction of the non-bonded atom pairs list. This is especially useful when simplified protein representations are used; *e.g.*, when non-bonded restraints need to be used on $C_{\alpha i} - C_{\alpha i+2}$ terms.

4.3.9 `energy_data.radii_factor` — scale atomic radii

Description: This is the scaling factor for the atom radii as read from the `$RADII_LIB` library file. The scaled radii are used only for the calculation of violations of the soft-sphere overlap restraints and by `model.write_data()`. Note that which radii to first read from the library file are determined by the `topology.submodel` variable.

4.3.10 `energy_data.lennard_jones_switch` — Lennard-Jones switching parameters

Description: These are the parameters f_1 and f_2 to the Lennard-Jones switching function, which smoothes the potential down to zero; see equation 6.63.

4.3.11 `energy_data.coulomb_switch` — Coulomb switching parameters

Description: These are the parameters f_1 and f_2 to the Coulomb switching function, which smoothes the potential down to zero; see equation 6.60.

4.3.12 `energy_data.relative_dielectric` — relative dielectric

Description: This sets the relative dielectric ϵ_r , used in the calculation of the Coulomb energy (equation 6.60).

4.3.13 `energy_data.covalent_cys` — use disulfide bridges in residue distance

Description: This modulates the effect of `residue_span_range`; if `True`, the disulfide bridges are taken into account when calculating the residue index difference between two atoms (*i.e.*, disulfides make some atom pairs closer in sequence). `energy_data.covalent_cys = True` is slow and only has an effect when certain statistical non-bonded potentials are used (*i.e.*, `energy_data.dynamic_modeller` is `True` and the non-bonded library has been derived considering the disulfide effect). Thus, it should generally be set to `False`.

4.3.14 `energy_data.nonbonded_sel_atoms` — control interaction with picked atoms

Description: When only a subset of all atoms is selected by the `model.pick_atoms()` or `model.pick_hot_atoms()` commands (set 1), this variable controls the interaction between the picked atoms and the rest of the system. When it is 2 (default), the non-bonded pairs will contain only selected atoms

(set 1). This means that the optimized atoms will not “feel” the rest of the protein through the non-bonded terms at all.

If `energy_data.nonbonded_sel_atoms` is 1, only one of the atoms in the non-bonded pair has to be a selected atom. This means that the selected region feels the rest of the system through the non-bonded terms, at the expense of longer CPU times. When all atoms are selected, this variable has no effect.

4.3.15 `energy_data.nlogn_use` — select non-bond list generation algorithm

Description: Before calculating dynamic non-bonded restraints, MODELLER determines which of the several routines is most appropriate and efficient for calculating the non-bonded atom pairs list. The user can use this variable to influence the selection, of either a straightforward $\mathcal{O}(n^2)$ search or a cell-based algorithm which has $n \log n$ dependency of CPU time *versus* size n . The latter algorithm is used when the maximal difference in residue indices of the atoms in the current dynamic restraints is larger than `energy_data.nlogn_use` or when the box size for this algorithm would have to be larger than 8Å.

4.4 The `io_data` class: coordinate file input parameters

The `io_data` class stores information useful in reading coordinate files. This is used both for reading models, and for reading coordinate files used for templates.

4.4.1 `io_data()` — create a new input parameters object

Description: This creates a new `io_data` object, with default parameters. You can, however, specify any of these parameters when you create the object:

```
io = io_data(hetatm=True)
```

Alternatively, you can set parameters in an existing object:

```
io.hetatm = True
```

Many commands use `io_data` objects. However, for convenience, the `environ` class also contains an `io_data` object, as `environ.io`. This is used as the default if you do not give an `io_data` parameter, so you can set this to change the input configuration used by all functions:

```
env = environ()
env.io.hetatm = True
```

Please note that `io_data.hetatm`, `io_data.hydrogen` and `io_data.water` are ignored when reading UHBD files. When reading CHARMM files, `io_data.hetatm` is ignored.

4.4.2 `io_data.hetatm` — whether to read HETATM records

Description: If set to `True`, then 'HETATM' records are read from PDB files. (By default, only 'ATOM' records are read.) Note that you will need to turn this on when using BLK residues, or if you want to use a PDB file containing ligands.

4.4.3 `io_data.hydrogen` — whether to read hydrogen atoms

Description: If set to `True`, then hydrogen atoms are read from PDB or CHARMM files. (By default, only heavy atoms are read.) You will need to turn this on when building all-atom models, although note that the `allhmodel` class does this for you automatically.

4.4.4 `io_data.water` — whether to read water molecules

Description: If set to `True`, then water molecules are read. (For PDB files, this is regardless of whether they are in 'ATOM' or 'HETATM' records.) Ordinarily, they are ignored. (See '`LIB/restyp.lib`' for the definition of a water molecule used by MODELLER.)

4.4.5 `io_data.atom_files_directory` — search path for coordinate files

Description: This can contain one or more directories, separated by colons (:) in which to look for coordinate files. If empty, only the current directory is searched.

4.5 The libraries class: stereochemical parameters and molecular topology

All molecular modeling programs generally need to know what are the atoms in all residue types, what are the atom pairs that are covalently bonded to each other (*i.e.*, molecular topology), and what are the ideal bond lengths, angles, dihedral angles, and improper dihedral angles (*i.e.*, internal coordinates and stereochemical restraints). This information is stored in the residue topology and parameter libraries, which are manipulated by MODELLER scripts using the `libraries` class.

Note that unlike other MODELLER objects, in this release of MODELLER there can be only one set of topology and parameter libraries. Thus, you will not need to create a `libraries` object directly, but instead deal with the `environ.libs` variable instead.

For commands dealing with generating, patching, and mutating molecular topology, see the `model` class reference, in section 4.6.

4.5.1 Modeling residues with non-existing or incomplete entries in the topology and parameter libraries

Defining new residue types is generally one of the more painful areas in developing and using a molecular modeling program. MODELLER has two quick-and-dirty solutions described in the next two sections that are often sufficient for comparative modeling involving new residue types. On the other hand, if you are willing to spend some time and define a new entry or complete an incomplete entry in the residue topology or parameter libraries, see the FAQ Section 1.8, Question 10.

Residues with defined topology, but with missing parameters

The parameter library is used by the `model.restraints.make()` command to construct bond, angle, dihedral angle, improper dihedral angle, and non-bonded Lennard-Jones restraints. If some parameters for these restraints are missing, they are guessed on the fly from the current Cartesian coordinates of the MODEL. Thus, when there are missing parameters, the MODEL coordinates must be defined before calling `model.restraints.make()`. The coordinates can be defined by the `model.build()` command (from the IC entries in the residue topology library), by the `model.read()` command (from an existing coordinate file for MODEL), or by the `model.transfer_xyz()` command (from template coordinate files aligned with MODEL). The bonds, angles, and improper dihedral angles are restrained by a harmonic potential with the mean equal to the value in the current structure and a force constant typical for chemical bonds, angles, and improper dihedral angles, respectively. The dihedral angles are restrained by a tri-modal cosine term with the mean equal to the angle in the current structure. A message detailing MODELLER's improvization is written to the `log` file.

Block (BLK) residues with undefined topology and parameters

The second relatively easy way of dealing with missing entries in the residue topology and/or parameters libraries is to use a “block” residue. These residues are restrained more or less as rigid bodies to the conformation of the equivalent residue(s) in the template(s). No chemical information is used. The template residues can themselves be defined as block residues. The symbol for the block residues is ‘BLK’ in the four- and three-letter codes and ‘.’ in the single-letter code. The atoms in a BLK residue include all uniquely named atoms from the equivalent residues in all the templates. The atom type of all BLK atoms is the CHARMM type ‘`undf`’. The IUPAC atom names (as opposed to the atom types) are the same as in the templates. The ‘`undf`’ atom type for all BLK atoms facilitates using the `model.pick_atoms()` command for generating restraints on the ‘BLK’ residues.

The ‘`undf`’ atoms are treated differently from the other atoms during preparation of dynamic restraints: No pairs of intra-BLK atoms are put on the dynamic non-bonded list. Only the “inter-BLK” atom pairs and “BLK–other” atom pairs are considered for the dynamic non-bonded restraints. The radius of all block atoms is obtained from the `$RADII.LIB` library using the block atom names (as written out to a PDB file), not the ‘`undf`’ atom type. All intra-BLK and inter-residue BLK restraints other than the non-bonded restraints have to be derived separately and explicitly by `model.restraints.make()` command using `restraint_type = 'distance'`. See script


```

# Example for: topology.make(), topology.write()

# This creates a topology library for heavy atoms from the
# CHARMM all-atom topology library:

env = environ()

tpl = env.libs.topology
# Read CHARMM all-atom topology library:
tpl.read(file='${LIB}/top.lib')

# Keep only heavy atoms (TOPOLOGY_MODEL = 3)
tpl.make(submodel=3)

# Write the resulting topology library to a new file:
tpl.write(file='top_heav.lib')

```

4.5.11 topology.submodel — select topology model type

Description: This is used to select the type of topology model currently in use; see **topology.make()** for more information. Note that it is not usually necessary to explicitly set **topology.submodel**, since it is set for you automatically when you read a topology file (all of the standard MODELLER topology files include a `'*MODELLER topology.submodel'` header from which this information is derived). See **topology.read()**.

4.5.12 topology.write() — write residue topology library

Options:

<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
<code>output_directory = <str : 1></code>	<code>''</code>	output directory

Description: This command writes a residue topology library to the specified file. It is usually used after **topology.make()**.

Example: See **topology.make()** command.

4.6 The model class: handling of atomic coordinates, and model building

The `model` class holds all information about a 3D model (such as its Cartesian coordinates, topology, and optimization information). It also provides methods for reading, writing and transforming the model.

4.6.1 `model()` — create a new 3D model

Options:

`env = <environ>` MODELLER environment

Description: This creates a new `model` object. If used with no parameters, the new model is empty (*i.e.*, it contains no atoms). However, if any keyword arguments are given, they are passed to the `model.read()` function to read in an initial model. See the `model.to_iupac()` example.

4.6.2 `model.seq_id` — sequence identity between the model and templates

Description: This is set by MODELLER during model building when `model.transfer_xyz()` is called. It is also read from or written to MODELLER-produced PDB files in a REMARK header.

You should set `model.seq_id` explicitly if you want to carry out GA341 assessment on a model not produced by MODELLER, as the sequence identity is not known in this case. See `model.assess_ga341()`.

4.6.3 `model.restraints` — all restraints which act on the model

Description: This provides the restraints themselves, and methods to manipulate them. See section 4.7 for more information.

4.6.4 `model.read()` — read coordinates for MODEL

Options:

<code>io = <io_data></code>		Options for reading atom files
<code>aln = <alignment></code>		Alignment to use for <code>model.segment</code>
<code>file = <str : 1></code>	<code>'default'</code>	name of the coordinates' file
<code>model.segment = <str : 2></code>	<code>'FIRST:0' 'LAST:'</code>	segment to be read in
<code>model.format = <str : 1></code>	<code>'PDB'</code>	selects input atom file format: <code>'PDB' 'CHARMM' 'UHBD'</code>

Description: This command reads the atomic coordinates, atom names, residue names, residue numbers, isotropic temperature factors and segment specifications for MODEL, assigns residue types, and defines the dihedral angles listed in the \$RESDIH_LIB library. For CHARMM and UHBD file formats, it also reads the atomic charges. However, it does not assign CHARMM and MODELLER atom types, internal coordinates, charges (in the case of the 'PDB' format), or patches (such as disulfides); to make these assignments, which are necessary for almost all energy commands, use `model.generate_topology()`. All real and pseudo atoms are selected. The PDB residue type 'HIS' is assigned the CHARMM residue type 'HSD', which is the neutral His with H on ND1. The PDB types 'ASP' and 'GLU' are assigned the corresponding charged CHARMM residue types, as are 'LYS' and 'ARG'. These conventions are relevant only if electrostatic terms and/or hydrogens are used.

`model.segment` sets the beginning and ending residue identifiers for the contiguous sequence of residues to be read from the PDB file (this option does not work yet for the other file formats). The format of residue identifiers is described in Section 4.9.1. In addition, the following rule applies: If there is no ':' in the first residue specification, the segment specification is taken from the alignment entry with the specified code (in

which case this command must be passed an alignment). Similarly, if there is no ‘:’ in the second residue specification, the PDB filename is taken from the alignment entry with the specified code. The two codes do not have to be the same. For example, `model_segment = '4ape' '4ape'` will take the segment specification and atom filename for entry `4ape` in the alignment.

Example: `examples/commands/read_model.py`

```
# Example for: model.read(), model.write()

# This will read a PDB file and write a CHARMM atom file without atomic charges
# or radii. For assigning charges and radii, see the all_hydrogen.py script.

env = environ()

mdl = model(env)
mdl.read(file='1fas')
mdl.write(file='1fas.crd', model_format='CHARMM')
mdl.write(file='1fas.cif', model_format='MMCIF')
```

4.6.5 `model.write()` — write MODEL

Options:

<code>file = <str : 1></code>	<code>'default'</code>	name of the coordinates' file
<code>output_directory = <str : 1></code>	<code>''</code>	output directory
<code>model_format = <str : 1></code>	<code>'PDB'</code>	selects output atom file type: <code>'PDB' 'CHARMM' 'UHBD' 'GRASP' 'MMCIF'</code>
<code>write_all_atoms = <bool : 1></code>	<code>True</code>	whether to write all atoms, even if unselected
<code>no_ter = <bool : 1></code>	<code>False</code>	whether to not write TER into PDB

Description: This command writes the current MODEL to a file in the selected format. If the file format is `'PDB'`, only the selected atoms are written out when `write_all_atoms = False`; otherwise all atoms are written out. `'MMCIF'` writes out files in the Macromolecular Crystallographic Information File (mmCIF) format.

The `'GRASP'` format is the same as the `'PDB'` format, except that it includes two special lines at the top of the file and the atomic radii and charges in the columns following the Cartesian coordinates of atoms. This format is useful for input to program GRASP, written by Anthony Nicholls in the group of Barry Honig at Columbia University [Nicholls *et al.*, 1991].

Example: See `model.read()` command.

4.6.6 `model.generate_topology()` — generate MODEL topology

Options:

<code>aln = <alignment></code>		Alignment between model and templates
<code>io = <io_data></code>		Options for reading atom files
<code>add_segment = <bool : 1></code>	<code>False</code>	whether to add the new segments to the list of segments
<code>patch_default = <bool : 1></code>	<code>True</code>	whether to do default NTER and CTER patching
<code>sequence = <str : 1></code>	<code>'undefined'</code>	protein code in the alignment whose topology is constructed

Requirements: topology and parameter libraries

Description: This command calculates MODEL's covalent topology (*i.e.*, atomic connectivity) and internal coordinates, and assigns CHARMM atom types, MODELLER atom types for non-bonded spline restraints, atomic charges, and atomic radii.

The sequence of the model to generate is taken from the sequence with code `sequence` in the alignment `aln`.

The new sequence is added to the list of segments of the MODEL if `add_segment` is `True`, otherwise this list is initiated.

A sequence in the alignment can use any non-patching residue listed in the single-character code column of the `$RESTYP_LIB` library (`'modlib/restyp.lib'`). Examples of non-standard residue types include water (`'w'`), zinc (`'z'`), calcium (`'3'`), heme (`'h'`), and many others. Patching residues must not be used here, but with the subsequent `model.patch()` commands. Unrecognized residues are ignored. A special allowed residue type is the chain break `'/'`. This can be used to construct a protein that consists of several chains separated by chain breaks. Chain breaks before a non-standard residue type (there are 23 standard residue types, including `'-'`, `'Asx'` and `'Glx'`) are inserted automatically and do not have to be specified explicitly in the sequence.

The `model.generate_topology()` command generates only the topology of the MODEL, not its Cartesian coordinates; the Cartesian coordinates are assigned by the `model.build()`, `model.transfer_xyz()`, or `model.read()` commands.

In general, the `model.generate_topology()` command has to be executed before any energy commands (`model.energy()`, `model.optimize()`, `model.pick_hot_atoms()`). The reason is that reading the Cartesian coordinates by the `model.read()` command does not generate all the data usually needed for energy evaluation. However, if the order and number of atoms in the input file correspond exactly to the order and number of atoms implied by the restraint atom indices and if you are not using dynamic restraints that rely on non-existing data, such as bond, angle, and dihedral angle lists, atomic charges, radii, Lennard-Jones parameters, MODELLER atom types, or CHARMM atom types (which are used to determine the atomic radii), it is sufficient to do only `model.read()` and omit `model.generate_topology()` before the energy commands. In short, if you use static restraints alone and if the atom file has the atoms in the correct order, you do not have to call `model.generate_topology()` before calculating energy.

The variables `alignment[x].atom_file`, `io` and `topology.submodel` are necessary only when the `'BLOCK'` residues are present in the sequence whose topology is generated. In that case, the template PDB files are read in.

Example: See `model.patch()` command.

4.6.7 `model.patch()` — patch MODEL topology

Options:

<code>residue_ids = <str : 0></code>	<code>''</code>	identifiers of the patched residues
<code>residue_type = <str : 1></code>	<code>'undefined'</code>	patching residue type

Description: This command uses a CHARMM patching residue to patch the topology of the MODEL. CHARMM patch rules are observed.

`residue_type` is the type of the patching residue (PRES entry in the topology library), such as `'DISU'`, `'NTER'`, `'CTER'`, *etc.* You do not have to apply explicitly the N- and C-terminal patches to protein chains because the `'NTER'` and `'CTER'` patches are applied automatically to the appropriate residue types at the termini of each chain at the end of each `model.generate_topology()` command.

`residue_ids` are residue identifiers of the patched residues (Section 4.9.1). The first residue is the patched residue 1, the second residue is the patched residue 2, *etc.*; for example, the `'DISU'` patching residue has two patched Cys residues while the `'ACE'` patching residue has only one patched residue. The order of the residue identifiers here has to match the definition of the patching residue in the topology library.

It is not allowed to patch an already patched residue. Since the N- and C-terminal residues of each chain are automatically patched with the 'NTER' and 'CTER' patching residues, respectively, a user who wants to patch the N- or C-terminal residues with other patches, should turn the default patching off before executing `model.generate_topology()`. This is achieved by setting `patch_default = False`.

Example: `examples/commands/patch.py`

```
# Example for: model.patch(), topology(), parameters.read()

# This will define a CYS-CYS disulfide bond between residues 3 and 22.

env = environ()
env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

# Read the sequence:
code = '1fas'
mdl = model(env, file=code)
# have two copies of the sequence in the alignment, for TRANSFER_XYZ later:
aln = alignment(env)
aln.append_model(mdl, atom_files=code, align_codes=code)
aln.append_model(mdl, atom_files=code+'.ini', align_codes=code+'-ini')
mdl.generate_topology(aln, sequence=code+'-ini')

# Create the disulfide bond:
mdl.patch(residue_type='DISU', residue_ids=('3', '22'))

# Get MODEL's coordinates from the template, using the alignment (1:1 here):
mdl.transfer_xyz(aln)
# Calculate missing coordinates using internal coordinates:
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')

# Create the stereochemical restraints
mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False)

# Calculate the energy to test the disulfide:
mdl.energy()
```

4.6.8 `model.patch_ss_templates()` — guess MODEL disulfides from templates

Options:

<code>aln = <alignment></code>	Alignment between model and templates
<code>io = <io_data></code>	Options for reading atom files

Description: This command defines and patches disulfide bonds in the MODEL using an alignment of the MODEL sequence with one or more template structures. The MODEL sequence has to be the last sequence in the alignment, `aln`. The template structures are all the other proteins in the alignment. All Cys–Cys pairs in the target sequence that are aligned with at least one template disulfide are defined as disulfide bonds themselves. The covalent connectivity is patched accordingly.

This command should be run after `model.generate_topology()` and before `model.restraints.make()` to ensure that the disulfides are restrained properly by the bond length, angle, and dihedral angle restraints and that no SG–SG non-bonded interactions are applied.

The disulfide bond, angle and dihedral angle restraints have their own physical restraint type separate from the other bond, angle and dihedral angle restraints (Table 4.5).

Example: examples/commands/patch_disulfides.py

```
# Example for: model.patch_ss_templates() and model.patch_ss()

# This will patch CYS-CYS disulfide bonds using disulfides in aligned templates:

log.verbose()
env = environ()
env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

# Read the sequence, calculate its topology, and coordinates:
aln = alignment(env, file='toxin.ali', align_codes=('2ctx', '2abx'))
# Superpose the two template structures without changing the alignment.
# This is for TRANSFER_XYZ to work properly. It relies on not reading
# the atom files again before TRANSFER_XYZ.
aln.malign3d(fit=False) # This is for TRANSFER_XYZ to work properly.
# Restore the alignment, and add in the model sequence, 1fas:
aln.clear()
aln.append(file='toxin.ali', align_codes=('2ctx', '2abx', '1fas'))
mdl = model(env)
mdl.generate_topology(aln, sequence='1fas')
mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=True, build_method='INTERNAL_COORDINATES')
mdl.write(file='1fas.noSS')
# Create the disulfide bonds using equivalent disulfide bonds in templates:
mdl.patch_ss_templates(aln)

# Create the stereochemical restraints
mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False)

# Calculate energy to test the disulfide restraints (bonds, angles, dihedrals):
mdl.energy()

mdl.read(file='1fas.noSS')
# Create the disulfide bonds guessing by coordinates
mdl.patch_ss()

# Create the stereochemical restraints
mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False)

# Calculate energy to test the disulfide restraints (bonds, angles, dihedrals):
mdl.energy()
```

4.6.9 model.patch_ss() — guess MODEL disulfides from model structure

Description: This command defines and patches disulfide bonds in MODEL using MODEL's current structure. A disulfide bridge is declared between all pairs of Cys residues whose SG-SG distances are less than 2.5Å. The covalent connectivity is patched accordingly.

This command should be run after **model.read()** and before optimization to ensure that the disulfides are fixed properly and that no SG-SG non-bonded interactions are applied.


```

# Add the mutated sequence to the alignment arrays (it is now the second
# sequence in the alignment):
aln.append_model mdl, align_codes='1fas-1')

# Generate molecular topology for the mutant:
mdl.generate_topology(aln, sequence='1fas-1')

# Transfer all the coordinates you can from the template native structure
# to the mutant (this works even if the order of atoms in the native PDB
# file is not standard):
mdl.transfer_xyz(aln)

# Build the remaining unknown coordinates for the mutant:
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')

# Write the mutant to a file:
mdl.write(file='1fas-1.atm')

```

4.6.11 model.build() — build MODEL coordinates from topology

Options:

<code>initialize_xyz = <bool : 1></code>	<code>True</code>	whether to use IC entries to calculate all coordinates
<code>build_method = <str : 1></code>	<code>'INTERNAL_COORDINATES'</code>	method for building coordinates: 'INTERNAL_COORDINATES' 'ONE_STICK' 'TWO_STICKS' '3D_INTERPOLATION'

Requirements: topology file & parameters file & MODEL topology

Description: This command builds Cartesian coordinates of the MODEL.

If `initialize_xyz` is `True`, all coordinates are built. Otherwise only the undefined coordinates are built. The latter is useful because some coordinates may be undefined after the `model.read()` or `model.transfer_xyz()` command. The undefined coordinates have a value of `-999`. when written to a PDB file.

If `build_method` is `'INTERNAL_COORDINATES'`, the Cartesian coordinates are built from the ideal values of the internal coordinates as obtained from the IC entries in the residue topology library. If an appropriate IC entry does not exist, the ideal value of the internal coordinate is calculated from the corresponding energy term in the parameter library. If some coordinates still cannot be built, they are set to values close to those of the neighboring atoms. If even this fails, they are set randomly.

If `build_method` is `'3D_INTERPOLATION'`, the Cartesian coordinates are built by linearly interpolating between the two defined atoms that span the contiguous undefined segment of atoms. In this mode, both the main-chain and sidechain conformations of all inserted residues are random and distorted. This build-up mode is useful because it may eliminate a knot and minimize the extended nature of the insertion obtained by `build_method = 'INTERNAL_COORDINATES'`. In the end, the coordinates of each of the interpolated atoms are slightly randomized ($\pm 0.2\text{\AA}$) to prevent numerical problems with colinear angles and colinear dihedral angles. If one or both of the spanning atoms are undefined, the `'ONE_STICK'` option (below) is used.

If `build_method` is `'ONE_STICK'`, the Cartesian coordinates are built by “growing” them linearly out of the N-terminal spanning atom (C-terminal atom for the undefined N-terminal), away from the gravity center of all the defined atoms. If there are no spanning atoms, the spanning atom is defined randomly.

If `build_method` is `'TWO_STICK'`, the loop is broken into two equal pieces and the `'ONE_STICK'` algorithm is applied to both halves of the loop separately.

Example: examples/commands/build_model.py

```

# Example for: model.build()
# This will build a model for a given sequence in an extended conformation.

env = environ()
env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

# Read the sequence from a file (does not have to be part of an alignment):
aln = alignment(env, file='toxin.ali', align_codes='1fas')
# Calculate its molecular topology:
mdl = model(env)
mdl.generate_topology(aln, sequence='1fas')
# Calculate its Cartesian coordinates using internal coordinates and
# parameters if necessary:
mdl.build(initialize_xyz=True, build_method='INTERNAL_COORDINATES')

# Write the coordinates to a PDB file:
mdl.write(file='1fas.ini')

```

Example: examples/commands/all_hydrogen.py

```

# Example for: model.generate_topology(), model.build()

# This will read a specified atom file, generate all hydrogen atoms,
# add atomic radii and charges, and write the model to a PDB file in
# the GRASP format. This can be used with GRASP to display electrostatic
# properties without assigning charges and radii in GRASP.

log.verbose()
env = environ()

env.libs.topology.read(file='$(LIB)/top_allh.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

code = '1fas'
mdl = model(env, file=code)

aln = alignment(env)
aln.append_model(mdl, align_codes=code, atom_files=code)
aln.append_model(mdl, align_codes=code+"-ini", atom_files=code+".ini")

mdl.generate_topology(aln, sequence=code)

# Have to patch the topology here to remove sulphhydryl hydrogens:
mdl.patch(residue_type='DISU', residue_ids=('17', '39'))
mdl.patch(residue_type='DISU', residue_ids=('3', '22'))
mdl.patch(residue_type='DISU', residue_ids=('53', '59'))
mdl.patch(residue_type='DISU', residue_ids=('41', '52'))

mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')

```

```
mdl.write(file='1fas.ini1', model_format='GRASP')
mdl.write(file='1fas.ini2', model_format='PDB')
```

4.6.12 `model.unbuild()` — undefine MODEL coordinates

Description: This command undefines all of the Cartesian coordinates of the MODEL.

4.6.13 `model.transfer_xyz()` — copy templates' coordinates to MODEL

Options:

<code>aln = <alignment></code>		Alignment between the model and templates
<code>io = <io_data></code>		Options for reading atom files
<code>cluster_cut = <float : 1></code>	1.0	definition of a cluster
<code>cluster_method = <str : 1></code>	'RMSD'	what distance function to use; 'RMSD' 'MAXIMAL_DISTANCE'

Description: This command transfers coordinates of the equivalent atoms and their isotropic temperature factors from the template structures to MODEL.

The target sequence must be the last protein in the alignment, `aln`, and has to be the same as the MODEL sequence. The template structures are all the other proteins in the alignment.

Before transferring coordinates, the template structures generally have to be explicitly least-squares superposed onto each other. This is most conveniently achieved with the `alignment.malign3d()` command called just before `model.transfer_xyz()`. This is an important difference relative to MODELLER-3, which did not require explicit superposition by the user. Note, however, that the `automodel` class script does this superposition automatically.

If `cluster_cut` is less than 0, the transferred coordinates for a given target atom are the average of the coordinates of all the equivalent template atoms. Otherwise, the transferred coordinates are the average of the templates in the largest cluster of the atoms. This cluster is obtained as follows (it only works when all templates and the target have exactly the same topology): For each residue position separately, calculate the maximal inter-template equivalent atom–atom distances (`cluster_method = 'MAXIMAL_DISTANCE'`) or atomic RMS deviation (`cluster_method = 'RMSD'`) for all template–template comparisons. Use the weighted pair-group average clustering method (the same as in the `environ.dendrogram()` command) to obtain the clustering tree for the given residue position. Find the clusters that contain residues joined above `cluster_cut` angstroms (1Å is a good value). Use the largest cluster in the averaging for the target coordinates. The number of residue positions at which each template contributes to the consensus is written to the log file ('The largest cluster occupancy'). Sometimes the first template contributes many more times than the rest of the templates. This results from having many residue positions where all “clusters” have one template only (the first cluster/template is then picked by default). This artifact can be corrected by specifying a larger `cluster_cut`.

Both kinds of averaging, but especially the cluster averaging, are useful for deriving a consensus model from an ensemble of models of the same sequence. If the consensus model is optimized by the conjugate gradients method, it frequently has a significantly lower value of the objective function than any of the contributing models. Thus, the construction of a consensus model can also be seen as part of an efficient optimization. The reason why consensus construction frequently results in better models is that the consensus model generally picks the best (*i.e.*, most frequent) conformation for the regions that are variable in the individual models, while it is very unlikely that a single model will have optimal conformation in all of the variable regions. The consensus construction may not work when two or more locally optimal conformations are inconsistent with each other (*e.g.*, because of the atom overlaps).

Two atoms are equivalent if they have exactly the same name and are in the equivalent residues. Note that the `$ATMEQV_LIB` library of equivalent residue–residue atom pairs, which is used in the construction of

homology-derived distance restraints, is not used here. The atom names in the target may not correspond to the atom names in the template files. In such a case, if you want to copy the template atoms' coordinates, you have to edit the atom names in the template atom files so that they correspond to the MODELLER atom names (which you can see in the `.ini` atom file). At least for water molecules, this is usually better than letting the optimizer deal with grossly incorrect starting positions.

The atoms with undefined coordinates in MODEL are flagged by setting the coordinates to `-999`. The coordinates of the undefined atoms of the MODEL can be set with the `model.build()` command, which relies on the internal coordinates specified in the residue topology library or on various types of geometric interpolation and extrapolation.

Example: `examples/commands/transfer_xyz.py`

```
# Example for: model.transfer_xyz()

# This will build a model for a given sequence by copying
# coordinates from aligned templates. When the templates
# have the same sequence as the target, this procedure ensures
# that the new model corresponds to the MODELLER topology library.

env = environ()
env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

# Read the sequence and calculate its topology:
aln = alignment(env, file='toxin.ali', align_codes=('2ctx', '1nbt'))
aln.malign3d(fit=False)
aln.append(file='toxin.ali', align_codes='1fas')
mdl = model(env)
mdl.generate_topology(aln, sequence='1fas')
# Assign the average of the equivalent template coordinates to MODEL:
mdl.transfer_xyz(aln)
# Get the remaining undefined coordinates from internal coordinates:
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')

# Write the final MODEL coordinates to a PDB file:
mdl.write(file='1fas.ini')
```

4.6.14 `model.res_num_from()` — residue numbers from MODEL2 to MODEL

Options:

<code>mdl = <model></code>	Model to take residue numbers from
<code>aln = <alignment></code>	Alignment between the two models

Description: This command transfers residue numbers and chain ids from `mdl` to MODEL. It uses the current alignment if present, otherwise a 1:1 correspondence is assumed. `mdl` and MODEL must correspond to the first and second protein in the alignment, `aln`, respectively.

Example: `examples/commands/transfer_res_numb.py`

```
# Example for: model.res_num_from()

# This will transfer residue numbers and chain ids from model2 to model.
```

```

log.level(output=1, notes=1, warnings=1, errors=1, memory=0)
env = environ()

# Read an alignment for the transfer
aln = alignment(env, file='toxin.ali', align_codes=('2ctx', '1fas'))
# Read the template and target models:
mdl2 = model(env, file='2ctx')
mdl = model(env, file='1fas')
# Transfer the residue and chain ids and write out the new MODEL:
mdl.res_num_from(mdl2, aln)
mdl.write(file='1fas.ini')

```

4.6.15 `model.rename_segments()` — rename MODEL segments

Options:

```

segment_ids = <str: 0>          ''          new segment ids
renumber_residues = <int: 0>   starting residue index for renumbering
                                residues

```

Description: This command re-labels residue numbers in each chain (*i.e.*, segment) so that they start with `renumber_residues[iseg]`. In addition, the single character PDB chain id's are also assigned: They are obtained from the corresponding elements of `segment_ids`. Thus, there should be as many elements in `segment_ids` and `renumber_residues` as there are chains in the current MODEL.

Example: `examples/commands/rename_segments.py`

```

# Example for: model.rename_segments()

# This will assign new PDB single-character chain id's to all the chains
# in the input PDB file (here there are two 'chains': protein and the HETATM
# water molecules).

# Read the MODEL with all HETATM and water records (so there are two 'chains'):
env = environ()
env.io.hetatm = True
env.io.water = True

mdl = model(env, file='1fas')
# Assign new segment names and write out the new model:
mdl.rename_segments(segment_ids=('X', 'Y'))
mdl.write(file='1fas.ini')

```

4.6.16 `model.pick_atoms()` — select atoms in MODEL

Options:

```

aln = <alignment>              Used for detecting loop regions
pick_atoms_set = <int: 1>      1          index of the selected atoms set: 1 | 2 |
                                3

```

selection_search = <str : 1>	'SEGMENT'	search method: 'SPHERE' 'SEGMENT' 'SPHERE_SEGMENT'
res_types = <str : 1>	'ALL'	residue type selection: 'ALL' 'HET' 'BLK' 'STD' CHARMM 4-letter codes
atom_types = <str : 1>	'ALL'	atom type selection: 'ALL' 'SDCH' 'MNCH' IUPAC atom names
selection_from = <str : 1>	'ALL'	selecting from: 'ALL' 'SELECTED'
selection_mode = <str : 1>	'ATOM'	selecting what: 'ATOM' 'RESIDUE'
selection_status = <str : 1>	'INITIALIZE'	what to do with selected atoms: 'ADD' 'REMOVE' 'INITIALIZE'
• For selection_search = 'SEGMENT':		
selection_segment = <str : 2>	'' ''	'RES:CHN' ids for the first and last residues in a chain/segment; or 'LOOPS'
gap_extension = <int : 2>	2 1	extend insertions/deletions for that many residues, in PICK_ATOMS
minmax_loop_length = <int : 2>	5 15	minimal/maximal length of a loop in PICK_ATOMS
• For selection_search = 'SPHERE':		
sphere_center = <str : 2>	'undefined' 'undefined'	'#RES1:C' 'ATOM_NAME'
sphere_radius = <float : 1>	10.0	sphere radius for atoms selection
selection_slab = <float : 5>	9999 9999 0 0 0	slab for atoms selection: 'dz1' 'dz2' 'xtrans' 'ytrans' 'ztrans'
• For selection_search = 'SPHERE_SEGMENT':		
selection_segment = <str : 2>	'' ''	'RES:CHN' ids for the first and last residues in a chain/segment; or 'LOOPS'
gap_extension = <int : 2>	2 1	extend insertions/deletions for that many residues, in PICK_ATOMS
sphere_radius = <float : 1>	10.0	sphere radius for atoms selection

Description: This command adds atoms to, removes atoms from, or initializes any one of the three independent sets of selected atoms of MODEL. There are three selection sets because it is convenient to have different sets used by different MODELLER commands.

`pick_atoms_set` specifies the set of selected atoms. Set 1 is used in the `model.restraints.pick()`, `model.rotate_dihedrals()`, `model.randomize_xyz()` and `model.mutate()` commands. Sets 2 and 3 are used in the `model.restraints.make()` command.

`selection_status` determines whether the selected atoms are added ('ADD'), removed ('REMOVE'), or a set is initialized and then the selected atoms are added ('INITIALIZE').

The selection of atoms is a hierarchical two level process. The first level of selection consists of specifying how the atoms will be scanned. The second level consists of selecting by the specified atom and residue names.

How the atoms are scanned is specified by setting the `selection_search` variable to either 'SEGMENT', 'SPHERE', or 'SPHERE_SEGMENT':

1. 'SEGMENT' mode: Only a single stretch of residues specified by the beginning and ending residue identifiers in `selection_segment` (Section 4.9.1) is scanned. Alternatively, if `selection_segment[1]` has the special value 'LOOPS' only residues in loops are scanned. Loops are defined as those residues in the MODEL that are aligned with only gap positions in the templates (MODEL has to be the last sequence in the current alignment), are within `gap_extension[1]` of an insertion in MODEL, or are within `gap_extension[2]` positions of a deletion in MODEL, and are not in a loop segment shorter than `minmax_loop_length[1]` or longer than `minmax_loop_length[2]`. This selection mode is useful for automatic selection of loops to be refined by the loop modeling procedure.
2. 'SPHERE' mode: Only those atoms that are closer than `sphere_radius` angstroms to the `sphere_center` atom, after the center atom was translated by (`xtrans`, `ytrans`, `ztrans`) angstroms specified in `selection_slab[3:5]`,

are scanned. If the first element of `sphere_center` is string 'INDEX', the second element is an integer atom index of the center atom; otherwise, the first and second element are the residue identifier (Section 4.9.1) and the IUPAC atom name, respectively. `selection_slab[1:2]` specifies the interval on the *Z*-axis relative to the *Z* coordinate of the translated central atom that imposes another condition on the selected atoms: $Z_{cen} + dz1 < Z + ztrans < Z_{cen} + dz2$. Larger *Z* values are in front, so *dz1* specifies the plane that is further away than the *dz2* plane. To pick any atoms, $dz1 < dz2$.

3. 'SPHERE_SEGMENT' mode: Only atoms within a sphere around the atoms in the specified segment of residues are scanned. This is useful, for example, when a neighborhood of a loop needs to be selected. As for the 'SEGMENT' mode, if `selection_segment[1]` has the special value 'LOOPS', only loop atoms are scanned for their neighbors.

If `selection_from` is 'SELECTED', scanning specified above is restricted only to the atoms that were already selected before calling `model.pick_atoms()`.

Once the method for scanning the atoms is specified, each of the scanned atoms is checked against the specified atom name(s) (`atom_types`) and residue name(s) (`res_types`). If `selection_mode` is 'RESIDUE', all atoms in a residue with at least one atom that matches both the residue and atom name criteria are selected. Otherwise, only those atoms that have both the specified residue and atom names are selected. The `res_types` and `atom_types` keywords can contain several residue and atom names in one quoted string or in several quoted strings. For example, both 'CA' 'N' and 'CA N' are valid specifications selecting the CA and N atoms. The following groups of residues and atoms are defined:

- If `res_types` contains word 'ALL', all residues will be selected.
- If `res_types` contains word 'HET', all 'HETATM' residues will be selected (*e.g.*, all residue types with the MODELLER residue code larger than 27; see library \$RESTYP_LIB).
- If `res_types` contains word 'BLK', all 'BLK' residue types will be selected (Section 4.5.1).
- If `res_types` contains word 'STD', all standard residue types will be selected. Standard residue types are all residue types but 'HETATM' and 'BLK' types.
- If `atom_types` contains word 'ALL', all atoms will be selected.
- If `atom_types` contains word 'MNCH', all mainchain atoms will be selected. Mainchain atoms are N, C, CA, O, and OXT.
- If `atom_types` contains word 'SDCH', all sidechain atoms will be selected. Sidechain atoms are all non-mainchain atoms, including non-mainchain atoms in 'HETATM' and 'BLK' residues.

Example: examples/commands/pick_atoms.py

```
# Example for: model.pick_atoms()

# This will pick various subsets of atoms in the MODEL and compare them
# with MODEL2.

env = environ()
log.level(1, 1, 1, 0)

# Set some defaults (the same as in top.ini):
env.selection_mode = 'ATOM' # only the selected atoms, not whole residues
env.selection_from = 'ALL' # scanning of all atoms, not selected atoms
env.selection_search = 'SEGMENT' # scan over a segment
env.selection_segment = ('FIRST:', 'LAST:') # the whole chain as a segment
env.res_types = 'ALL' # all residue types
env.pick_atoms_set = 1 # put the selected atoms in set 1
env.selection_status = 'INITIALIZE' # select only the selected atoms

# Read the models and the alignment:
mdl = model(env, file='ifas')
```

```

mdl2 = model(env, file='2ctx')
aln = alignment(env, file='toxin.ali', align_codes=('1fas', '2ctx'))
aln.write(file='toxin.pap', alignment_format='PAP')

# Pick and superpose mainchain atoms:
mdl.pick_atoms(aln, atom_types='MNCH')
mdl.superpose(mdl2, aln)

# Pick and superpose sidechain atoms:
mdl.pick_atoms(aln, atom_types='SDCH')
mdl.superpose(mdl2, aln)

# Pick and superpose CA and CB atoms:
mdl.pick_atoms(aln, atom_types='CA CB')
mdl.superpose(mdl2, aln)

# Pick and superpose all atoms:
mdl.pick_atoms(aln, atom_types='ALL')
mdl.superpose(mdl2, aln)

# Pick and superpose CA and CB atoms in one segment only:
mdl.pick_atoms(aln, atom_types='CA CB', selection_segment=('2:', '10:'))
mdl.superpose(mdl2, aln)

# Pick and superpose all atoms within 6 angstroms of the 'CA' atom in residue '10:':
mdl.pick_atoms(aln, atom_types='ALL', sphere_radius=6.0,
               selection_search='SPHERE', sphere_center=('10:', 'CA'))
mdl.superpose(mdl2, aln)

# Pick and superpose all atoms within 6 angstroms of any atom in
# segment 2: to 10:
mdl.pick_atoms(aln, atom_types='ALL', selection_segment=('2:', '10:'),
               selection_search='SPHERE_SEGMENT', sphere_radius=6.0)
mdl.superpose(mdl2, aln)

# Pick and superpose all atoms in all loops (ie residues within 2 positions
# of any gap in the alignment):
mdl.pick_atoms(aln, atom_types='ALL', selection_segment=('LOOPS', ''),
               selection_search='SEGMENT', gap_extension=(2, 2))
mdl.superpose(mdl2, aln)

# Pick and superpose all atoms within 6 angstroms of all loops (ie residues
# within 2 positions of any gap in the alignment):
mdl.pick_atoms(aln, atom_types='ALL', selection_segment=('LOOPS', ''),
               selection_search='SPHERE_SEGMENT', sphere_radius=6.0,
               gap_extension=(2, 2))
mdl.superpose(mdl2, aln)

```

4.6.17 model.pick_hot_atoms() — pick atoms violating restraints

Options:

edat = <energy_data>

objective function parameters

<code>viol_report_cut = <float : 35></code>	4.5 4.5 4.5 4.5 4.5 4.5 4.5	cutoffs for selecting violated restraints
	4.5 4.5 4.5 4.5 4.5 4.5 999	
	999 999 999 4.5 4.5 4.5 4.5	
	4.5 4.5 999 6.5 4.5 4.5 4.5	
	4.5 4.5 999 999 999 4.5 4.5	
<code>pick_hot_cutoff = <float : 1></code>	4.0	radius for picking hot atoms
<code>selection_mode = <str : 1></code>	'ATOM'	selecting what: 'ATOM' 'RESIDUE'
<code>extend_hot_spot = <int : 1></code>	0	whether to extend hot spots
<code>residue_span_range = <int : 2></code>	0 99999	range of residues spanning the allowed distances; for MAKE_RESTRAINTS, PICK_RESTRAINTS, non-bonded dynamic pairs

Description: This command selects those selected atoms (set 1) in the MODEL that should be optimized to remove hot spots in the MODEL; only selected restraints are considered.

More precisely, the command first flags violated selected atoms. An atom is violated if it is part of a violated restraint. A restraint of physical group i (Table 4.5) is violated when its relative deviation from the optimal value is larger than specified in `viol_report_cut[i]`. For restraints that are based on probability density functions, relative violation is defined as the difference between the actual and the ideal values divided by the standard deviation ('relative heavy violation'); energy based restraints have *ad hoc* definition of violations (Table 4.3).

The command then flags those selected atoms that are within the `pick_hot_cutoff` angstroms of any of the already flagged atoms.

Next, if `selection_mode` is 'RESIDUE', all atoms in the residues that have at least one atom flagged are also flagged. In addition, the contiguous segments of flagged residues are extended for `extend_hot_spot` residues on either side.

This command is usually followed by the `model.restraints.pick()` and `model.optimize()` commands to select all the restraints that operate on selected (hot) atoms and optimize positions of these hot atoms.

Example: `examples/commands/pick_hot_atoms.py`

```
# Example for: model.pick_hot_atoms()

# This will pick atoms violated by some restraints (bond length restraints here),
# select restraints operating on violated atoms, and calculate the energy for
# the selected restraints only (note that a list of violated restraints
# can be obtained by the ENERGY command alone, without preceding it with
# PICK_HOT_ATOMS).

env = environ()
env.edat.dynamic_sphere = False
env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

# Read the sequence, calculate its topology and coordinates:
code = '1fas'
mdl = model(env, file=code)
aln = alignment(env)
aln.append_model(mdl, atom_files=code, align_codes=code)
aln.append_model(mdl, atom_files=code+'.ini', align_codes=code+'-ini')
mdl.generate_topology(aln, sequence=code+'-ini')
mdl.transfer_xyz(aln)
# Just to get some violations:
mdl.randomize_xyz(deviation=0.03)
# Create the bond length restraints and ignore the hard sphere overlap:
```

```
mdl.restraints.make(aln, restraint_type='bond', spline_on_site=False)
# Pick hot atoms and the corresponding violated and neighbouring restraints:
mdl.pick_hot_atoms(pick_hot_cutoff=4.0)
mdl.restraints.pick(add_restraints=False)
# Calculate the energy of the selected restraints and write them out in detail:
mdl.energy(output='VERY_LONG')
```

4.6.18 `model.randomize_xyz()` — randomize MODEL coordinates

Options:

```
deviation = <float : 1>          0.0          coordinate randomization amplitude in
                                angstroms
```

Description: This command randomizes the Cartesian coordinates of the selected atoms (set 1) in MODEL. If deviation is positive, the coordinates are randomized by the *addition* of a random number uniformly distributed in the interval from $-deviation$ to $+deviation$ angstroms. If deviation is negative, the coordinates are *assigned* a random value uniformly distributed in the interval from $-deviation$ to $+deviation$ angstroms.

Example: `examples/commands/randomize_xyz.py`

```
# Example for: model.randomize_xyz()

# This will randomize the X,Y,Z of the model:

env = environ()

mdl = model(env, file='1fas')

# Change existing X,Y,Z for +- 4 angstroms:
mdl.randomize_xyz(deviation=4.0)
mdl.write(file='1fas.ini1')

# Assign X,Y,Z in the range from -100 to 100 angstroms:
mdl.randomize_xyz(deviation=-100.0)
mdl.write(file='1fas.ini2')
```

4.6.19 `model.to_iupac()` — standardize certain dihedral angles

Description: This routine swaps specific pairs of atoms within some residues of MODEL so that certain dihedral angles are within $\pm 90^\circ$, satisfying the IUPAC convention [IUPAC-IUB, 1970, Kendrew *et al.*, 1970]. These residues, pairs of atoms, and dihedral angles are:

- Phe, Tyr: (CD1, CD2), (CE1, CE2); χ_2 ;
- Asp: (OD1, OD2); χ_2 ;
- Glu: (OE1, OE2); χ_3 ;
- Arg: (NH1, NH2); χ_4 .

It is possible that for distorted sidechains, neither of the two possibilities satisfies the IUPAC convention. In such a case, a warning message is written to the `log` file.

Example: `examples/commands/iupac_model.py`

```
# This will swap certain atom names in some planar sidechains to satisfy
# the IUPAC convention.

env = environ()
log.level(1, 1, 1, 1, 0)

mdl = model(env, file='2abx')
mdl.to_iupac()
mdl.write(file='2abx.iup')
```

4.6.20 `model.reorder_atoms()` — standardize order of MODEL atoms**Requirements:** topology library**Description:** This routine reorders atoms within the residues of MODEL so that they follow the order in the current residue topology library.**Example:** `examples/commands/reorder_atoms.py`

```
# Example for: model.reorder_atoms()

# This will standardize the order of atoms in the model.

env = environ()

# Order the atoms according to a topology library:
env.libs.topology.read(file='$(LIB)/top_heav.lib')

mdl = model(env, file='1fas')
mdl.reorder_atoms()
mdl.write(file='1fas.ini')
```

4.6.21 `model.rotate_dihedrals()` — change dihedral angles**Options:**

<code>dihedrals = <str : 0></code>	<code>'PHI' 'PSI' 'CHI1' 'CHI2'</code> <code>'CHI3' 'CHI4'</code>	dihedral angle type selection: <code>'phi' 'psi' 'omega' 'chi1' 'chi2' 'chi3' 'chi4' 'chi5' 'alpha'</code>
<code>change = <str : 1></code>	<code>'RANDOMIZE'</code>	what to do: <code>'RANDOMIZE' 'OPTIMIZE'</code>
<code>deviation = <float : 1></code>	<code>0.0</code>	amplitude of dihedral angle randomization

Requirements:

for `change='OPTIMIZE'`: topology & restraints
 for `change='RANDOMIZE'`: topology

Description: This command changes the dihedral angles in MODEL.

`change` selects an optimization (when equal to `'OPTIMIZE'`) or randomization (when equal to `'RANDOMIZE'`):

1. When optimizing, this command finds the first selected restraint that restrains the specified dihedral angle of each selected residue. It then sets the value of that dihedral to the most likely value. A residue is selected if any of its atoms is in the set 1 of selected atoms.
2. When randomizing, the command changes the specified dihedral angle of each selected residue by adding a random value distributed uniformly from $-deviation$ to $+deviation$ degrees.

dihedrals can be either a vector of dihedral angle names or a single string containing all the dihedral angle names separated by blanks. The dihedral angles involved in cyclic structures are not changed (*e.g.*, sidechain dihedral angles in disulfide bonds and prolines). The dihedral angles that can be changed are listed at the top of the \$RESDIH_LIB library: `alpha`, `phi`, `psi`, `omega`, `chi1`, `chi2`, `chi3`, `chi4`, `chi5`. Dihedral angle `'alpha'` is the virtual C_α dihedral angle defined by four consecutive C_α atoms.

The bond connectivity of the MODEL has to exist before this command is executed. If you read in the model by `model.read()`, the bond connectivity is defined by subsequent calls to `topology.append()` and `model.generate_topology()` (also make sure that `sequence` entry does not exist in the alignment or that no alignment is in memory).

Example: examples/commands/rotate_dihedrals.py

```
# Example for: model.rotate_dihedrals()

# This will optimize and randomize dihedrals in a MODEL
env = environ()
env.libs.topology.read(file='${LIB}/top_heav.lib')
env.libs.parameters.read(file='${LIB}/par.lib')

# Select dihedral angle types for optimization and randomization:
dih = 'phi psi omega chi1 chi2 chi3 chi4 chi5'

# Read the sequence, get its topology and coordinates:
code='1fas'
mdl = model(env, file=code)
aln = alignment(env)
aln.append_model(mdl, align_codes=code, atom_files=code)
aln.append_model(mdl, align_codes=code+'_ini', atom_files=code+'_ini')
mdl.generate_topology(aln, sequence=code+'_ini')
mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')
mdl.rotate_dihedrals(change='RANDOMIZE', deviation=90.0, dihedrals=dih)
mdl.write(file='1fas.ini1')

# Get restraints from somewhere and optimize dihedrals:
mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False)
mdl.rotate_dihedrals(change='OPTIMIZE', deviation=90.0, dihedrals=dih)
mdl.write(file='1fas.ini2')
```

4.6.22 model.orient() — center and orient MODEL

Description: This command translates the MODEL so that its gravity center is at the origin of the coordinate system and that the three principal axes of the model's inertia ellipsoid correspond to the x , y , and z axes of the coordinate system. It may even be used for approximate superposition if molecules have a similar non-spherical shape. Information about the principal axes is written to the `log` file.

Example: examples/commands/orient_model.py

```
# Example for: model.orient()

# This will orient the model along the principal axes of the inertia ellipsoid:

env = environ()
mdl = model(env)
mdl.read(file='1fas')
mdl.orient()
mdl.write(file='1fas.ini')
```

4.6.23 model.rotate() — rotate and translate MODEL

Options:

translation = \langle float : 3	0.0 0.0 0.0	translation vector for MODEL
rotation_matrix = \langle float : 9	1 0 0 0 1 0 0 0 1	rotation matrix for MODEL
rotation_angle = \langle float : 1	0.0	rotation of MODEL around axis [degrees]
rotation_axis = \langle float : 3	1.0 0.0 0.0	rotation axis for MODEL

Description: This command transforms the Cartesian coordinates of MODEL.

Translation is specified by a translation vector `translation` and is done first.

Rotation is specified by a rotation matrix `rotation_matrix` that is given as a vector of 9 elements (three rows times three columns), with column index running first: $a_{11}, a_{12}, a_{13}, a_{21}, a_{22}, a_{23}, a_{31}, a_{32}, a_{33}$. The rotation matrix pre-multiplies the Cartesian coordinate vectors. The matrix corresponds to the view matrix of QUANTA and to the rotation matrix of MOLSCRIPT.

The second kind of rotation is specified by a screw transformation, given by the `rotation_axis` axis and `rotation_angle` rotation around the axis (in degrees). This is done last.

If only some transformations are desired, set the other values to 0.

Example: `examples/commands/rotate_model.py`

```
# Example for: model.rotate()

# This will orient a model as specified (no change in this example):

# Read the structure and transform it:
env = environ()
mdl = model(env, file='1fas')
mdl.rotate(translation=(0, 0, 0),
           rotation_matrix=(1, 0, 0,
                             0, 1, 0,
                             0, 0, 1),
           rotation_angle=0, rotation_axis=(1, 1, 1))
mdl.write(file='1fas.ini')
```

4.6.24 `model.write_data()` — write derivative MODEL data

Options:

<code>edat = <energy_data></code>		objective function parameters
<code>file = <str : 1></code>	<code>'default'</code>	root of output filename(s)
<code>accessibility_type = <int : 1></code>	8	type of solvent accessibility: 1–10
<code>surftyp = <int : 1></code>	1	Surface Type for accessibility calculations 1= contact; 2=surface
<code>psa_integration_step = <float : 1></code>	0.1	integration step for WRITE_DATA
<code>probe_radius = <float : 1></code>	1.4	probe_radius for WRITE_DATA
<code>neighbor_cutoff = <float : 1></code>	6.0	for defining atom-atom contacts in WRITE_DATA
<code>output = <str : 1></code>	<code>'LONG'</code>	what to calculate and write out: <code>'ALL'</code> <code>'PSA'</code> <code>'ATOMIC_SOL'</code> <code>'NGH'</code> <code>'DIH'</code> <code>'SSM'</code> <code>'CRV'</code> <code>'CAV'</code> <code>'CROSS-SECTIONS'</code>

Requirements: topology file

Description: This command writes the selected types of data about the MODEL to a corresponding file and to the ‘fourth’ column of the model. The root of the output filenames is given by the file variable. In addition to the output files, the B_{ISO} field of the model (‘fourth column’ in the PDB file) will be assigned the last selected property from the following list: atomic or residue accessibility, dihedral type `accessibility_type` (from 1 to 9 for α , Φ , Ψ , ω , χ_1 , χ_2 , χ_3 , χ_4 , and χ_5 ; where α is the virtual dihedral angle between four successive C_α atoms), number of residue neighbors, the secondary structure type, and the local mainchain curvature. For accessibility, when `output` contains `ATOMIC_SOL`, atomic accessibilities in \AA^2 are assigned to B_{ISO} ; otherwise residue accessibility of type `accessibility_type` (from 1 to 10, for the columns in the `.psa` file) is assigned. If `surftyp` is 1, contact accessibility is calculated; if 2, surface accessibility is returned.

The data to be calculated are specified by concatenating the corresponding keywords in the `output` variable:

- `'ALL'`: All types of data are written to the corresponding files.
- `'PSA'`: The atomic and residue solvent accessibilities are written to the `.sol` and `.psa` files, respectively. The algorithm for the solvent contact areas is described in [Richmond & Richards, 1978]. The normalization for the fractional areas is carried out as described in [Hubbard & Blundell, 1987], with the normalization factors courtesy of Simon Hubbard (personal communication). The single reference is Šali & Overington, 1994. Accessibilities are calculated with scaled radii from the `$MODELS_LIB` library, as specified by `topology.submodel`. The radii are scaled by `energy_data.radii_factor`, which should usually be set to 1.
- `'CAV'`: The protein and internal cavity volumes are written out. The calculation on a grid is used. The grid unit is specified by `grid_unit` in angstroms (say 1.4 \AA). The radii are scaled by `energy_data.radii_factor`, which should usually be set to 1. The cross-sections are written to file `file.cav` when `output` contains `CROSS-SECTIONS`. The `number_of_steps` is the number of small shifts along x, y, and z that are used in the averaging of the protein and cavity volumes with respect to small changes in the relative position of the protein and the grid; the total number of calculations is therefore equal to the third power of `number_of_steps`. If `orient` is `True`, the structure is oriented before the volume calculation such that the moment of inertia are parallel to the x, y, and z coordinate axes (this orientation minimizes the size of the grid). However, the coordinates of the MODEL are not changed upon exit from this routine (you need to use `model.orient()` to change the orientation of the MODEL).
- `'NGH'`: Residue neighbors of each residue are listed to a `.ngh` file. The MODELLER definition of a residue–residue contact used in restraints derivation is applied [Šali & Blundell, 1993]: Any pair of residues that has any pair of atoms within 6 \AA of each other are in contact.
- `'DIH'`: All the dihedral angle types defined in the `$RESDIH_LIB` library (virtual C_α , mainchain, and sidechain dihedral angles) are written to a `.dih` file.

- 'SSM': Secondary structure assignments are written to a `.ssm` file. The algorithm for secondary structure assignment depends on the C_α positions only and is based on the distance matrix idea described in [Richards & Kundrot, 1988]. For each secondary structure type, a 'library' C_α distance matrix was calculated by averaging distance matrices for several secondary structure segments from a few high resolution protein structures. Program DSSP was used to assign these secondary structure segments [Kabsch & Sander, 1983]. Outlier distances were omitted from the averaging. Currently, there are only two matrices: one for the α -helix (secondary structure type 2) and one for the β -strand (type 1). The algorithm for secondary structure assignment is as follows:
 1. For each secondary structure type (begin with a helix, which can thus overwrite parts of strand if they overlap):
 - Define the degree of the current secondary structure fit for each C_α atom by DRMS deviation (P_1) and maximal distance difference (P_2) obtained by comparing the library distance matrix with the distance matrix for a segment starting at the given C_α position;
 - Assign the current secondary structure type to all C_α 's in all segments whose DRMS deviation and maximal distance difference are less than some cutoffs ($P_1 < cut_1$, $P_2 < cut_2$) and are not already assigned to 'earlier' secondary structure types;
 2. Split kinked contiguous segments of the same type into separate segments: Kinking residues have both DRMS and maximal distance difference beyond their respective cutoffs ($P_1 > cut_3$, $P_2 > cut_4$). The actual single kink residue separating the two new segments of the same type is the central kinking residue. Note: we are assuming that there are no multiple kinks within one contiguous segment of residues of the same secondary structure type. The kink residue type is `-2`.
 3. If the current secondary structure type is β -strand: Eliminate those runs of strand residues that are not close enough to other strand residues separated by at least two other residues: P_3 is minimal distance to a non-neighboring residue of the strand type ($P_3 < cut_3$). Currently, only one pass of this elimination is done, but could be repeated until self-consistency.
 4. Eliminate those segments that are shorter than the cutoff (cut_6) length (e.g., 5 or 6).
 5. Remove the isolated kinking residues (those that occur on their own or begin or end a segment).
- 'CRV': Local mainchain curvatures are written to a `.crv` file. Local mainchain curvature at residue i is defined as the angle between the least-squares lines through C_α atoms $i - 3$ to i and i to $i + 3$.

Example: `examples/commands/write_data.py`

```
# Example for: model.write_data()

# This will calculate solvent accessibility, dihedral angles, and
# residue-residue neighbors for a structure in the PDB file.

log.verbose()

# Get topology library for radii and the model without waters and HETATMs:
env = environ()
env.io.hetatm = False
env.io.water = False

env.libs.topology.read(file='${LIB}/top_heav.lib')
mdl = model(env, file='1fas')

# Calculate residue solvent accessibilities, dihedral angles, and
# residue-residue neighbors:
myedat = energy_data()
myedat.radii_factor = 1.0 # The default is 0.82 (for soft-sphere restraints)
mdl.write_data(file='1fas', edat=myedat, output='PSA DIH NGH SSM CRV')
```

4.6.25 `model.write_pdb_xref()` — write residue number/index correspondence

Options:

<code>aln = <alignment></code>		Alignment between model and templates
<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
<code>output_directory = <str : 1></code>	<code>''</code>	output directory
<code>model_segment = <str : 2></code>	<code>'FIRST:@' 'LAST:'</code>	segment to be read in

Description: This command writes the correspondence between the PDB residue numbers and residue indices for the selected part of the MODEL. It is more useful than one would think because of its interaction with the alignment data and the option to use wild characters to specify the beginning and ending residues.

Example: `examples/commands/write_pdb_xref.py`

```
# Example for: model.write_pdb_xref()

# This writes out information useful for relating PDB residue numbers with
# residue indices.

log.verbose()
env = environ()

mdl = model(env, file='2abx.atm', model_segment=('FIRST:@', 'END:'))
aln = alignment(env)
aln.append_model(mdl, align_codes='2abx')

# if not found, 1, NRES used
mdl.write_pdb_xref(aln, file='2abx.xref1', model_segment=('FIRST:@', 'END:'))

# You can use string resid's
mdl.write_pdb_xref(aln, file='2abx.xref2', model_segment=('!:', '50:'))

# You can use integer resid's
mdl.write_pdb_xref(aln, file='2abx.xref3', model_segment=('!2', '!50'))

# You can even use the alignment specs, but not with codes that start with '!'
mdl.write_pdb_xref(aln, file='2abx.xref4', model_segment=('2abx', '2abx'))

# You can mix the specs
mdl.write_pdb_xref(aln, file='2abx.xref5', model_segment=('!2', '50:'))
mdl.write_pdb_xref(aln, file='2abx.xref6', model_segment=('!2', 'END:'))
```

4.6.26 `model.make_region()` — define a random surface patch of atoms

Options:

<code>atom_accessibility = <float : 1></code>	<code>1.0</code>	accessible atoms for MAKE_REGION
<code>region_size = <int : 1></code>	<code>20</code>	size of exposed region in MAKE_REGION

Description: This command defines a contiguous patch of exposed atoms of the specified size. First, the exposed atoms in MODEL are identified by using the `atom_accessibility` cutoff (in \AA^2). The seed atom is picked

randomly among the exposed atoms. The patch is expanded by iteratively adding the exposed atom that is closest to the gravity center of the currently selected patch atoms. Thus, the patch is defined deterministically once the seed atom is picked. The patch is defined by setting the fourth column parameter (B_{iso}) to 1 for the patch atoms and to 0 for the remaining atoms. The “temperature” color option of Rasmol can be used to display the patch graphically.

To obtain surface patches that look good in visual inspection, it is necessary to use a non-obvious scaling factor for atomic radii and probe radius for solvent calculation by `model.write_data()`, as well as the accessibility cutoff for `model.make_region()`.

Example: `examples/commands/make_region.py`

```
# Example for: model.make_region()

# This will define a random contiguous patch of atoms on a surface of the
# protein.

env = environ(rand_seed=-18343)
log.level(1, 1, 1, 1, 0)

# Read the PDB file
mdl = model(env)
mdl.read(file='../atom_files/pdb1fdn.ent')

# Calculate atomic accessibilities with appropriate probe_radius
myedat = energy_data()
myedat.radii_factor = 1.6
mdl.write_data(file='1fdn', edat=myedat, output='PSA ATOMIC_SOL',
               psa_integration_step=0.05, probe_radius=0.1)

# Get the "random" patch of exposed atoms on the surface
mdl.make_region(atom_accessibility=0.5, region_size=35)

# Write out a PDB file with the patch indicated by Biso = 1:
mdl.write(file='1fdn.reg')
```

4.6.27 `model.color()` — color MODEL according to alignment

Options:

`aln = <alignment>` Alignment between the model and another sequence

Description: This command colors MODEL according to a given alignment, `aln`, between MODEL and a sequence. MODEL has to be the first protein in the alignment. The second protein can be any sequence, with or without known structure. The MODEL can be displayed on the screen, colored by ‘the fourth parameter’ and inspected for the structural context of deletions and insertions. This is useful in optimizing the alignment for comparative modeling. The isotropic temperature factors in MODEL are set as follows:

- 0, for those regions that have residues in both MODEL and the sequence (blue in RASMOL; light green in QUANTA);
- 1, for the two residues that span regions occurring in the sequence but not in MODEL (green in RASMOL; pink in QUANTA);
- 2, regions that occur in MODEL but are deleted from the sequence (red in RASMOL; bright green in QUANTA).

Example: examples/commands/color_aln_model.py

```

# Example for: model.color()

# Two demos:
#
# 1) Use a given alignment to color a structure according to
#     insertions and deletions in a pairwise alignment.
#
# 2) Superpose two 3D structure and do (1).

env = environ()

# Demo 1:
mdl = model(env)
aln = alignment(env)
mdl.read(file='1nbt')
aln.append(file='toxin.ali', align_codes=('1nbt', '1fas'), remove_gaps=True)
mdl.color(aln)
mdl.write(file='1nbt-1.clr')

# Demo 2:
aln = alignment(env)
segs = {'1nbt':('1:A', '66:A'), '1fas':('1:', '61:')}
for code in ('1nbt', '1fas'):
    mdl.read(file=code, model_segment=segs[code])
    aln.append_model(mdl, align_codes=code, atom_files=code)
aln.align(gap_penalties_1d=(-600, -400))
aln.malign3d(gap_penalties_3d=(0, 3.0))
aln.write(file='color_aln_model.pap', alignment_format='PAP')

mdl.read(file='1nbt', model_segment=segs['1nbt'])
mdl.color(aln)
mdl.write(file='1nbt-2.clr')

```

4.6.28 model.superpose() — superpose MODEL2 on MODEL given alignment**Options:**

<code>mdl = <model></code>		Input model for superposition
<code>aln = <alignment></code>		Alignment between the two models
<code>fit = <bool : 1></code>	True	whether to superpose
<code>superpose_refine = <bool : 1></code>	False	whether to refine the superposition
<code>rms_cutoffs = <float : 11></code>	3.5 3.5 60 60 15 60 60 60 60 60 60	only the first element is used for calculating the cutoff RMS and DRMS measures
<code>reference_atom = <str : 1></code>	''	reference atom name in SUPERPOSE
<code>reference_distance = <float : 1></code>	3.5	cutoff for selecting reference positions in SUPERPOSE
<code>refine_local = <bool : 1></code>	True	whether to refine superposition using locally similar substructures
<code>swap_atoms_in_res = <str : 1></code>	''	minimize RMS by swapping atoms in these residues (1 char code: 'DEFHLNQRVY')

Description: This command superposes `mdl` on `MODEL`, without changing the alignment, `aln`.

`MODEL` must be the first sequence in the alignment; `mdl` must be the second sequence in the alignment. The equivalent atoms are those selected atoms (set 1) of the `MODEL` that have equivalently named atoms in `mdl`; the atom equivalences are defined in library `$ATMEQV_LIB`. Use the `model.pick_atoms()` command to select the desired atoms for superposition. By default, all atoms are selected.

No fitting is done if `fit = False`.

The `rms_cutoffs[1]` element is the cutoff used in calculating the cutoff RMS deviations; *i.e.*, those position and distance RMS deviations that are defined on the equivalent atoms which are less than `rms_cutoffs[1]` angstroms away from each other (as superposed using all aligned positions) and those equivalent distances which are less than `rms_cutoffs[1]` angstroms different from each other, respectively.

If `superpose_refine` is `True` the refinement of the superposition is done by repeating the fitting with only those aligned pairs of atoms that are within `rms_cutoffs[1]` of each other until there is no change in the number of equivalent positions. This refinement can only remove compared positions, not add them like `alignment.align3d()` can do. This is useful for comparing equivalent parts of two structures with a fixed alignment but omitting divergent parts from the superposition and RMS deviation calculation; *e.g.*, comparing a model with the X-ray structure.

If `superpose_refine` is `False` and `reference_atom` is non-blank, only those pairs of equivalently named selected atoms from aligned residues are superposed that come from residues whose `reference_atom` atoms are closer than `reference_distance` Å to each other.

When `MODEL` and `mdl` have exactly the same atoms in the same order, one can set `swap_atoms_in_res` to any combination of single character amino acid residue codes in `DEFHLNQRVY`. Certain atoms (see below) in the specified sidechains of `mdl` are then swapped to minimize their RMS deviation relative to `MODEL`. The labelling resulting in the lowest RMS deviation is retained. The following swaps are attempted:

Residue	Swap(s)
D	OD1, OD2
E	OE1, OE2
F	CD1, CD2 CE1, CE2
H	ND1, CD2 NE2, CE1
N	OD1, ND2
Q	OE1, NE2
R	NH1, NH2
V	CG1, CG2
Y	CD1, CD2 CE1, CE2

Example: `examples/commands/superpose.py`

```
# Example for: model.superpose()

# This will use a given alignment to superpose Calpha atoms of
# one structure (2ctx) on the other (1fas).

env = environ()

mdl = model(env, file='1fas')
mdl2 = model(env, file='2ctx')
aln = alignment(env, file='toxin.ali', align_codes=('1fas', '2ctx'))
mdl.pick_atoms(aln, pick_atoms_set=1, atom_types='CA')
mdl.superpose(mdl2, aln)
mdl2.write(file='2ctx.fit')
```

Example: examples/commands/align3d.py

```

# Example for: alignment.align3d(), model.superpose()

# This will align 3D structures of two proteins:

log.verbose()
env = environ()

# First example: read sequences from a sequence file:
aln = alignment(env)
aln.append(file='toxin.ali', align_codes=['1fas', '2ctx'])
aln.align(gap_penalties_1d=[-600, -400])
aln.align3d(gap_penalties_3d=[0, 4.0])
aln.write(file='toxin-str.ali')

# Second example: read sequences from PDB files to eliminate the
# need for the toxin.ali sequence file:
mdl = model(env)
aln = alignment(env)
for code in ['1fas', '2ctx']:
    mdl.read(file=code)
    aln.append_model(mdl, align_codes=code, atom_files=code)
aln.align(gap_penalties_1d=(-600, -400))
aln.align3d(gap_penalties_3d=(0, 2.0))
aln.write(file='toxin-str.ali')

# And now superpose the two structures using current alignment to get
# various RMS's:
mdl = model(env, file='1fas')
mdl.pick_atoms(aln, atom_types='CA')
mdl2 = model(env, file='2ctx')
mdl.superpose(mdl2, aln)

```

Example: examples/commands/swap_atoms_in_res.py

```

# This script illustrates the use of the swap_atoms_in_res
# argument to the model.superpose() command:

# Need to make sure that the topologies of the two molecules
# superposed are exactly the same:

env = environ()
env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

atfil = '../atom_files/1fdx.atm'
mdl = model(env, file=atfil)
aln = alignment(env)
aln.append_model(mdl, align_codes='1fdx', atom_files=atfil)
aln.append_model(mdl, align_codes='1fdx', atom_files=atfil)
mdl.generate_topology(aln, sequence='1fdx')
mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')

```

```
mdl2 = model(env, file='./1fdx.swap.atm')
mdl.superpose(mdl2, aln, swap_atoms_in_res='')
mdl.superpose(mdl2, aln, swap_atoms_in_res='DEFHLNQRVY', fit=False)
mdl.superpose(mdl2, aln, swap_atoms_in_res='', fit=True)
```

4.6.29 model.make_chains() — Fetch sequences from PDB file

Options:

file = <str : 1>	'default'	partial or complete filename
structure_types = <str : 1>	'structure'	'structure structureX structureN structureM structureF structureE structureU'
minimal_resolution = <float : 1>	99.0	resolution threshold to select chains
minimal_chain_length = <int : 1>	30	length cut-off for selecting chains
max_nonstdres = <int : 1>	10	max. number of non-standard residues
chop_nonstd_terminii = <bool : 1>	True	whether or not to chop non-standard N- and/or C-terminal residue in MAKE_CHAINS
minimal_stdres = <int : 1>	30	min. number of standard residues
alignment_format = <str : 1>	'PIR'	format of the alignment file: 'PIR' 'PAP' 'QUANTA' 'INSIGHT' 'FASTA'

Description: This command fetches the sequences of the various chains found in the PDB file that has been read into memory (see `model.read()`).

The sequence of every chain that matches the input criteria is written out to separate files.

`structure_types` refers to the experimental method used to determine the structure. The following types are recognized: 'structureX' for x-ray, 'structureN' for NMR and 'structureM' for model, 'structureE' for electron microscopy, 'structureF' for fiber diffraction, 'structureU' for neutron diffraction, 'structure' for any structure.

`minimal_resolution` refers to the cut-off value of the experimental resolution of the structure. Structures with resolutions larger than this threshold are not processed.

`minimal_chain_length` refers to the lower limit of the chain length. Chains whose lengths are smaller than this value are not processed.

`max_nonstdres` sets the maximum limit of non-standard residues that is tolerated.

`chop_nonstd_terminii` removes the non-standard residues found at the terminii of the chains.

`minimal_stdres` sets the minimum number of standard residues that are required to process the chain. Chains that don't have at least this number of standard residues are not written out.

`alignment_format` specifies the format of the output file.

Example: examples/commands/make_chains.py

```
# Example for: model.make_chains()

# This will read a PDB file (segment), and write out all of its chains
# satisfying the listed conditions into separate alignment files in the
# PIR format.

env = environ()
```

```
mdl = model(env, file='../atom_files/pdb1lzd.ent')
mdl.make_chains(file='1lzd', minimal_chain_length=30, minimal_resolution=2.0,
               minimal_stdres=30, chop_nonstd_terminii=True,
               structure_types='structureN structureX',
               alignment_format='PIR')
```

4.6.30 model.schedule.make() — create optimization schedule

Options:

library_schedule = <int : 1>	1	selects schedule from the \$SCHED_LIB library
schedule_scale = <float : 35>	1 1 1 1 1 1 1 1 1 1 1 1 1 1	factors for physical restraint types in scaling the schedule
	1 1 1 1 1 1 1	

Description: This command constructs an optimization schedule for the variable target function method for the current MODEL.

The template for construction of the schedule is the library_schedule-th entry in library file \$SCHED_LIB.

The usual schedule for the variable target function part of optimization in comparative modeling is as follows. The residue range (**model.restraints.pick()** and Section 4.7.3) is increased with increasingly larger steps until the protein length is reached. The scaling of homology-derived and bonded stereochemical restraints increases from a small value to 1 in the initial few steps to allow for imperfect starting geometries, especially those that result from **model.randomize_xyz()** and long insertions or deletions. The soft-sphere overlap restraints are slowly introduced only in the last four steps of the variable target function method to save CPU time and increase the radius of convergence. In comparative modeling by the 'model' script in the default mode, the variable target function method is usually followed by simulated annealing with molecular dynamics. In this last stage, all homology-derived and stereochemical restraints are generally used with the scaling factors of 1. There are a number of variables defined in the **automodel** class that can be used to influence the thoroughness of both the variable target function and molecular dynamics parts of the optimization (Chapter 2).

The scaling factors for all physical restraint groups, in all schedule steps, are multiplied by the corresponding scalar in **schedule_scale** (1 by default). This is useful when template-derived fold restraints have to be weakened relative to some external restraints, so that the fold can actually reflect these external restraints, even when they are quite different from the template-derived restraints.

This command is an alternative to the **model.schedule.read()** command.

Use the **model.schedule.write()** command to find out what the calculated schedule is. The schedule file written by the 'model' routine has an extension **.sch**.

Example: examples/commands/make_schedule.py

```
# Example for: model.schedule.make(), model.schedule.write(),
#             model.schedule.read()

# This will create an VTFM optimization schedule for a model
# and write it to a file.

env = environ()
mdl = model(env, file='1fas')
mdl.schedule.make(library_schedule=1)

# Write the schedule to a file:
```

```
mdl.schedule.write(file='1fas.sch')

# Read it in just for fun:
mdl.schedule.read(file='1fas.sch')
```

4.6.31 model.schedule.step — current optimization schedule step

Description: This is the current step in the optimization schedule. See `model.schedule.make()` for more information.

4.6.32 model.schedule.read() — read optimization schedule

Options:

<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
<code>schedule_scale = <float : 35></code>	<code>1 1 1 1 1 1 1 1 1 1 1 1 1 1 1</code>	factors for physical restraint types in
	<code>1 1 1 1 1 1 1 1 1 1 1 1 1 1 1</code>	scaling the schedule
	<code>1 1 1 1 1 1 1</code>	

Description: This command reads a text file that contains an optimization schedule for the variable target function method.

Each line in the file contains in free format the parameters for a single step of the variable target function method. These parameters are: step index (not used by the program), optimization method, maximal difference in residue indices of atoms restrained by the selected restraints (`model.restraints.pick()` and Section 4.7.3), and the scaling factors for all types of restraints. The smaller the scaling factor, the weaker the corresponding restraint.

See `model.schedule.make()` for explanation of `schedule_scale`.

Example: See `model.schedule.make()` command.

4.6.33 model.schedule.write() — write optimization schedule

Options:

<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
<code>output_directory = <str : 1></code>	<code>''</code>	output directory

Description: This command writes out the schedule for the variable target function method. This schedule file can then be read by the `model.schedule.read()` command.

Example: See `model.schedule.make()` command.

4.6.34 `model.energy()` — evaluate MODEL given restraints**Options:**

<code>edat = <energy_data></code>		objective function parameters
<code>viol_report_cut = <float : 35></code>	4.5 4.5 4.5 4.5 4.5 4.5 4.5 4.5 4.5 4.5 4.5 4.5 4.5 999 999 999 999 4.5 4.5 4.5 4.5 4.5 4.5 999 6.5 4.5 4.5 4.5 4.5 4.5 999 999 999 4.5 4.5	cutoffs for reporting relative violations
<code>viol_report_cut2 = <float : 35></code>	2.0 2.0 2.0 2.0 2.0 2.0 2.0 2.0 2.0 2.0 2.0 2.0 2.0 2.0	
<code>output = <str : 1></code>	'LONG'	'SHORT' 'LONG' 'VERY_LONG' 'GRADIENT' 'SYMMETRY' 'ENERGY_PROFILE' 'VIOLATIONS_PROFILE' 'NO_REPORT'
<code>normalize_profile = <bool : 1></code>	False	whether to normalize energy/violations profiles or not, by the number of terms per residue
<code>smoothing_window = <int : 1></code>	3	profiles are smoothed over $2*SW + 1$ residues
<code>schedule_scale = <float : 35></code>	1 1	factors for physical restraint types in scaling the schedule
<code>file = <str : 1></code>	'default'	partial or complete filename
<code>asgl_output = <bool : 1></code>	False	whether to write output for ASGL
<code>residue_span_range = <int : 2></code>	0 99999	range of residues spanning the allowed distances; for MAKE_RESTRAINTS, PICK_RESTRAINTS, non-bonded dynamic pairs

Output: molpdf**Requirements:** restraints

Description: The main purpose of this command is to compare spatial features of the current MODEL with the selected restraints in order to determine the violations of the molecular pdf. It lists variable amounts of information about the values of the basis, feature, and molecular pdf's for the current MODEL. All arguments that affect the value of the molecular pdf are also relevant for the `model.energy()` command.

Within this routine only, the scaling factors for the physical restraint types are obtained from the `model.schedule.step` step of the current schedule, multiplied by `schedule_scale` (the original values are returned upon exit from the routine). This allows easy reporting of only a selected subset of all restraints.

Most of the output goes to the `log` file. The output of the `model.energy()` command has to be examined carefully, at least at the end of the optimization, when the final model is produced. Additional output files, for the ASGL plotting program are created if `asgl_output = True` (undocumented).

output selects various kinds of output information:

- 'LONG' writes restraint violations one per line to the `log` file.
- 'VERY_LONG' writes the most detailed examination of the selected basis and feature pdf's to the `log` file, using several lines of output for each restraint.
- 'NO_REPORT' suppresses the output of the violated restraints summary (unless profiles are also requested) and also that of nonbond clashes.

- 'GRADIENT' writes the 'force' gradients for the currently selected restraints to the isotropic temperature factors for each atom of the current MODEL.
- 'SYMMETRY' writes a comparison of equivalent distances involved in the definition of the symmetry enforcing term to the log file.

`viol_report_cut` is a vector with one real number for each physical restraint type. A restraint is reported when its 'heavy relative violation' is larger than the corresponding cutoff. The heavy relative violation is calculated by finding the global minimum of a feature according to the restraint, taking the difference between the actual feature in the model and this global minimum, and then normalizing the difference by the standard deviation of the global minimum. The 'minimal violation' of a restraint is defined as the difference from the local minimum closest to the value of the feature in the model (with the exception of the spline restraints; see next paragraph).

`viol_report_cut2` is similar to `viol_report_cut`, except that it contains cutoffs for restraint 'energies', not heavy relative violations.

The meaning of various other reported properties of the violated restraints is briefly described in the log file. Note that for multi-modal restraints that are described by cubic splines (by default, all multimodal homology-derived restraints), only one optimal value is defined, not the local and global minimum as for the multi-modal Gaussian restraints. As a result, the minimal violations and heaviest violations are the same. For interpreting the seriousness of violations, use the following rule of thumb: There should be at most a few small violations (*e.g.*, 4 standard deviations) for all monomodal restraints. In comparative modeling, the monomodal restraints include the stereochemical restraints and distance restraints when only one homologous structure is used. For the multimodal restraints, there are usually many violations reported because the heaviest violations are used in deciding whether or not to report a violation. In comparative modeling, the multimodal restraints include the χ_i restraints, (Φ , Ψ) binormal restraints and distance restraints when more than one template is used. See also Section 1.8, Question 15.

For profiles:

This command calculates residue energies or heavy relative violations, depending on `output`, for all physical restraint types (there are NPHYCNS of them). Relative heavy violations (Table 4.3) are used because only *relative* violations of different features are comparable. In both cases, the residue sum is the sum over all restraints that have at least one atom in a given residue. The contribution of each restraint is counted exactly once for each residue, without any weighting. Restraints spanning more than one residue contribute equally to all of them. Thus, the sum of residue energies is generally larger than molecular pdf. The command also calculates the sum of the NPHYCNS contributions for each residue and writes all NPHYCNS+1 columns to a file suitable for plotting by ASGL.

If `normalize_profile` is True the profile for each residue is normalized by the number of terms applying to each residue.

All the curves are smoothed by the running window averaging method if `smoothing_window` is larger than 0: The window is centered on residue i and extends for $(\text{smoothing_window}/2) - 1$ residues on each side. Thus, `smoothing_window` has to be an even number (or it is made such by the program automatically). The only exceptions are the two termini, where a smaller number of residues are available for smoothing. The relative weight of residue j when calculating the smoothed value at residue i is $(\text{smoothing_window}/2 - |j - i|)$.

The energy or the violations profile is written to the fourth column of the MODEL atomic records (atomic isotropic temperature factors for X-ray structures). Note that all the atoms in one residue get the same number. This output is useful for exploring the violations on a graphics terminal.

See description of `model.optimize()` for the other variables.

Example: `examples/commands/energy.py`

```
# Example for: model.energy()

# This will calculate the stereochemical energy (bonds,
# angles, dihedrals, impropers) for a given model.

env = environ()
```

```

env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

aln = alignment(env)
code = "1fas"
mdl = model(env, file=code)
aln.append_model(mdl, atom_files=code, align_codes=code)
aln.append_model(mdl, atom_files=code+'.ini', align_codes=code+'-ini')

mdl.generate_topology(aln, sequence=code+'-ini')
# Must patch disulfides here to calculate the non-bonded
# energy properly. Also, when you use hydrogens, disulfides
# must always be patched so that sulfhydryl hydrogens are
# removed from the model.
for ids in [ ('17', '39'),
             ('3', '22'),
             ('53', '59'),
             ('41', '52') ]:
    mdl.patch(residue_type='DISU', residue_ids=ids)

mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')

mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False)
molpdf = mdl.energy(edat=energy_data(dynamic_sphere=True))

```

4.6.35 model.optimize() — optimize MODEL given restraints

Options:

edat = <energy_data>		objective function parameters
optimization_method = <int : 1>	999	type of optimization method: 1 3
residue_span_range = <int : 2>	0 99999	range of residues spanning the allowed distances; for MAKE_RESTRAINTS, PICK_RESTRAINTS, non-bonded dynamic pairs
trace_output = <int : 1>	0	modulus for writing information about optimization iterations: 0 for nothing
max_iterations = <int : 1>	200	maximal iterations in optimization
output = <str : 1>	'LONG'	'NO_REPORT' 'REPORT'
• For conjugate gradients:		
min_atom_shift = <float : 1>	0.010	minimal atomic shift for the optimization convergence test
• For molecular dynamics:		
md_time_step = <float : 1>	4.0	time step for MD in fs
init_velocities = <bool : 1>	True	whether to initialize velocities before MD
temperature = <float : 1>	293.0	temperature for MD simulation in K
equilibrate = <int : 1>	999999	equilibrate during MD every that many steps
md_return = <str : 1>	'FINAL'	return MODEL with 'MINIMAL' energy or 'FINAL' MODEL
cap_atom_shift = <float : 1>	0.2	limit for atomic shifts in optimization

Output: molpdf

Requirements: restraints

Description: This command performs a number of optimizing iterations using a selected optimization method (6.2). One call to `model.optimize()` corresponds to a single step of the variable target function method. The whole variable target function method is implemented by the `automodel` class. The molecular pdf is optimized with respect to the selected coordinates of the current MODEL; the optimized coordinates are returned as the current MODEL.

Some output may be generated during optimization; for example, a value of the molecular pdf, average and maximal atomic shifts are written to the current tracing file every `trace_output` iterations of the optimizer if `trace_output` is larger than 0 (see the `model.switch_trace()` command).

In addition, a summary of the optimization results is written to the `log` file after optimization, unless `output` contains string 'NO_REPORT'.

`optimization_method = 1` selects a conjugate gradients optimization method. `optimization_method = 3` selects a molecular dynamics optimization at a fixed temperature. The conjugate gradients optimizer is a modified version of the Beale restart conjugate gradients method [Shanno & Phua, 1980, Shanno & Phua, 1982]. The molecular dynamics routine is the most basic version of the iterative solver of the Newton's equations of motion. The integrator uses the Verlet algorithm [Verlet, 1967]. All atomic masses are set to that of carbon 12. A brief description of the algorithms is given in Section 6.2.

`model.schedule.step` is the variable target function step. It selects some of the optimization parameters; it refers to the line in the schedule file which specifies (1) the optimization method (1=Conjugate Gradients, 3=Molecular Dynamics); (2) maximal number of residues that the restraints are allowed to span (Section 4.7.3); (3) the individual scaling factors for all the physical restraint types. `optimization_method` overrides the schedule specification if it is within a defined range.

`edat` is an `energy_data` object containing objective function parameters. See section 4.3 for more information. It is useful in some simulations to be able to set `energy_data.contact_shell` to something large (e.g., 8Å) and `energy_data.update_dynamic` to 999999.9, so that the pairs list is prepared only at the beginning of the optimization. However, you have to make sure that the potential energy is not invisibly pumped into the system by making contacts that are not on the list of non-bonded pairs (see below).

`residue_span_range` determines what atom pairs can possibly occur in the dynamic non-bonded atom pairs list (see `model.restraints.make()`). `residue_span_sign` is ignored in `model.optimize()`. The dynamic restraints include soft-sphere overlap, Lennard-Jones, electrostatic restraints, and general spline restraints. The first three types of restraints can also be generated as static restraints by `model.restraints.make()`.

The automatically generated dynamic restraints are always deleted after a command that calculates them is finished (`model.optimize()`, `model.energy()`, `model.pick_hot_atoms()`); you have to use `model.restraints.make()` to calculate equivalent static restraints if you want to write the 'dynamic' restraints to a file.

`min_atom_shift` is a convergence criterion for the conjugate gradients optimization. When the maximal atomic shift is less than the specified value, the optimization is finished regardless of the number of optimization cycles or function value and its change.

`max_iterations` is used to prevent a waste of CPU time in the conjugate gradients optimization. When that many cycles are done, the optimization is finished regardless of the maximal atomic shift.

The molecular dynamics optimizer pretends that the natural logarithm of the molecular pdf is energy in kcal/mole. `md_time_step` is the time step in femtoseconds. `temperature` is the temperature of the system in degrees Kelvin. `max_iterations` determines the number of MD steps. If `md_return` is 'FINAL' the last structure is returned as the MODEL. If `md_return` is 'MINIMAL' then the structure with the lowest value of the objective function on the whole trajectory is returned as the MODEL. Rescaling of velocities is done every equilibration steps to match the specified temperature. Atomic shifts along one axis are limited by `cap_atom_shift`. This value should be smaller than `energy_data.update_dynamic`. If `init_velocities = True`, the velocity arrays are

initialized, otherwise they are not. In that case, the final velocities from the previous run are used as the initial velocities for the current run.

`molpdf` contains the value of the objective function at the end of optimization.

An exception is raised if optimization is aborted because dynamic restraints could not be calculated as a result of a system being too large. It is up to the calling script to ensure that sensible action is taken; *e.g.*, skipping the rest of modeling for the model that resulted in an impossible function evaluation. This option is useful when calculating several independent models and you do not want one bad model to abort the whole calculation. A probable reason for an interrupted optimization is that it was far from convergence by the time the calculation of dynamic restraints was first requested. Two possible solutions are: (1) optimize more thoroughly (*i.e.* slowly) and (2) use a different contact pairs routine (set `energy_data.nlogn_use = 9999`).

Example: `examples/commands/optimize.py`

```
# Example for: model.optimize(), misc.switch_trace()

# This will optimize stereochemistry of a given model, including
# non-bonded contacts.

env = environ()
env.edat.dynamic_sphere = True

env.libs.topology.read(file='${LIB}/top_heav.lib')
env.libs.parameters.read(file='${LIB}/par.lib')

code = '1fas'
aln = alignment(env)
mdl = model(env, file=code)
aln.append_model(mdl, align_codes=code, atom_files=code)
aln.append_model(mdl, align_codes=code+'-ini', atom_files=code+'.ini')
mdl.generate_topology(aln, sequence=code+'-ini')
mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')
mdl.write(file=code+'.ini')

# Generate the restraints:
mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False)
mdl.restraints.write(file=code+'.rsr')

mpdf = mdl.energy()
mdl.switch_trace(file=code+'.trc')
mdl.optimize(optimization_method=1, max_iterations=20, output='REPORT',
             trace_output=1)
mdl.optimize(optimization_method=3, temperature=300, max_iterations=50,
             output='REPORT', trace_output=1)
mdl.optimize(optimization_method=1, max_iterations=20, output='REPORT',
             trace_output=1)
mpdf = mdl.energy()

mdl.write(file=code+'.B')
```

4.6.36 `model.switch_trace()` — open new optimization trace file

Options:

`file = <str : 1>` 'default' partial or complete filename

Column	Description
1	iteration number within one step of the variable target function method
2	number of function evaluations within one step of VTFM
3	objective function value
4	average atomic shift
5	maximal atomic shift
6	proportional to the gradient
7	kinetic energy
8	temperature for molecular dynamics optimization
9	total energy (kinetic and potential; potential = objective function value)

Table 4.2: *Columns in an optimization trace file.*

```
output_directory = <str : 1>          ''          output directory
```

Description: This command specifies the file for the subsequent optimization tracing output. It is useful for separating tracing output for different models constructed in a single run of MODELLER. The tracing output is only produced if `trace_output` is larger than 0. The tracing file includes the iteration number, number of function evaluations, function value, average and maximal atomic shifts, the size of the gradient vector, kinetic energy (for molecular dynamics ‘optimization’ only), temperature (MD only) and total energy. This is written out in every `trace_output`-th cycle of whatever optimization method is used, starting with the state just before the optimization (iteration 0).

When using the `model` script for comparative modeling, there is one `.D` file for each `.B` file with a model. The `.D` files contain information about the progress of optimization, from the beginning to the end. The most important column is column 3, which contains the value of the objective function, which is being optimized, as a function of the iteration step (every 10 steps, by default). Thus, the best model, according to MODELLER, is the one that has the lowest number in the third column of the last line of its `.D` file. This value is also written out in the `REMARK` record of the PDB file containing the model and in the `log` file.

Example: See `model.optimize()` command.

4.6.37 `model.debug_function()` — test code self-consistency

Options:

```
edat = <energy_data>          objective function parameters
debug_function_cutoff = <float : 3>    0.01 0.001 0.1    cutoffs for reporting differences between numerical and analytical derivatives: absolute, relative errors, factor_for_indiv_rstrs
detailed_debugging = <bool : 1>      False          whether to evaluate energy and derivatives wrt each restraint
```

Description: This command checks the self-consistency of the code for the objective function and its derivatives by calculating and comparing numeric and analytical derivatives. All the parameters influencing the evaluation of the molecular pdf are also relevant (see `model.energy()`). The derivative is reported if both the absolute difference and the fractional difference between the two kinds of evaluations are larger than `debug_function_cutoff[1]` and `debug_function_cutoff[2]`, respectively.

When `detailed_debugging` is `True`, the analytic and numeric derivatives of each restraint with respect to atomic positions are also compared for the atoms ‘violated’ by the whole molecular pdf. The absolute cutoff

for writing out the discrepancies is scaled by `debug_function_cutoff[3]`; the relative cutoff remains the same as before.

When MODELLER is compiled in double precision, this test reports a smaller number of discrepancies.

Example: examples/commands/debug_function.py

```
# Example for: model.debug_function()

# This will use the MODELLER automodel class to construct homology
# restraints for 1fas. It will then use model.debug_function() to test
# the source code for the function and derivatives calculation
# by comparing analytical and numerical first derivatives (note that
# automodel is a derived class of model, so all 'model' methods will work
# on 'automodel'). Some discrepancies may be reported but ignore them here.

from modeller.automodel import *    # Load the automodel class

log.verbose()
env = environ()

a = automodel(env, alnfile = 'debug_function.ali',
              knowns = ('2ctx', '1nbt'), sequence = '1fas')
a.spline_on_site = False
a.make(exit_stage=1)

# To assign 0 weights to restraints whose numerical derivatives
# code does not work (i.e., splines for angles and dihedrals):
a.schedule.read(file='debug_function.sched')
a.energy(output='SHORT')
a.debug_function(debug_function_cutoff=(15.00, 0.10, 0.1),
                detailed_debugging=True)
```

4.6.38 model.assess_ga341() — assess a model with the GA341 method

Output: molpdf

Description: This command assesses the quality of the model using the GA341 method method [Melo *et al.*, 2002]. The method uses the percentage sequence identity between the template and the model as a parameter. MODELLER-produced PDB files contain this information in a 'REMARK'; in the case of other PDB files, you should supply this by setting `model.seq_id`.

If the model contains multiple chains, only the first is evaluated by this method; if you wish to evaluate the model in a different chain, you should write out that single chain into a new model first. (The method was parameterized for use with single-chain proteins, so its use for multi-chain models is not recommended.)

Example: examples/commands/assess_ga341.py

```
# Example for: model.assess_ga341()

env = environ()

# Read a model previously generated by Modeller's automodel class
mdl = model(env, file='../atom_files/1fdx.B99990001.pdb')
```

```
# Set template-model sequence identity. (Not needed in this case, since
# this is written by Modeller into the .pdb file.)
mdl.seq_id = 37.037

score = mdl.assess_ga341()
```

4.6.39 `model.assess_dope()` — assess a model with the DOPE method

Output: molpdf

Description: This command assesses the quality of the model using the DOPE (Discrete Optimized Protein Energy) method. This is a statistical potential optimized for model assessment. As with `model.assess_ga341()`, this should be used for single chains only.

DOPE uses the standard MODELLER energy function, so any of the arguments accepted by `model.energy()` can also be used here. (Note also that the model's topology must be set up in order to calculate the energy - see the example for `model.energy()`.)

Only the DOPE energy itself is returned by this command (all other components of the MODELLER energy function, such as stereochemical restraints, Lennard-Jones interactions, homology-derived restraints, *etc.* are ignored) unless you manually set `schedule_scale`. See `model.energy()` for more details.

The DOPE model score is designed for selecting the best structure from a collection of models built by MODELLER. (For example, you could build multiple `automodel` models by setting `automodel.ending_model`, and select the model that returns the lowest DOPE score.) Currently the score is unnormalized with respect to the protein size and has an arbitrary scale, therefore scores from different proteins cannot be compared directly. A normalization scheme is under development.

4.7 Calculation of spatial restraints

This Chapter explains how the restraints are represented in a restraint file and also describes commands for reading, writing, generating, and manipulating restraints. See Section 6.3 for equations defining the restraints and their derivatives with respect to atomic positions. See Section 4.6 for commands for calculating the objective function and Section 6.2 for optimization methods. See the original papers for the most detailed definition and description of the restraints [Šali & Blundell, 1993, Šali & Overington, 1994].

4.7.1 Specification of restraints

Static and dynamic restraints

Static restraints are read from the restraints file or are generated by the `model.restraints.make()` command. All other restraints are *dynamic* restraints and are created on the fly; they currently include restraints on non-bonded atom pairs.

Formats of the restraints file

Restraints may be read from a restraints file in two formats, MODELLER or USER. The files in the MODELLER and USER formats have to begin with the lines 'MODELLER5 VERSION: MODELLER FORMAT' and 'MODELLER5 VERSION: USER FORMAT', respectively. In both formats, there is one entry per line. The format is free, except that the first character has to be at the beginning of the line. There are three different entry types in the MODELLER format:

```
R Form Modality Feature Group Numb_atoms Numb_parameters 0 Atom_indices Parameters
E Atom_index_1 Atom_index_2
P Pseudo_atom_index Pseudo_atom_type Numb_real_atoms Real_atom_indices
```

For example,

```
R 3 1 1 1 2 2 0 437 28 1.5000 0.1000
E 120 540
P 1 3 3 120 121 122
```

When the line starts with 'R', it contains a restraint, 'E' indicates a pair of atoms to be excluded from the calculation of the dynamic non-bonded pairs list, and 'P' indicates a pseudo atom definition (Section 4.7.2).

The USER format recognizes only the R entries. The fields of a line in the USER format are:

```
Id Form Modality Feature Group Numb_atoms Numb_parameters 0 Parameters Atom_ids
```

(Note that `Parameters` and `Atom_ids` are in opposite orders in the USER and MODELLER formats.)

For example,

```
R 3 1 1 1 2 2 0 1.5000 0.1000 NH#:1:A CA:2:A
```

The seven integer indices used to specify various restraint properties are listed in Tables 4.3–4.5. They are: `Form` specifies the mathematical form of the restraint. `Modality` should be viewed as the argument to `Form`. It specifies the number of single Gaussians in a poly-Gaussian pdf, periodicity n of the cosine in the cosine potential, and the number of spline points for cubic splines. Only certain combinations of `Form` and `Modality` are possible. Any `Feature` can be used with any `Form/Modality` pair. `Group` or “physical feature type” groups restraints for reporting purposes in `model.energy()`, *etc.* The number of atoms and parameters for the restraint are specified by `Numb_atoms` and `Numb_prms`, respectively. The seventh integer index can be ignored. `Atom_indices` and `Parameters` have to match the hard-wired conventions. The format of the atom id is `ATOM_NAME:RESIDUE_#[:CHAIN_ID]`, where `ATOM_NAME` is the four character IUPAC atom name as found in a PDB file, `RESIDUE_#` is a five character residue number as it occurs in the PDB file of a model, and the optional `CHAIN_ID` is the single character chain id as it occurs in the PDB file. For example, the carbonyl oxygen (O) in residue '10A' in chain 'A' is specified by 'O:10A:A'; if the chain has no chain id, the name would be only 'O:10A'.

4.7.2 Specification of pseudo atoms

There are virtual and pseudo atoms. A virtual atom is an atom that occurs in the actual molecule, but whose position is not represented explicitly in the MODEL and topology file. A pseudo atom is a position that does not correspond to an actual atom in a molecule, but is some sort of an average of positions of real atoms. MODELLER follows GROMOS definitions for the seven types of pseudo and virtual atoms: gravity center, V41, V31, P2, V42, P3, and P6. These names are constructed using the following rules: 'V' and 'P' indicate virtual and pseudo atoms, respectively. The second digit indicates the number of substituents on the central atom (for 'V') and the number of protons whose positions are averaged (for 'P'). The last digit indicates the number of protons on the central atom (for 'V').

GROMOS TYPE	ROUTINE NAME	#DEF ATM	DESCRIPTION
1	PSD	N	gravity center
2	VCH1	4	virtual aliphatic proton on a tetrahedral carbon (->CH), defined by the central C and the three other substituents;
3	VCH1A	3	virtual aromatic proton on a trigonal carbon (=CH), defined by the central C and the two C atoms bonded to the central C;
4	PCH2	3	pseudo aliphatic proton on a tetrahedral carbon (>CH2) not assigned stereospecifically; its position is between the two real protons; defined by the central C and the other two substituents;
5	VCH2	3	virtual aliphatic proton on a tetrahedral carbon (>CH2) assigned stereospecifically; defined by the central tetrahedral atom and the other two substituents on it;
6	PCH31	2	pseudo aliphatic proton on a tetrahedral carbon (-CH3), defined by the central C and the heavy atom X in X-CH3; its position is the average of the three real protons;
7	PCH32	3	pseudo aliphatic proton between two unassigned -CH3 groups; defined by X in CH3 - X - CH3 and the two C atoms from the two CH3 groups (Val, Leu!); its position is the average of the six real protons;
0	-	-	delta and epsilon protons on rapidly flipping aromatic rings should refer directly to real gamma and delta C atoms, respectively.

In a restraints file, pseudo atoms are indexed from NATM+1 to NATM+NPSEUDO where NPSEUDO is the number of pseudo atoms. The restraints (the R entries) are exactly the same as for the real atoms, except that the pseudo atom integer indices are used (indices are larger than NATM). The pseudo atoms are defined in the P entries:

```
P i j k a1 a2 a3
```

where *i* is atom index of pseudo atom *i*, *j* is the type of the pseudo atom *i* (see the table above), *k* is the number of real atoms defining the current pseudo atom (3 in this case), and *a1 a2 a3* are the integer indices of real atoms defining the current pseudo atom.

For example, if you want to define a pseudo atom which is a gravity center of atoms 4, 7, and 10, and there are 101 real atoms in the protein:

P 102 1 3 4 7 10

#	Form	Parameters	Violation	Reference
1	left Gaussian (harmonic lower bound)	f, σ	$(f - f)/\sigma$	Eq. 6.55
2	right Gaussian (harmonic upper bound)	\bar{f}, σ	$(f - \bar{f})/\sigma$	Eq. 6.56
3	single Gaussian (harmonic potential)	\bar{f}, σ	$(f - \bar{f})/\sigma$	Eq. 6.39
4	multiple Gaussian	$(\omega_i)_n, (\bar{f}_i)_n, (\sigma_i)_n$	$\max_{\omega_i} (f - \bar{f}_i)/\sigma_i$	Eq. 6.41
5	Lennard-Jones potential	A, B	0.0	Eq. 6.63
6	Coulomb point-to-point potential	q_1, q_2	0.0	Eq. 6.60
7	Cosine potential	a, b	c	Eq. 6.57
8	undefined			
9	multiple binormal	$(\omega_i)_n, (\bar{f}_{1i}, \bar{f}_{2i})_n, (\sigma_{1i}, \sigma_{2i})_n, (\rho)_n$	$\max_{\omega_i} \sqrt{-\frac{1}{2(1-\rho_i^2)} \left[\left(\frac{f_1 - \bar{f}_{1i}}{\sigma_{1i}} \right)^2 - 2\rho_i \frac{f_1 - \bar{f}_{1i}}{\sigma_{1i}} \frac{f_2 - \bar{f}_{2i}}{\sigma_{2i}} + \left(\frac{f_2 - \bar{f}_{2i}}{\sigma_{2i}} \right)^2 \right]}$	Eq. 6.51
10	cubic spline	p_i , for $i = 1, 6 + n$	$(f - f_{min})/\sigma$	Eq. 6.70

Table 4.3: *List of mathematical forms of restraints.* The parameters and their order in the restraint file are also given (**Params** above). $(\dots)_n$ indicates that (\dots) is repeated n times, where n is specified by the second integer parameter of the restraint, **modality** (see above). **Modality** also defines periodicity of the cosine restraint, corresponding to parameter n in Eq 6.57, and the number of interpolating points for the spline restraint (Eq. 6.70). Feature f can generally be either a measure of solvent exposure (undocumented), a distance, an angle, or a dihedral angle, with the exception of restraint form 9 that only works with a pair of dihedral angles. The angle unit in the restraints file is radians. The internal angle unit of MODELLER is radians, too. Column 'Violation' defines the "relative heavy violations" used in **model.pick_hot_atoms()**. For cubic splines, f_{min} is the feature value that results in the smallest value of the restraint and σ is the standard deviation of the Gaussian function fitted locally around f_{min} . The parameters p_i for a spline restraint are: the scaling factor (p_1), the smallest value at which interpolation is done, x_1 (p_2), the largest interpolating value x_n (p_3), the interval between interpolating points, Δx (p_4), the first derivative at x_1 (p_5), the first derivative at x_n (p_6). The following n values are the values of the restraint at the interpolating x_i points. The MODELLER-4 format has additional n values, which are the second derivatives of the restraint at the interpolating x_i points.

Index	Feature
1	distance
2	angle
3	dihedral angle
4	a pair of dihedral angles (points 1–4 and 5–8)
5	distance between gravity centers of two groups of atoms
6	minimal distance between several pairs of atoms
7	atomic area exposed to solvent in \AA^2
8	atomic density (number of atoms within <code>contact_shell</code>)
9	x coordinate
10	y coordinate
11	z coordinate
12	difference between two dihedral angles (1–4 and 5–8)

Table 4.4: *List of feature types that can be restrained.*

Index	Group
1	Bond length potential
2	Bond angle potential
3	Stereochemical cosine dihedral potential
4	Stereochemical improper dihedral potential
5	soft-sphere overlap restraints
6	Lennard-Jones 6-12 potential
7	Coulomb point-point electrostatic potential
8	H-bonding potential
9	Distance restraints 1 (C_{α} - C_{α})
10	Distance restraints 2 (N-O)
11	Mainchain Φ dihedral restraints
12	Mainchain Ψ dihedral restraints
13	Mainchain ω dihedral restraints
14	Sidechain χ_1 dihedral restraints
15	Sidechain χ_2 dihedral restraints
16	Sidechain χ_3 dihedral restraints
17	Sidechain χ_4 dihedral restraints
18	Disulfide distance restraints
19	Disulfide angle restraints
20	Disulfide dihedral angle restraints
21	X lower bound distance restraints
22	X upper bound distance restraints
23	Distance restraints 3 (SDCH-MNCH)
24	Sidechain χ_5 dihedral restraints
25	(Φ , Ψ) binomial dihedral restraints
26	Distance restraints 4 (SDCH-SDCH)
27	Distance restraints 5 (X-Y)
28	NMR distance restraints 6 (X-Y)
29	NMR distance restraints 7 (X-Y)
30	Minimal distance restraints
31	Non-bonded spline restraints
32	Atomic accessibility restraints
33	Atom density restraints
34	Absolute position restraints
35	Dihedral angle difference restraints

Table 4.5: *List of “physical” restraint types.*

4.7.3 model.restraints.make() — make restraints

Options:

edat = \langle energy_data \rangle		objective function parameters
aln = \langle alignment \rangle		Template-model alignment; for homology-derived restraints only
io = \langle io_data \rangle		Options for reading atom files
restraint_type = \langle str : 1 \rangle	'STEREO'	restraint type to be calculated: 'STEREO' 'BOND' 'ANGLE' 'IMPROPER' 'DIHEDRAL' 'SPHERE' 'SPHERE14' 'LJ' 'LJ14' 'COULOMB' 'COULOMB14' 'ALPHA' 'STRAND' 'SHEET' 'DISTANCE' 'USER_DISTANCE' 'NONB_PAIR_SPLINE' 'PHI-PSI_BINORMAL' 'PHI_DIHEDRAL' 'PSI_DIHEDRAL' 'OMEGA_DIHEDRAL' 'CHI1_DIHEDRAL' 'CHI2_DIHEDRAL' 'CHI3_DIHEDRAL' 'CHI4_DIHEDRAL'
dih_lib_only = \langle bool : 1 \rangle	False	whether to use only library, not homologs for dihedral angle rsrs
mnch_lib = \langle int : 1 \rangle	1	which MNCH lib to use in MAKE_RESTRAINTS
intersegment = \langle bool : 1 \rangle	True	whether to restrain inter-segment non-bonded pairs
residue_grouping = \langle int : 1 \rangle	1	
maximal_distance = \langle float : 1 \rangle	999	maximal distance for distance restraints
residue_span_range = \langle int : 2 \rangle	0 99999	range of residues spanning the allowed distances; for MAKE_RESTRAINTS, PICK_RESTRAINTS, non-bonded dynamic pairs
residue_span_sign = \langle bool : 1 \rangle	True	whether to do $N*(N-1)/2$ loop for atom pairs in MAKE_RESTRAINTS
restraint_sel_atoms = \langle int : 1 \rangle	1	RESTRAINT_TYPE = 'distance' a restraint other than non-bonded pair has to have at least as many selected atoms
accessibility_type = \langle int : 1 \rangle	8	type of solvent accessibility: 1–10
distance_rsr_model = \langle int : 1 \rangle	1	the model for calculating distance restraints: 1–7
restraint_group = \langle int : 1 \rangle	26	physical restraint group
restraint_stdev = \langle float : 2 \rangle	0.0 1.0	transforming factors for standard deviations ($y=a+bx$) in models 1–6 or standard deviation for model 7 (a)
restraint_stdev2 = \langle float : 3 \rangle	0 0 0	transforming standard deviation in models 3–6: $S' = S + [a + b \max(0, c-g)]$
restraint_parameters = \langle float : 0 \rangle	3 1 3 3 4 2 0 0.0 0.087	restraint parameters for 'USER_DISTANCE'
basis_pdf_weight = \langle str : 1 \rangle	'LOCAL'	a method for calculation of basis pdf weights: 'LOCAL' 'GLOBAL'
basis_relative_weight = \langle float : 1 \rangle	0.05	the cutoff weight of basis pdf's for their removal
residue_ids = \langle str : 0 \rangle	''	residue id (number:chnid)
spline_on_site = \langle bool : 1 \rangle	False	whether to convert restraints to splines

<code>spline_dx = <float : 1></code>	0.5	interval size for splining restraints
<code>spline_min_points = <int : 1></code>	5	have at least as many intervals in a spline
<code>spline_range = <float : 1></code>	4.0	range of the splines
<code>sheet_h_bonds = <int : 1></code>	7	specify hydrogen bonds in a beta-sheet

Requirements: topology & parameters [`&` picked atoms sets 2 and 3]

Description: This command calculates and selects new restraints of a specified type. See the original papers for the most detailed definition and description of the restraints [Šali & Blundell, 1993, Šali & Overington, 1994]. The calculation of restraints of all types is now (partly) limited to the selected atoms only (either set 1, or 2 and 3; see below). The new restraints are added to any currently present.

`restraint_type` selects the types of the generated restraints. Only one restraint type can be selected at a time, except for the stereochemical restraints (`BOND`, `ANGLE`, `DIHEDRAL`, `IMPROPER`) that can all be calculated at the same time. It is useful to distinguish between the stereochemical restraints and homology-derived restraints. The stereochemical restraints are obtained from libraries that depend on atom and/or residue types only (*e.g.*, CHARMM 22 force field [MacKerell *et al.*, 1998] or statistical potentials), and do not require the alignment `aln` with template structures. In contrast, the homology-derived restraints are calculated from related protein structures, which correspond to all but the last sequence in the alignment `aln` (the target). These templates are read from coordinate files, which are the only data files required. All restraints are added to the existing restraints, even if they duplicate them (but see the comment for the 'OMEGA' restraints below).

Stereochemical restraints:

- '`BOND`'. This calculates covalent bond restraints (harmonic terms). It relies on the list of the atom-atom bonds for MODEL, prepared previously by the `model.generate_topology()` command. The mean values and force constants are obtained from the parameter library in memory. Only those bonds are restrained that have all or at least `restraint_sel_atoms` in the selected atom set 1.
- '`ANGLE`'. This calculates covalent angle restraints (harmonic terms). It relies on the list of the atom-atom-atom bonds for MODEL, prepared previously by the `model.generate_topology()` command. The mean values and force constants are obtained from the parameter library in memory. Only those angles are restrained that have all or at least `restraint_sel_atoms` in the selected atom set 1.
- '`DIHEDRAL`'. This calculates covalent dihedral angle restraints (cosine terms). It relies on the list of the atom-atom-atom-atom dihedral angles for MODEL, prepared previously by the `model.generate_topology()` command. The minima, phases, and force constants are obtained from the parameter library in memory. Only those dihedral angles are restrained that have all or at least `restraint_sel_atoms` in the selected atom set 1.
- '`IMPROPER`'. This calculates improper dihedral angle restraints (harmonic terms). It relies on the list of the improper dihedral angles for MODEL, prepared previously by the `model.generate_topology()` command. The mean values and force constants are obtained from the parameter library in memory. Only those improper angles are restrained that have all or at least `restraint_sel_atoms` in the selected atom set 1.
- '`STEREO`'. This implies all '`BOND`', '`ANGLE`', '`DIHEDRAL`', and '`IMPROPER`' restraints.
- '`SPHERE14`'. This constructs soft-sphere overlap restraints (lower harmonic bounds) for atom pairs separated by exactly three bonds (1-4 pairs). It relies on atom radii from the '`$RADII14_LIB`' library. Only those non-bonded pairs are restrained that have all or at least `energy_data.nonbonded_sel_atoms` in the selected atom set 1. They must also satisfy the `residue_span_range` & `residue_span_sign` criterion.
- '`LJ14`'. This constructs 1-4 Lennard-Jones restraints using the modified 1-4 Lennard-Jones parameters from the CHARMM parameter library. There is no way to calculate '`LJ14`' as dynamic restraints. Only those non-bonded pairs are restrained that have all or at least `energy_data.nonbonded_sel_atoms` in the selected atom set 1. They must also satisfy the `residue_span_range` & `residue_span_sign` criterion.
- '`COULOMB14`'. This constructs 1-4 Coulomb restraints by relying on the atomic charges from the CHARMM topology library. There is no way to calculate '`COULOMB14`' as dynamic restraints. Only those non-bonded pairs are restrained that have all or at least `energy_data.nonbonded_sel_atoms` in the selected atom set 1. They must also satisfy the `residue_span_range` & `residue_span_sign` criterion.

- **'SPHERE'**. This constructs soft-sphere overlap restraints (lower harmonic bounds) for all atom pairs that are not in bonds, angles, dihedral angles, improper dihedral angles, nor are explicitly excluded by the **'E'** entries read from a restraint file or added by the **model.restraints.add()** command. Only those non-bonded pairs are restrained that have all or at least **energy_data.nonbonded_sel_atoms** in the selected atom set 1. They must also satisfy the **residue_span_range** & **residue_span_sign** criterion. Note that this makes these restraints static (*i.e.*, not dynamic) and that you must set **energy_data.dynamic_sphere** to **False** before evaluating the molecular pdf if you want to avoid duplicated restraints. These restraints should usually not be combined with the Lennard-Jones (**'LJ'**) restraints.
When **intersegment** is **True**, the inter-segment non-bonded restraints are also constructed; otherwise, the segments do not feel each other *via* the non-bonded restraints. This option does not apply to the **model.optimize()** command where information about segments is not used at all (*i.e.*, **model.optimize()** behaves as if **intersegment = True**).
- **'LJ'**. This constructs Lennard-Jones restraints for all atom pairs that are not in bonds, angles, dihedral angles, improper dihedral angles, nor are explicitly excluded by the **'E'** entries read from a restraint file or added by the **model.restraints.add()** command. Only those non-bonded pairs are restrained that have all or at least **energy_data.nonbonded_sel_atoms** in the selected atom set 1. They must also satisfy the **residue_span_range** & **residue_span_sign** criterion. Note that this command makes the non-bonded restraints static (*i.e.*, not dynamic) and that you must set **energy_data.dynamic_lennard** to **False** before evaluating the molecular pdf if you want to avoid duplicated restraints. Note that CHARMM uses both **'LJ14'** and **'LJ'**. For large molecules, it is better to calculate **'LJ'** as dynamic restraints because you can use distance cutoff **energy_data.contact_shell** in **model.optimize()** to reduce significantly the number of non-bonded atom pairs.
- **'COULOMB'**. This constructs Coulomb restraints for all atom pairs that are not in bonds, angles, dihedral angles, improper dihedral angles, nor are explicitly excluded by the **'E'** entries read from a restraint file or added by the **model.restraints.add()** command. Only those non-bonded pairs are restrained that have all or at least **energy_data.nonbonded_sel_atoms** in the selected atom set 1. They must also satisfy the **residue_span_range** & **residue_span_sign** criterion. Note that this command makes the non-bonded restraints static (*i.e.*, not dynamic) and that you must set **energy_data.dynamic_coulomb** to **False** before evaluating the molecular pdf if you want to avoid duplicated restraints. Note that CHARMM uses both **'COULOMB14'** and **'COULOMB'**. For large molecules, it is better to calculate **'COULOMB'** as dynamic restraints because you can use distance cutoff **energy_data.contact_shell** in **model.optimize()** to reduce significantly the number of non-bonded atom pairs.
- **'ALPHA'**. This makes restraints enforcing an α -helix (mainchain conformation class "A") for the residue segment specified by the two **residue_ids** (Section 4.9.1). The helix is restrained by Φ, Ψ binormal restraints, N-O hydrogen bonds, $C_\alpha-C_\alpha$ distances for $i-j \in \{2-9\}$, $C_\alpha-O$ distances for $i-j \in \{2-9\}$, and O-O distances for $i-j \in \{2-6\}$. These target distances were all obtained from a regular α -helix in one of the high-resolution myoglobin structures. A convenient way to add **'ALPHA'**, **'STRAND'**, or **'SHEET'** restraints to the calculation by the **'model'** script is to include them in the **special_restraints** routine (Section 1.8, Question 12). Note that at least the non-hydrogen mainchain atoms topology model is required although the same functionality could also be provided for the C_α -only topology with small changes to the source code.
- **'STRAND'**. This makes restraints enforcing an extended strand conformation for the residue segment specified by the two **residue_ids** (Section 4.9.1). This is achieved by applying Φ, Ψ binormal restraints only. These binormal restraints force the mainchain conformation into class "B", except for the Pro residues which are restrained to class "P" [Šali & Blundell, 1993].
- **'SHEET'**. This calculates H-bonding restraints for a pair of β -strands. **atom_ids** specifies the two atom identifiers (Section 4.7.1) defining the first H-bond in the β -sheet ladder. **sheet_h_bonds** specifies the number of H-bonds to be added. The parallel and anti-parallel sheets are selected by a positive and negative integer in **sheet_h_bonds**, respectively. In a parallel sheet, hydrogen bonds start at the first or the second term of the following series (depending on **atom_ids**): 1N:1O, 1O:3N, 3N:3O, 3O:5N, *etc.* For an anti-parallel sheet, the corresponding series is 1N:3O, 1O:3N, 3N:1O, 3O:1N, *etc.*; note that the residue indices are always decreasing for the second strand. The extended structure of the individual strands must be enforced separately by the **'STRAND'** restraints if so desired.
- **'USER_DISTANCE'**. This makes distance restraints between pairs of atoms from set 2 and 3 (inter-set

only), using the value of `restraint_parameters`. Only distances satisfying the `residue_span_range` criterion are restrained. This command is useful for making non-specific “compactization” restraints.

Homology-derived restraints:

For these restraints, the input alignment `aln` must be given.

- ‘DISTANCE’. This makes distance restraints that are generated for all pairs of atoms i, j where atom i is from selected set 2 and atom j is from selected set 3 (as defined by the `model.pick_atoms()` command). The atoms also have to be within the residue spanning range specified by `residue_span_range = r1 r2`, such that the residue index difference $r1 \leq |ir2 - ir1| \leq r2$ when `residue_span_sign = False` and $r1 \leq (ir2 - ir1) \leq r2$ when `residue_span_sign = True`. Moreover, for a restraint to be created, at least one distance in the template structures must be less than `maximal_distance` (in Å). The mean of this basis pdf is equal to the template distance and its standard deviation σ is calculated from an analytic model specified by `distance_rsr_model`. Use model 5 for C_α - C_α distances and model 6 for N-O distances. For models 1 through 6, this standard deviation is transformed by $\sigma' = a + b * (\sigma + W_g)$ where a and b are given by `restraint_stdev` and W_g is a gap weighting function of the form $W_g = 0.6 * \max(0, 4 - g_{av})$. g_{av} is the average distance of the two residues involved in the restraint from a gap. For models 3 through 6, this is additionally transformed by $\sigma'' = \sigma' + \sum_i [d + e * \max(0, f - g_i)]$ where the sum is over each of the atoms i involved in the distance, d e and f are given by `restraint_stdev2`, and g_i is the distance of each residue from a gap. The first six models are polynomials and depend on several structural features of the template and its similarity to the target. The polynomial coefficients are specified in library file ‘\$PARAMS_LIB’. When “polynomial model” 7 is selected, the standard deviation of restraints is set to constant a . Each basis pdf in the distance pdf corresponds to one template structure with an equivalent distance. The weights of basis pdf’s depend on local sequence similarity between the target and the templates when `basis_pdf_weight = 'LOCAL'` and on global sequence identity when `basis_pdf_weight = 'GLOBAL'`. In addition, the atom pairs restrained by homology-derived restraints must by default not be in a chemical bond, chemical angle, dihedral angle, or on an excluded pairs list. This behavior can be changed by resetting `energy_data.excl_local` (see `model.optimize()`).
- ‘CHI1_DIHEDRAL’, ‘CHI2_DIHEDRAL’, ‘CHI3_DIHEDRAL’, ‘CHI4_DIHEDRAL’, ‘PHI_DIHEDRAL’, ‘PSI_DIHEDRAL’, ‘OMEGA_DIHEDRAL’, ‘PHI-PSI_BINORMAL’ are the mainchain and sidechain dihedral angle restraints. Only those dihedral angles are restrained that have all or at least `energy_data.nonbonded_sel_atoms` in the selected atom set 1. The means and standard deviations for the dihedral Gaussian restraints are obtained from the \$RESDIH_LIB and \$MNCH?_LIB libraries and their weights from the MDT tables, which are read in as specified by `MDT_LIB` in `$LIB/libs.lib`. The large MDT tables give the conditional weights for each possible dihedral angle class, as a function of all possible combinations of features on which a particular class depends. If `dih_lib_only` is ON or there is no equivalent residue in any of the templates, the weights for the dihedral angle classes depend only on the residue type and are obtained from the ‘\$RESDIH_LIB’ and ‘\$MNCH?_LIB’ libraries; the `dih_lib_only` argument allows one to force the calculation of the “homology-derived” mainchain and sidechain dihedral angle restraints that ignore template information. `basis_pdf_weight` has the same effect as for the distance pdf’s. When MODELLER’s ‘OMEGA’ restraints are calculated, the currently existing restraints on atoms ‘O C +N +CA’ in all residues are automatically deleted. These deleted restraints correspond to the *improper* dihedral angles involving the ω atoms. They are deleted because they could be “frustrated” by the new ‘OMEGA’ restraints. No action is taken with regard to any of the previously existing, possibly duplicated dihedral angle restraints. Thus, to avoid restraint duplication, including that of the ‘OMEGA’ restraints, call the `model.restraints.condense()` command after all the restraints are calculated.

`basis_relative_weight` is the cutoff for removing weak basis pdf’s from poly-Gaussian feature pdf’s: a basis pdf whose weight is less than the `basis_relative_weight` fraction of the largest weight is deleted.

Example: `examples/commands/make_restraints.py`

```
# Example for: model.restraints.make(), model.restraints.spline(),
#               model.restraints.write()
```

```

# This will compare energies of bond length restraints expressed
# by harmonic potential and by cubic spline.

log.verbose()
env = environ()
env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

code = '1fas'
aln = alignment(env)
mdl = model(env, file=code, model_segment=('1:', '61:'))
aln.append_model(mdl, atom_files=code, align_codes=code)
aln.append_model(mdl, atom_files=code+'.ini', align_codes=code+'-ini')
mdl.generate_topology(aln, sequence=code+'-ini')
mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')
mdl.write(file=code+'.ini')

mdl.restraints.make(aln, restraint_type='bond', spline_on_site=False)
mdl.restraints.write(file=code+'-1.rsr')
edat = energy_data(dynamic_sphere=False)
mdl.energy(edat=edat)

mdl.restraints.spline(spline_range=5.0, spline_dx=0.005,
                     spline_select=(3, 1, 1), edat=edat)
mdl.restraints.condense()
mdl.restraints.write(file=code+'-2.rsr')
mdl.energy(edat=edat)

```

4.7.4 `model.symmetry.define()` — define similar segments

Options:

<code>symmetry_weight = <float : 1></code>	1.0	the weight of the symmetry objective function term
<code>add_symmetry = <bool : 2></code>	False True	whether to add segment pair, add atoms to segment pair

Description: This command allows defining pairs of segments that will be restrained to be the same during optimization of the objective function. This is achieved by adding the sum of squares of the differences between the equivalent distances (similar to distance RMS deviation) to the objective function being optimized, separately for each pair of segments defined by `model.symmetry.define()`. The value of this term is reported in the log file by the `model.energy()` command, which also reports the individual contributions to the term when output contains word 'SYMMETRY'. In each call of the `model.symmetry.define()` command, the list of such segments is either initiated, extended by a new pair of segments, or the last defined pair of segments is extended by adding new atoms.

`symmetry_weight` specifies the atomic weights to be used in the calculation of the symmetry term (Eq. 6.72).

The two segments correspond to the selected sets 2 and 3 (obtained by the `model.restraints.pick()` command). They must have the same number of atoms.

A pair of segments can be either added to the list (`add_symmetry[0] = True`) or the list can be initialized (`add_symmetry[0] = False`).

If `add_symmetry[1] = True`, the currently selected atoms are added to the last segment pair in the segment pairs list, otherwise a new segment pair is started.

Example: examples/commands/define_symmetry.py

```

# Example for: model.symmetry.define()

# This will force two copies of 1fas to have similar mainchain
# conformation.

log.level(1, 1, 1, 1, 0)
env = environ()
env.libs.topology.read(file='${LIB}/top_heav.lib')
env.libs.parameters.read(file='${LIB}/par.lib')

def defsym(mdl, aln, seg1, seg2):
    for (set, seg) in [(2, seg1), (3, seg2)]:
        mdl.pick_atoms(aln, pick_atoms_set=set, selection_segment=seg,
                       atom_types='MNCH', selection_status='INITIALIZE',
                       selection_search='SEGMENT')
        mdl.symmetry.define(symmetry_weight=1.0, add_symmetry=(True, False))

# Generate two copies of a segment:
code = '2abx'
mdl = model(env, file=code, model_segment=('1:A', '74:B'))
aln = alignment(env)
aln.append_model(mdl, align_codes=code, atom_files=code)
aln.append_model(mdl, align_codes=code+'_ini', atom_files=code+'_ini')
mdl.generate_topology(aln, sequence=code+'_ini')
mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')
mdl.rename_segments(segment_ids=('A', 'B'), renumber_residues=(1, 1))

myedat = energy_data(dynamic_sphere = False)
mdl.energy(edat=myedat)
mdl.randomize_xyz(deviation=6.0)
# Define the two segments (chains in this case) to be identical:
defsym(mdl, aln, seg1=('1:A', '74:A'), seg2=('1:B', '74:B'))

# Make them identical by optimizing the initial randomized structure
# without any other restraints:
mdl.energy(edat=myedat)
mdl.write(file='define_symmetry-1.atm')
mdl.optimize(max_iterations=300, edat=myedat)
mdl.write(file='define_symmetry-2.atm')
mdl.energy(edat=myedat)

# Now optimize with stereochemical restraints so that the
# result is not so distorted a structure (still distorted
# because optimization is not thorough):
myedat.dynamic_sphere = True
mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False,
                   edat=myedat)
mdl.randomize_xyz(deviation=3.0)
for method in (1, 3, 1): # 1 = conjugate gradients, 3 = molecular dynamics
    mdl.optimize(max_iterations=300, md_return='FINAL', edat=myedat,
                 optimization_method=method, output='REPORT')
mdl.write(file='define_symmetry-3.atm')
mdl.energy(edat=myedat)

```

```
# Create a blank alignment so that superpose uses its 1:1 default
aln = alignment(env)

mdl = model(env, file='define_symmetry-3.atm', model_segment=('1:A', '74:A'))
mdl2 = model(env, file='define_symmetry-3.atm', model_segment=('1:B', '74:B'))
mdl.pick_atoms(aln, atom_types='MNCH')
mdl.superpose(mdl2, aln)
```

4.7.5 `model.restraints.unpick_all()` — unselect all restraints

Description: This unselects all of the current restraints.

4.7.6 `model.restraints.clear()` — delete all restraints

Description: This deletes all of the current restraints.

4.7.7 `model.restraints.pick()` — pick restraints for selected atoms

Options:

<code>residue_span_range = <int : 2></code>	0 99999	range of residues spanning the allowed distances; for MAKE_RESTRAINTS, PICK_RESTRAINTS, non-bonded dynamic pairs
<code>restraints_filter = <float : 35></code>	999 999 999 999 999 999 999 999 999 999 999 999 999 999	keep restraints?
<code>restraint_sel_atoms = <int : 1></code>	1	a restraint has to have at least as many selected atoms

Description: This command selects some or all of the restraints currently in memory.

The selection is added to any existing selected atoms; if instead you want to select only these atoms, call `model.restraints.unpick_all()` first.

This command runs over all restraints in memory, including the currently unselected restraints. Be careful about this: If you have some unselected restraints in memory, `model.restraints.pick()` may select them; to prevent this, do `model.restraints.condense()` before calling `model.restraints.pick()`.

A static restraint is selected if all or at least `restraint_sel_atoms` of its atoms are selected (set 1), if it is strong enough based on its standard deviations or force constants (see the next paragraph), and if it does not span fewer residues than `residue_span_range[0]`, or more than `residue_span_range[1]`. (If `residue_span_range` is set to -999 -999, then the residue span range is instead taken to be between 0 and the number given in the third column of the current schedule; see `model.schedule.make()`. Restraints which act on only a single atom are not subject to this range check.) Note that here `restraint_sel_atoms` is used for *all* restraints, while the `model.restraints.make()` and `model.optimize()` commands use it for all restraint types *except* non-bonded pairs. (`energy_data.nonbonded_sel_atoms` is used for non-bonded pairs by these routines.)

To decide if a restraint is strong enough, the current standard deviations or force constants are compared with the corresponding `restraints_filter[physical_restraint_type]`. A harmonic restraint, lower and upper bounds, and multi-modal Gaussian restraints are selected if the (smallest) standard deviation is less than the corresponding `restraints_filter[i]`. The cosine energy term is selected if its force constant is larger than the corresponding

restraints_filter[i]. If restraints_filter[i] = -999, a restraint of type *i* is always selected. Restraints of the other physical_restraint_types are always selected (Coulomb, Lennard-Jones, binormal, and spline). The restraints_filter angles have to be specified in radians.

Example: examples/commands/pick_restraints.py

```
# Example for: model.restraints.pick(), model.restraints.condense()

# This will pick only restraints that include at least one
# CA atom and write them to a file.

log.verbose()
env = environ()
env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

code = '1fas'
mdl = model(env, file=code)
aln = alignment(env)
aln.append_model(mdl, atom_files=code, align_codes=code)
aln.append_model(mdl, atom_files=code+'.ini', align_codes=code+'.ini')
mdl.generate_topology(aln, sequence=code+'.ini')
mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')

mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False)
mdl.energy()

mdl.pick_atoms(aln, atom_types='CA N C O')
mdl.restraints.pick(add_restraints=False, restraint_sel_atoms=1)
# Delete the unselected restraints from memory:
mdl.restraints.condense()
mdl.energy()

mdl.restraints.write(file='1fas.rsr')
```

4.7.8 model.restraints.condense() — remove unselected restraints

Description: This command removes all the unselected restraints from memory. In addition, it also removes those cosine dihedral angle restraints (restraint_type = 'DIHEDRAL') that operate on the same atoms as any other restraints on a dihedral angle or a pair of dihedral angles. Such restraints include the MODELLER 'PHI_DIHEDRAL', 'PSI_DIHEDRAL', 'OMEGA_DIHEDRAL', 'CHI1_DIHEDRAL', 'CHI2_DIHEDRAL', 'CHI3_DIHEDRAL', 'CHI4_DIHEDRAL', 'PHI_PSI_CLASS', 'MRFP_DIHEDRAL', and 'PHI_PSI_BINORMAL' dihedral angle restraints, as well as the 2nd, 3rd, *etc.* cosine dihedral angle restraints on the same atoms; the improper dihedral angle restraints are not considered here. For this command to work properly, the cosine dihedral angle restraints must be constructed before any other dihedral angle restraints. This functionality is needed because some of the CHARMM cosine terms are sometimes duplicated by other CHARMM cosine terms as well as by MODELLER homology-derived mainchain and sidechain dihedral and bi-dihedral angle restraints. In the standard `_model` script, the redundant CHARMM terms are always removed.

Example: See `model.read()` command.

4.7.9 `model.restraints.add()` — add restraint

Options:

<code>atom_ids = <str : 0></code>	<code>''</code>	<code>atom</code>	<code>ids:</code>
		<code>'atom:residue_id[:chain_id]'</code>	
<code>restraint_parameters = <float : 0></code>	<code>3 1 3 3 4 2 0 0.0 0.087</code>	<code>restraint parameters</code>	

Description: This command adds a specified restraint to the end of the restraints list and selects it. It can also add an excluded pair or a pseudo atom definition to the respective lists, depending on the dimension of `restraint_parameters` (Section 4.7.1). This command is useful for specifying *cis*-peptide bonds from your own scripts. The angles have to be in radians.

Example: `examples/commands/add_restraint.py`

```
# Example for: model.restraints.add(), model.restraints.unpick()

# This will enforce cis conformation for Pro-56.

# Make a model and stereochemical restraints:

def cispeptide(rsr, ids1, ids2):
    # Delete the old restraint on the same atoms:
    rsr.unpick(atom_ids=ids1)
    # Add the new restraint:
    rsr.add(atom_ids=ids1, restraint_parameters=(3, 1, 3, 3, 4, 2, 0,
                                                3.141593, 0.087))

    rsr.unpick(atom_ids=ids2)
    rsr.add(atom_ids=ids2, restraint_parameters=(3, 1, 3, 3, 4, 2, 0, 0.0,
                                                0.087))

log.level(output=1, notes=1, warnings=1, errors=1, memory=0)
env = environ()

env.libs.topology.read(file='$(LIB)/top_heav.lib')
env.libs.parameters.read(file='$(LIB)/par.lib')

aln = alignment(env)
code = '1fas'
mdl = model(env, file=code)
aln.append_model(mdl, align_codes=code, atom_files=code)
aln.append_model(mdl, align_codes='1fas-ini', atom_files='1fas.ini')

mdl.generate_topology(aln=aln, sequence='1fas-ini', add_segment=False)
mdl.transfer_xyz(aln=aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')
mdl.restraints.make(aln=aln, restraint_type='stereo', spline_on_site=False)

# Change the Pro-56 restraint from trans to cis:
cispeptide(mdl.restraints, ids1=( 'O:56', 'C:56', 'N:57', 'CA:57'),
          ids2=('CA:56', 'C:56', 'N:57', 'CA:57'))

mdl.restraints.write(file='1fas.rsr')
mdl.energy()
```

4.7.10 `model.restraints.unpick()` — unselect restraint

Options:

```
atom_ids = (str : 0)          ''          atom          ids:
                                'atom:residue_id[:chain_id]'
```

Description: This command scans the currently selected restraints to find all the restraints that operate on the specified atoms (Section 4.7.1) and then unselects them. The order of the atom names in `atom_ids` does not matter: All restraints that contain all and only the specified atoms are unselected. This means that it is not possible to distinguish between the dihedral angle and improper dihedral angle restraints on the same four atoms.

The command only unselects the restraints found. To completely remove all the unselected restraints from memory, use `model.restraints.condense()`. The `model.restraints.unpick()` command is useful in specifying *cis*-peptide bonds in your own scripts.

Example: See `model.restraints.add()` command.

4.7.11 `model.restraints.reindex()` — renumber MODEL2 restraints for MODEL

Options:

```
mdl = (model)                Input model, with old restraint number-
                                ing
```

Requirements: restraints

Description: This command renumbers atom indices in all restraints in memory. It is expected that the input restraints refer to `mdl`; the re-indexed restraints will correspond to MODEL. Only those restraints that have all atoms in MODEL will be selected. You can remove the others by `model.restraints.condense()`. This command is useful when the old restraints have to be used while changing from one topology model to another.

Example: `examples/commands/reindex_restraints.py`

```
# Example for: model.restraints.reindex()

# This will reindex restraints obtained previously for a simpler topology so
# that they will now apply to a more complicated topology.

env = environ()
tpl = env.libs.topology
par = env.libs.parameters

# Generate the model for the simpler topology (CA only in this case):
tpl.read(file='${LIB}/top_ca.lib')
par.read(file='${LIB}/par_ca.lib')

code = '1fas'
mdl = model(env, file=code)
aln = alignment(env)
aln.append_model(mdl, atom_files=code, align_codes=code)
aln.append_model(mdl, atom_files=code+'.ca', align_codes=code+'-ca')
mdl.generate_topology(aln, sequence=code+'-ca')
mdl.transfer_xyz(aln)
mdl.build(initialize_xyz=False, build_method='INTERNAL_COORDINATES')
mdl.write(file=code+'.ca')
```

```

# Generate the restraints for the simpler topology:
mdl.restraints.make(aln, restraint_type='stereo', spline_on_site=False)
mdl.restraints.write(file='1fas-ca.rsr')
mdl.energy()

# Generate the model for the more complicated topology:
tpl.read(file='${LIB}/top_heav.lib')
par.read(file='${LIB}/par.lib')

mdl.read(file=code)
aln = alignment(env)
aln.append_model(mdl, atom_files=code, align_codes=code)
aln.append_model(mdl, atom_files=code+'.ini', align_codes=code+'-ini')
mdl.generate_topology(aln, sequence=code+'-ini')
mdl.transfer_xyz(aln)
mdl.write(file='1fas.ini')

mdl2 = model(env, file='1fas.ca')
mdl.restraints.reindex(mdl2)
mdl.restraints.write(file='1fas.rsr')
mdl.energy()

```

4.7.12 `model.restraints.spline()` — approximate restraints by splines

Options:

<code>edat = <energy_data></code>		objective function parameters
<code>spline_dx = <float : 1></code>	0.5	interval size for splining restraints
<code>spline_min_points = <int : 1></code>	5	have at least as many intervals in a spline
<code>spline_range = <float : 1></code>	4.0	range of the splines
<code>spline_select = <int : 3></code>	4 1 9	specification of the restraints to be splined: form feature group
<code>output = <str : 1></code>	'LONG'	'SHORT' 'LONG' 'VERY_LONG' 'GRADIENT' 'SYMMETRY' 'ENERGY_PROFILE' 'VIOLATIONS_PROFILE'
<code>residue_span_range = <int : 2></code>	0 99999	range of residues spanning the allowed distances; for MAKE_RESTRAINTS, PICK_RESTRAINTS, non-bonded dynamic pairs

Description: This command calculates and selects new restraints that are a spline approximation of the selected restraints of the specified type. It unselects the approximated restraints.

The type of the approximated restraints is specified by `spline_select` and is defined by the mathematical form (Gaussian, *etc*), feature type (distance, *etc*), and physical restraint group (sidechain χ_1 , *etc*) (the first, third, and fourth integer numbers in the restraint specification).

The restraint is approximated in a certain range only, determined differently for different mathematical forms. For example, the poly-Gaussian range is from $m - \text{spline_range} \times \sigma_m$ to $M + \text{spline_range} \times \sigma_M$, where m and M are the minimal and maximal means of the basis pdfs, and σ_m and σ_M are their corresponding standard deviations.

The spline points are distributed evenly over this range with an interval of `spline_dx`. `spline_dx` should be equal to the scale of the peaks of the restraint that you want to approximate reliably. The value of the restraint

4.8 The `group_restraints` class: restraints on atom groups

The `group_restraints` class holds classifications of atoms into classes/groups, and restraints which act on certain atom groups. Such restraints are used for `energy_data.dynamic_modeller` restraints, such as the statistical potential used for loop modeling.

4.8.1 `group_restraints()` — create a new set of group restraints

Options:

<code>env = <environ></code>	MODELLER environment
<code>classes = <str : 1></code>	file of atom classes
<code>parameters = <str : 1></code>	file of group restraint parameters

Description: This creates a new set of group restraints. The set is initialized by reading in `classes`, a file containing a classification of residue:atom pairs into groups. If the `parameters` argument is also given, this is used to read in a file of restraint parameters, using `group_restraints.append()`.

Example: `examples/commands/read_atom_classes.py`

```
# Example for: group_restraints()

# This will read an atom classification for non-bonded statistical potentials
# of mean force.

env = environ()
gprsr = group_restraints(env, classes='$(LIB)/atmcls-melo.lib',
                        parameters='$(LIB)/melo1-dist.lib')
```

4.8.2 `group_restraints.append()` — read group restraint parameters

Options:

<code>file = <str : 1></code>	'default'	partial or complete filename
-------------------------------------	-----------	------------------------------

Description: This reads a set of parameters from `file`, which should act on the atom classes previously defined. Any parameters read are added to any already in this object (to clear them all, simply create a new object).

4.9 The alignment class: comparison of sequences and structures

This section describes the format of the alignment file and commands for reading, writing, making, analyzing and using the alignments of sequences and structures (pairwise and multiple). For the underlying dynamic programming methods see Section 6.1.

4.9.1 Alignment file format

The preferred format for comparative modeling is related to the PIR database format:

C; A sample alignment in the PIR format; used in tutorial

```
>P1;5fd1
structureX:5fd1:1      : :106  : :ferredoxin:Azotobacter vinelandii: 1.90: 0.19
AFVVTDNCIKCKYTDCVEVCPVDCFYEGPNFLVIHPDECIDCALCEPECQAIFSEDEVPEDMQEFIQLNAELA
EVWPNITEKKDPLPDAEDWDGKGLQHLER*

>P1;1fdx
sequence:1fdx:1      : :54   : :ferredoxin:Peptococcus aerogenes: 2.00:-1.00
AYVINDSC--IACGACKPECPVNIQGS--IYAIDADSCIDCGSCASVCPVGAPNPED-----
-----*
```

The first line of each sequence entry specifies the protein code after the >P1; line identifier. The line identifier must occur at the beginning of the line. For example, 1fdx is the protein code of the first entry in the alignment above. The protein code corresponds to the `alignment[x].code` variable.

The second line of each entry contains information necessary to extract atomic coordinates of the segment from the original PDB coordinate set. The fields in this line are separated by colon characters, ':'. The fields are as follows:

- Field 1: A specification of whether or not 3D structure is available and of the type of the method used to obtain the structure (`structureX`, X-ray; `structureN`, NMR; `structureM`, model; `sequence`, sequence). Only `structure` is also a valid value.
- Field 2: The PDB code. While the protein code in the first line of an entry, which is used to identify the entry, must be unique for all proteins in the file, the PDB code in this field, which is used to get structural data, does not have to be unique. It is a good idea to use the PDB code with an optional chain identifier as the protein code. The PDB code corresponds to the `alignment[x].atom_file` variable and can also contain the full atom filename, directory included.
- Fields 3–6: The residue identifiers (see below) for the first (fields 3–4) and last residue (fields 5–6) of the sequence in the subsequent lines. There is no need to edit the coordinate file if a contiguous sequence of residues is required — simply specify the beginning and ending residues of the required contiguous region of the chain. If the beginning residue is not found, no segment is read in. If the ending residue identifier is not found in the coordinate file, the last residue in the coordinate file is used. By default, the whole file is read in.
 The unspecified beginning and ending residue numbers and chain id's for a `structure` entry in an alignment file are taken automatically from the corresponding atom file, if possible. The first matching sequence in the atom file that also satisfies the explicitly specified residue numbers and chain id's is used. A residue number is not specified when a blank character or a dot, '.', is given. A chain id is not specified when a dot, '.', is given. This slight difference between residue and chain id's is necessary because a blank character is a valid chain id.
- Field 7: Protein name. Optional.
- Field 8: Source of the protein. Optional.
- Field 9: Resolution of the crystallographic analysis. Optional.

Field 10: R-factor of the crystallographic analysis. Optional.

A residue identifier consists of a residue number and an optional chain identifier. They must be separated by a colon, ':'. For example, '10I:A' is residue number '10I' in chain 'A', and '6' or '6:' is residue number '6' in a chain without a name. Free format can be used, that is the blank characters are ignored. The residue number is a string of up to 5 characters long, as found in the PDB atom file and consists of the PDB residue number proper (22X,A4 in the PDB ATOM record) and PDB residue insertion code (26X, A1). The chain identifier is a single character, as found in the PDB atom file (21X,A1).

The residue number for the first position (resID1) in the `model_segment` range '`resID1:chainID1 resID2:chainID2`' can be either a real residue number or 'FIRST' (which indicates the first residue in a matching chain). The residue number for the second position (resID2) in the `model_segment` range can be either: (1) a real residue number; (2) 'LAST' (which indicates the last residue in a matching chain); or 'END' (which indicates the last residue in the PDB file). The chain id for either position in the `model_segment` range (chainID1 or chainID2) can be either: (1) a real chain id (including a blank/space/null/empty); or '@', which matches any chain id.

Examples, assuming a two chain PDB file (chains A and B):

- '15:A 75:A' reads residues 15 to 75 in chain A.
- 'FIRST:@ 75:@' reads the first 75 residues in chain A (the first chain).
- 'FIRST:@ LAST:@' reads all residues in chain A, assuming 'FIRST' is not a real number of the non-first residue in chain A.
- '10:@ LAST:'' reads all residues from 10 in chain A to the end of the file (chain id for the last residue is irrelevant), again assuming 'LAST' is not a real residue number of a non-last residue.
- 'FIRST:@ END:'' reads the whole file no matter what, the chainID is ignored completely.

For the `selection_segment` the string containing '@' will match any residue number and chainID. For example, '@:A' is the first residue in chain 'A' and '@:@' is the first residue in the coordinate file. The last chain can not be specified in a general way, except if it is the last residue in the file.

When an alignment file is used in conjunction with structural information, the first two fields must be filled in, the rest of them can be empty or even missing entirely. If the alignment is not used in conjunction with structural data, all but the first field can be empty. This means that in comparative modeling, the template structures must have at least the first two fields specified while the target sequence must only have the first field filled in. Thus, a simple second line of an entry in an alignment file in the 'PIR' format is

```
structure:pdb_file:.....
```

This entry will result in reading from PDB file `pdb_file` the structure segment corresponding to the sequence in the subsequent lines of the alignment entry.

The fields that do not exist are assigned blank values. Thus,

```
structure:pdb_file
```

is equivalent to

```
structure:pdb_file: : : : : :
```

which will achieve what was probably intended (read in the structure segment from file `pdb_file` that corresponds to the sequence in the subsequent lines of the alignment entry) only if the chain id is a blank character.

Each sequence must be terminated by the terminating character, '*'.

When the first character of the sequence line is the terminating character, '*', the sequence is obtained from the specified PDB coordinate file (Section 4.1.3).

Chain breaks are indicated by '/'. There should not be more than one chain break character to indicate a single chain break (use gap characters instead, '-'). All residue types specified in \$RESTYP_LIB, but not patching residue types, are allowed; there are on the order of 100 residue types specified in the \$RESTYP_LIB library. To add your own residue types to this library, see Section 1.8, Question 10.

The alignment file can contain any number of blank lines between the protein entries. Comment lines can occur outside protein entries and must begin with the identifiers 'C;' or 'R;' as the first two characters in the line.

An alignment file is also used to input non-aligned sequences.

Example: examples/commands/aln_append_model.py

```
# This demonstrates one way to generate an initial alignment between two
# PDB sequences. It can later be edited by hand.

# Set Modeller environment (including search patch for model.read())
env = environ()
env.io.atom_files_directory = "./:../atom_files/"

# Create a new empty alignment and model:
aln = alignment(env)
mdl = model(env)

# Read the whole 1fdx atom file
code='1fdx'
mdl.read(file=code, model_segment=('FIRST:@', 'END:'))

# Add the model sequence to the alignment
aln.append_model(mdl, align_codes=code, atom_files=code)

# Read 5fd1 atom file from 1-63, and add to alignment
code='5fd1'
mdl.read(file=code, model_segment=('1:', '63:'))
aln.append_model(mdl, align_codes=code, atom_files=code)

# Align them by sequence
aln.malign(gap_penalties_1d=(-500, -300))
aln.write(file='fer1-seq.ali')

# Align them by structure
aln.malign3d(gap_penalties_3d=(0.0, 2.0))

# check the alignment for its suitability for modeling
aln.check()

aln.write(file='fer1.ali')
```

4.9.2 alignment() — create a new alignment

Options:

env = <environ>

MODELLER environment

Description: This creates a new `alignment` object; by default, this contains no sequences. If any keyword arguments are given, they are passed to the `alignment.append()` or `alignment.append_profile()` functions to create the initial alignment.

4.9.3 `alignment.append()` — read sequences and/or their alignment

Options:

<code>io = <io_data></code>		Options for reading atom files
<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
<code>align_codes = <str : 0></code>	<code>'all'</code>	codes of proteins in the alignment
<code>atom_files = <str : 0></code>	<code>''</code>	complete or partial atom filenames
<code>alignment_format = <str : 1></code>	<code>'PIR'</code>	format of the alignment file: <code>'PIR'</code> <code>'PAP'</code> <code>'QUANTA'</code> <code>'INSIGHT'</code> <code>'FASTA'</code>
<code>remove_gaps = <bool : 1></code>	<code>True</code>	whether to remove all-gap positions in input alignment
<code>close_file = <bool : 1></code>	<code>True</code>	whether or not to close the alignment file at the end of <code>READ_ALIGNMENT</code>
<code>rewind_file = <bool : 1></code>	<code>False</code>	whether or not to rewind the alignment file at the start of <code>READ_ALIGNMENT</code>

Output: `end_of_file`

Description: This command reads the sequence(s) and/or their alignment from a text file. Only sequences with the specified codes are read in; `align_codes = 'all'` can be used to read all sequences. The sequences are added to any currently in the alignment.

There are several alignment formats:

1. The `'PIR'` format resembles that of the PIR sequence database. It is described in Section 4.9.1 and is used for comparative modeling because it allows for additional data about the proteins that are useful for automated access to the atomic coordinates.
2. The `'FASTA'` format resembles the `'PIR'` format but has a missing second 'comment' line and a missing star at the end of each sequence.
3. The `'PAP'` format is nicer to look at but contains less information and is not used by other programs. When used in conjunction with PDB files, the PDB files must contain exactly the residues in the sequences in the `'PAP'` file; *i.e.*, it is not possible to use only a segment of a PDB file. In addition, the `'PAP'` protein codes must be expandable into proper PDB atom filenames, as described in Section 4.1.3. The protein sequence can now start in any column (this was limited to column 11 before release 5).
4. The `'QUANTA'` format can be used to communicate with the QUANTA program. You are not supposed to mix `'QUANTA'` format with any other format because the `'QUANTA'` format contains residue numbers which do not occur in the other formats and are difficult to guess correctly. MODELLER can write out alignments in the `'QUANTA'` format but cannot read them in.
5. The `'INSIGHT'` format is very similar to the `'PAP'` format and can sometimes be used to communicate with the INSIGHTII program. When used in conjunction with PDB files, the same rules as for the `'PAP'` format apply.

If `remove_gaps = True`, positions with gaps in all selected sequences are removed from the alignment.

The `io` argument is required since PIR files can contain empty sequences; in this case, the sequence is read from the corresponding PDB file.

Ordinarily, the alignment file is closed at the end of this command. However, when reading `'PIR'` or `'FASTA'` format files, if `close_file` is `False`, then the file is left open. Subsequent calls to `alignment.append()` will then resume at this point in the file, provided they set `rewind_file` to `False`. The `end_of_file` variable is set to 1 if MODELLER reached the end of the `'PIR'` or `'FASTA'` file during the read, or 0 otherwise.

Example: examples/commands/read_alignment.py

```
# Example for: alignment.append(), alignment.write(),
#             alignment.check()

# Read an alignment, write it out in the 'PAP' format, and
# check the alignment of the N-1 structures as well as the
# alignment of the N-th sequence with each of the N-1 structures.

log.level(output=1, notes=1, warnings=1, errors=1, memory=0)
env = environ()

aln = alignment(env)
aln.append(file='toxin.ali', align_codes='all')
aln.write(file='toxin.pap', alignment_format='PAP')
aln.write(file='toxin.fasta', alignment_format='FASTA')
aln.check()
```

4.9.4 alignment.clear() — delete all sequences from the alignment

Description: This deletes all of the sequences from the alignment. It is not exactly the same as deleting the alignment object and creating a new one, since any structural data already read in remains in the alignment object. This is useful if the sequences are reread and the structural information needs to be reused.

4.9.5 alignment.check() — check alignment for modeling**Options:**

io = {io_data}

Options for reading atom files

Description: This command evaluates an alignment to be used for comparative modeling. It uses two criteria:

First, it checks the alignment of the template structures (all but the last entry in the alignment): For each pairwise superposition of the templates, it reports those equivalent pairs of C_{α} atoms that are more than 6Å away from each other. Such pairs are almost certainly misaligned. The pairwise superpositions are done using the C_{α} atoms and the given alignment.

Second, the command checks the alignment of the target sequence (the last entry in the alignment) with each of the templates: For all consecutive pairs of C_{α} atoms in the target, it calculates the distance between the two equivalent C_{α} atoms in each of the templates. If the distance is longer than 8Å, it is reported. In such a case, the alignment between the template and the target is almost certainly incorrect.

Example: examples/commands/read_alignment.py

```
# Example for: alignment.append(), alignment.write(),
#             alignment.check()

# Read an alignment, write it out in the 'PAP' format, and
# check the alignment of the N-1 structures as well as the
# alignment of the N-th sequence with each of the N-1 structures.

log.level(output=1, notes=1, warnings=1, errors=1, memory=0)
env = environ()
```

```
aln = alignment(env)
aln.append(file='toxin.ali', align_codes='all')
aln.write(file='toxin.pap', alignment_format='PAP')
aln.write(file='toxin.fasta', alignment_format='FASTA')
aln.check()
```

4.9.6 alignment.compare_with() — compare two alignments

Options:

aln = *<alignment>* Alignment to compare with

Description: This command compares two pairwise alignments read by the **alignment.append()** commands. The alignment of the first sequence with the second sequence in ALIGNMENT is evaluated with respect to 'aln'. The numbers are not symmetric; *i.e.*, they will change if the sequences or alignments are swapped. The output in the log file is self-explanatory.

Example: examples/commands/compare_alignments.py

```
# Example for: alignment.compare_with(), alignment.append_model()

# Compare two alignments of two proteins each. In this case, the first
# alignment is a sequence-sequence alignment and the second alignment
# is a structure-structure alignment.

log.level(1, 1, 1, 1, 0)
env = environ()

# Generate and save sequence-sequence alignment:
aln = alignment(env)
for code in ('1fas', '2ctx'):
    mdl = model(env, file=code)
    aln.append_model(mdl=mdl, align_codes=code, atom_files=code)
aln.align(gap_penalties_1d=(-600, -400))
aln.write(file='toxin-seq.ali')

# Generate and save structure-structure alignment:
aln.align3d(gap_penalties_3d=(0, 2.0))
aln.write(file='toxin-str.ali')

# Compare the two pairwise alignments:
aln = alignment(env, file='toxin-seq.ali', align_codes='all')
aln2 = alignment(env, file='toxin-str.ali', align_codes='all')
aln.compare_with(aln2)
```

4.9.7 alignment.append_model() — copy model sequence and coordinates to alignment

Options:

mdl = *<model>* model to read sequence from

```
align_codes = <str : 0>          'all'          codes of proteins in the alignment
atom_files = <str : 0>          ''           complete or partial atom filenames
```

Description: This command adds the sequence and coordinates of the given model, `mdl`, to the end of the current alignment.

You should additionally set `align_codes` and `atom_files` to the PDB ID and file name, respectively. This information is added to the alignment with the new sequence. (Alternatively, you can set this information later by setting `alignment[x].code` and `alignment[x].atom_file`, where 'x' is the index of the sequence in the alignment, starting at zero.)

Example: See `alignment.compare_with()` command.

4.9.8 `alignment.append_profile()` — add profile sequences to the alignment

Options:

```
prf = <profile>                input profile
```

Description: This adds all the sequences from the given profile, `prf`, to the alignment. It is similar in operation to `profile.to_alignment()`.

4.9.9 `alignment.write()` — write sequences and/or their alignment

Options:

```
file = <str : 1>                'default'    partial or complete filename
output_directory = <str : 1>    ''           output directory
alignment_format = <str : 1>    'PIR'        format of the alignment file: 'PIR'
                                     | 'PAP' | 'QUANTA' | 'INSIGHT' |
                                     'FASTA'
align_block = <int : 1>         0            the last sequence in the first block of se-
                                     quences
align_alignment = <bool : 1>    False        writing out an alignment of alignments
                                     (for *)
alignment_features = <str : 1>  'INDICES CONSERVATION'
                                     what alignment features to write out:
                                     'ACCURACY' | 'HELIX' | 'BETA' |
                                     'ACCESSIBILITY' | 'STRAIGHTNESS' |
                                     'CONSERVATION' | 'INDICES' | 'ALL' |
                                     'GAPS'
```

Description: This command writes the whole alignment to a text file.

The 'PAP' format, which corresponds to a relatively nice looking alignment, has several additional formatting options that can be selected by the `alignment_features` variable. This scalar variable can contain any combination of the following keywords:

- 'INDICES', the alignment position indices;
- 'CONSERVATION', a star for each absolutely conserved position;
- 'ACCURACY', the alignment accuracy indices, scaled between 0–9, as calculated by `alignment.consensus()`;
- 'HELIX', average content of helical residues for structures 1 – `align_block` at each position, 0 for 0% and 9 for 100%, as calculated by `alignment.align2d()`.

- 'BETA', average content of β -strand residues for structures 1 – align_block at each position, 0 for 0% and 9 for 100%, as calculated by `alignment.align2d()`.
- 'ACCESSIBILITY', average relative sidechain buriedness for structures 1 – align_block, 0 for 0% (100% accessibility) and 9 for 100% (0% accessibility), as calculated by `alignment.align2d()`;
- 'STRAIGHTNESS', average mainchain straightness structures 1 – align_block at each position 0 for 0% and 9 for 100%, as calculated by `alignment.align2d()`.

Options 'HELIX', 'BETA', 'ACCESSIBILITY', and 'STRAIGHTNESS' are valid only after executing command `alignment.align2d()`, where the corresponding quantities are defined. They refer to the 3D profile defined for the first align_block structures (run `alignment.align2d()` with `fit = False` to prepare these structural data without changing the input alignment). Similarly, the 'ACCURACY' option is valid only after the `alignment.consensus()` command.

`align_alignment` and `align_block` are used to ensure correct indication of identical alignment positions, depending on whether sequences or two blocks of sequences were aligned: For sequences (`align_alignment = False` and `align_block` is ignored), a '*' indicating a conserved position is printed where all sequences have the same residue type. For blocks (`align_alignment = True` and `align_block` indicates the last sequence of the first block), a '*' is printed only where the two blocks have the same order of residue types (there has to be the same number of sequences in both blocks). The blocks option is useful when comparing two alignments, possibly aligned by the `alignment.align()` command.

Example: See `alignment.append()` command.

4.9.10 alignment.edit() — edit overhangs in alignment

Options:

<code>io = <io_data></code>		Options for reading atom files
<code>edit_align_codes = <str : 0></code>	'last'	codes of proteins in the alignment to be edited
<code>base_align_codes = <str : 0></code>	'rest'	codes of proteins in the alignment to be used as the base
<code>min_base_entries = <int : 1></code>	1	minimal number of templates in EDIT_ALIGNMENT
<code>overhang = <int : 1></code>	0	un-penalized overhangs in protein comparisons

Description: This command edits the overhangs in the alignment.

`edit_align_codes` specifies the alignment codes for the alignment entries whose overhangs are to be cut; in addition, `all` or `last` can be used.

`base_align_codes` specifies the alignment codes for the alignment entries that are used to determine the extent of the overhangs to be cut from the edited entries; in addition, `all` or `rest` (relative to `edit_align_codes`) can be used.

The same entries can be cut and used for determining the base.

The base of the alignment is determined by the first and last alignment positions that have at least `min_base_entries` entries that started by that position, beginning from the first and last alignment positions, respectively.

The cuts are shortened by `overhang` residues respectively, so that reasonably short termini can be easily modeled ab initio if desired.

The `io` argument is used because the beginning and ending residue numbers for the 'structure' entries in the alignment are renumbered automatically by reading the appropriate atom files.

Example: `examples/commands/edit_alignment.py`

```

# Example for: alignment.edit()

# Read an alignment, write it out in the 'PAP' format, with overhangs cut.

log.level(1, 1, 1, 1, 0)
env = environ()
env.io.atom_files_directory = './:../atom_files'

aln = alignment(env, file='overhang.ali', align_codes='all',
                alignment_format='PIR')

# Cut overhangs in the 1is4 sequence that are longer than 3 residues
# relative to the longest remaining entry in the alignment:
aln.edit(edit_align_codes='1is4', base_align_codes='rest',
         min_base_entries=1, overhang=3)
aln.write(file='overhang-1.pir', alignment_format='PIR')
aln.write(file='overhang-1.pap', alignment_format='PAP')

```

4.9.11 alignment.describe() — describe proteins

Options:

io = `<io_data>`

Options for reading atom files

Description: This command outputs basic data about the proteins in the current alignment (*e.g.* as read in by `alignment.append()`). This command is useful for preparation before comparative modeling because it summarizes disulfides, *cis*-prolines, charges, chain breaks, *etc.* Results which depend only on the amino acid sequences are still written out even if some atom files do not exist.

Example: `examples/commands/describe.py`

```

# Example for: alignment.describe()

# Describe the sequences and structures in the alignment.

env = environ()
aln = alignment(env, file='toxin.ali', align_codes=('2ctx', '2abx'))
aln.describe()

```

4.9.12 alignment.id_table() — calculate percentage sequence identities

Options:

matrix_file = `<str : 1>`

'family.mat'

the filename of the pairwise distance matrix

output_directory = `<str : 1>`

''

output directory

Description: This command calculates percentage residue identities for all pairs of sequences in the current alignment. The percentage residue identity is defined as the number of identical residues divided by the length of the shorter sequence.

In addition to the output in the log file, this routine creates file `output_directory/matrix_file` with pairwise sequence distances that can be used directly as the input to the tree making programs of the PHYLIP package, such as KITSCH [Felsenstein, 1985], and also for the `environ.dendrogram()` and `environ.principal_components()` commands. A more general version of this command, which allows a user specified measure for residue–residue differences is `alignment.compare_sequences()`.

Example: `examples/commands/id_table.py`

```
# Example for: alignment.id_table(), alignment.compare_sequences(),
#             misc.principal_components(), misc.dendrogram()

# Pairwise sequence identity between sequences in the alignment.

env = environ()
# Read all entries in this alignment:
aln = alignment(env, file='toxin.ali')

# Calculate pairwise sequence identities:
aln.id_table(matrix_file='toxin_id.mat')

# Calculate pairwise sequence similarities:
mdl = model(env, file='2ctx', model_segment=('1:', '71:'))
aln.compare_sequences(mdl, rr_file='$(LIB)/as1.sim.mat', max_gaps_match=1,
                    matrix_file='toxin.mat', variability_file='toxin.var')
mdl.write(file='2ctx.var')

# Do principal components clustering using sequence similarities:
env.principal_components(matrix_file='toxin.mat', file='toxin.princ')

# Dendrogram in the log file:
env.dendrogram(matrix_file='toxin.mat')
```

4.9.13 `alignment.compare_sequences()` — compare sequences in alignment

Options:

<code>mdl = <model></code>		Model
<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>matrix_file = <str : 1></code>	<code>'family.mat'</code>	the filename of the pairwise distance matrix
<code>variability_file = <str : 1></code>	<code>'undefined'</code>	output filename
<code>output_directory = <str : 1></code>	<code>''</code>	output directory
<code>max_gaps_match = <int : 1></code>	<code>1</code>	

Description: The pairwise similarity of sequences in the current alignment is evaluated using a user specified residue–residue scores file.

The residue–residue scores, including gap–residue, and gap–gap scores, are read from file `rr_file`. The sequence pair score is equal to the average pairwise residue–residue score for all alignment positions that have at most `max_gaps_match` gaps (1 by default). If the gap–residue and gap–gap scores are not defined in `matrix_file`, they are set to the worst and best residue–residue score, respectively. If `matrix_file` is a similarity matrix, it is converted into a distance matrix ($x' = -x + x_{\max}$).

The comparison matrix is written in the PHYLIP format to file `matrix_file`.

The family variability as a function of alignment position is calculated as the RMS deviation of all residue – residue scores at a given position, but only for those pairs of residues that have at most `max_gaps_match` gaps (0, 1, or 2). The variability is written to file `variability_file`, as is the number of pairwise comparisons contributing to each positional variability. The variability is also written into the fourth column of the current model, which must correspond to the first sequence in the alignment.

Example: See `alignment.id_table()` command.

4.9.14 `alignment.align()` — align two (blocks of) sequences

Options:

<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>gap_penalties_1d = <float : 2></code>	900 50	gap creation and extension penalties for sequence/sequence alignment
<code>align_block = <int : 1></code>	0	the last sequence in the first block of sequences
<code>off_diagonal = <int : 1></code>	100	to speed up the alignment
<code>matrix_offset = <float : 1></code>	0.00	substitution matrix offset for local alignment
<code>overhang = <int : 1></code>	0	un-penalized overhangs in protein comparisons
<code>max_gap_length = <int : 1></code>	999999	maximal length of gap in protein comparisons
<code>local_alignment = <bool : 1></code>	False	whether to do local as opposed to global alignment
<code>align_what = <str : 1></code>	'BLOCK'	what to align in ALIGN; 'BLOCK' 'ALIGNMENT' 'LAST' 'PROFILE'
<code>read_weights = <bool : 1></code>	False	whether to read the whole NxM weight matrix for ALIGN*
<code>write_weights = <bool : 1></code>	False	whether to write the whole NxM weight matrix for ALIGN*
<code>input_weights_file = <str : 1></code>	''	External weight matrix input to MOD-ELLER (SALIGN/ALIGN)
<code>output_weights_file = <str : 1></code>	''	File into which the weight file is written (iff WRITE_WEIGHTS = 'on')
<code>weigh_sequences = <bool : 1></code>	False	whether or not to weigh sequences in a profile
<code>smooth_prof_weight = <float : 1></code>	10	for smoothing the profile aa frequency with a prior

Description: This command aligns two blocks of sequences.

The two blocks of sequences to be aligned are sequences 1 to `align_block` and `align_block+1` to the last sequence. The sequences within the two blocks should already be aligned; their alignment does not change.

The command can do either the global (similar to [Needleman & Wunsch, 1970]; `local_alignment = False`) or local dynamic programming alignment (similar to [Smith & Waterman, 1981]; `local_alignment = True`).

For the global alignment, set overhang length `overhang` to more than 0 so that the corresponding number of residues at either of the four termini won't be penalized by any gap penalties (this makes it a pseudo local alignment).

To speed up the calculation, set `off_diagonal` to a number smaller than the shortest sequence length. The alignments matching residues i and j with $|i - j| > \text{off_diagonal}$ are not considered at all in the search for the best alignment.

The gap initiation and extension penalties are specified by `gap_penalties_1d`. The default values of -900 -50 for the `'as1.sim.mat'` similarity matrix were found to be optimal for pairwise alignments of sequences that share from 30% to 45% sequence identity (RS and AŠ, in preparation).

The residue type – residue type scores are read from file `rr_file`. The routine automatically determines whether it has to maximize similarity or minimize distance.

`matrix_offset` applies to local alignment only and influences its length. `matrix_offset` should be somewhere between the lowest and highest residue–residue scores. A smaller value of this parameter will make the local alignments shorter when distance is minimized, and longer when similarity is maximized. This works as follows: The recursively constructed dynamic programming comparison matrix is reset to 0 at position i, j when the current alignment score becomes larger (distance) or smaller (similarity) than `matrix_offset`. Note that this is equivalent to the usual shifting of the residue–residue scoring matrix in the sense that there are two combinations of `gap_penalties_1d` and `matrix_offset` values that will give exactly the same alignments irrespective of whether the matrix is actually offset (with 0 used to restart local alignments in dynamic programming) or the matrix is not offset but `matrix_offset` is used as the cutoff for restarting local alignments in dynamic programming. For the same reason, the matrix offset does not have any effect on the global alignments if the gap extension penalty is also shifted for half of the matrix offset.

The position–position score is an average residue–residue score for all possible pairwise comparisons between the two blocks ($n \times m$ comparisons are done, where n and m are the number of sequences in the two blocks, respectively). The first exception to this is when `align_what` is set to 'ALIGNMENT', in which case the two alignments defined by `align_block` are aligned; *i.e.*, the score is obtained by comparing only equivalent positions between the two alignment blocks (only n comparisons are done, where n is the number of sequences in each of the two blocks). This option is useful in combination with `alignment.compare_with()` and `alignment.write()` for evaluation of various alignment parameters and methods. The second exception is when `align_what` is set to 'LAST', in which case only the last sequences in the two blocks are used to get the scores. In 'BLOCK', 'ALIGNMENT', and 'LAST' comparisons, penalty for a comparison of a gap with a residue during the calculation of the scoring matrix is obtained from the score file (gap–gap match should have a score of 0.0).

Only the 20 standard residue types, plus Asx (changes to Asn) and Glx (changes to Gln) are recognized. Every other unrecognized residue, except for a gap and a chain break, changes to Gly for comparison purposes.

For the time being, this and the other alignment commands (`alignment.malign()`, `alignment.align2d()`, `alignment.align3d()`, and `alignment.malign3d()`) remove chain break information from the `CALN` array, which means that chain breaks are not retained when the alignment is written to a file after executing these commands.

Example: `examples/commands/align.py`

```
# Example for: alignment.align()

# This will read two sequences, align them, and write the alignment
# to a file:

log.verbose()
env = environ()

aln = alignment(env)
aln.append(file='toxin.ali', align_codes=('1fas', '2ctx'))

# The as1.sim.mat similarity matrix is used by default:
aln.align(gap_penalties_1d=(-600, -400))
aln.write(file='toxin-seq.ali')
```

4.9.15 `alignment.align2d()` — align sequences with structures

Options:

`io = <io_data>`

Options for reading atom files

<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>gap_penalties_1d = <float : 2></code>	900 50	gap creation and extension penalties for sequence/sequence alignment
<code>gap_penalties_2d = <float : 9></code>	0.35 1.2 0.9 1.2 0.6 8.6 1.2 0 0	gap penalties for sequence/structure alignment: helix, beta, accessibility, straightness, and CA-CA distance factor, dst min, dst power, t, structure_profile ; best U,V=-450,0
<code>align_block = <int : 1></code>	0	the last sequence in the first block of sequences
<code>max_gap_length = <int : 1></code>	999999	maximal length of gap in protein comparisons
<code>off_diagonal = <int : 1></code>	100	to speed up the alignment
<code>matrix_offset = <float : 1></code>	0.00	substitution matrix offset for local alignment
<code>overhang = <int : 1></code>	0	un-penalized overhangs in protein comparisons
<code>local_alignment = <bool : 1></code>	False	whether to do local as opposed to global alignment
<code>align_what = <str : 1></code>	'BLOCK'	what to align in ALIGN; 'BLOCK' 'ALIGNMENT' 'LAST' 'PROFILE'
<code>subopt_offset = <float : 1></code>	0.0	offset for residue-residue score in getting suboptimals in ALIGN/ALIGN2D
<code>fit = <bool : 1></code>	True	whether to align
<code>read_weights = <bool : 1></code>	False	whether to read the whole NxM weight matrix for ALIGN*
<code>write_weights = <bool : 1></code>	False	whether to write the whole NxM weight matrix for ALIGN*
<code>input_weights_file = <str : 1></code>	''	External weight matrix input to MODELLER (SALIGN/ALIGN)
<code>output_weights_file = <str : 1></code>	''	File into which the weight file is written (iff WRITE_WEIGHTS = 'on')
<code>weigh_sequences = <bool : 1></code>	False	whether or not to weigh sequences in a profile
<code>smooth_prof_weight = <float : 1></code>	10	for smoothing the profile aa frequency with a prior
<code>read_profile = <bool : 1></code>	False	whether to read str profile for ALIGN2D
<code>input_profile_file = <str : 1></code>	''	multiple sequence alignment read into MODELLER for profile-profile alignments
<code>write_profile = <bool : 1></code>	False	whether to write str profile for ALIGN2D
<code>output_profile_file = <str : 1></code>	''	

Description: This command aligns a block of sequences (second block) with a block of structures (first block). It is the same as the `alignment.align()` command except that a variable gap opening penalty is used. This gap penalty depends on the 3D structure of all sequences in block 1. The variable gap penalty can favor gaps in exposed regions, avoid gaps within secondary structure elements, favor gaps in curved parts of the main-chain, and minimize the distance between the two C_α positions spanning a gap. The `alignment.align2d()` command is preferred for aligning a sequence with structure(s) in comparative modeling because it tends to place gaps in a better structural context. See Section 6.1.2 for the dynamic programming algorithm that implements the variable gap penalty. `gap_penalties_2d` specifies parameters ω_H , ω_S , ω_B , ω_C , ω_D , d_o , γ , t and ω_{SC} . (Section 6.1.2). The default gap penalties `gap_penalties_1d` (-450, 0) and `gap_penalties_2d` (0.35, 1.2, 0.9, 1.2, 0.6, 8.6, 1.2, 0.0, 0.0) as well as the `rr_file` substitution matrix ('as1.sim.mat') were found to

be optimal in pairwise alignments of structures and sequences sharing from 30% to 45% sequence identity (MSM, MAM-R, RS and AS, in preparation).

— move to back

The linear gap penalty function for inserting a gap in block 1 of structures is: $g = f_1(H, S, B, C, SC)u + lv$ where u and v are the usual gap opening and extension penalties, l is gap length, and f_1 is a function that is at least 1, but can be larger to make gap opening more difficult in the following circumstances: between two consecutive (*i.e.*, $i, i + 1$) helical positions, two consecutive β -strand positions, two consecutive buried positions, or two consecutive positions where the mainchain is locally straight. This function is $f_1 = 1 + [\omega_H H_i H_{i+1} + \omega_S S_i S_{i+1} + \omega_B B_i B_{i+1} + \omega_C C_i C_{i+1} + \omega_S C S C_i S C_{i+1}]$, H_i is the fraction of helical residues at position i in block 1, S_i is the fraction of β -strand residues at position i in block 1, B_i is the average relative sidechain buriedness of residues at position i in block 1, C_i is the average straightness of residues at position i in block 1, and SC_i is the structural conservedness at position i in block 1. See Section 4.6.24 for the definition of these features. The original straightness is modified here by assigning maximal straightness of 1 to all residues in a helix or a β -strand. The structural conservedness of the residues in block 1 are imported from an external source "input_profile_file". The structural conservedness at a particular position gives the likelihood of the occurrence of a gap when structurally similar regions from all known protein structures are aligned structurally.

The linear gap penalty function for opening a gap in block 2 of sequences is: $g = f_2(H, S, B, C, D, SC)u + lv$ where f_2 is a function that is at least 1, but can be larger to make the gap opening in block 2 more difficult in the following circumstances: when the first gap position is aligned with a helical residue, a β -strand residue, a buried residue, extended mainchain, or when the whole gap in block 2 is spanned by two residues in block 1 that are far apart in space. This function is $f_2 = 1 + [\omega_H H_i + \omega_S S_i + \omega_B B_i + \omega_C C_i + \omega_D \sqrt{d - d_o} + \omega_S C S C_i]$. d is the distance between the two C_α atoms spanning the gap, averaged over all structures in block 1 and d_o is the distance that is small enough to correspond to no increase in the opening gap penalty (*e.g.*, 8.6Å).

When `fit` is `False`, no alignment is done and the routine returns only the average structural information, which can be written out by the `alignment.write()` command.

Example: examples/commands/align2d.py

```
# Demonstrating ALIGN2D, aligning with variable gap penalty

log.verbose()
env = environ()

env.libs.topology.read(file='$(LIB)/top_heav.lib')

# Read aligned structure(s):
aln = alignment(env)
aln.append(file='toxin.ali', align_codes='2ctx')
aln_block = len(aln)

# Read aligned sequence(s):
aln.append(file='toxin.ali', align_codes='1nbt')

# Structure sensitive variable gap penalty sequence-sequence alignment:
aln.align2d(overhang=0, gap_penalties_1d=(-450, 0),
            gap_penalties_2d=(0.35, 1.2, 0.9, 1.2, 0.6, 8.6, 1.2, 0., 0.),
            align_block=aln_block)

aln.write(file='align2d.ali', alignment_format='PIR')
aln.write(file='align2d.pap', alignment_format='PAP',
          alignment_features='INDICES HELIX BETA STRAIGHTNESS ' + \
                              'ACCESSIBILITY CONSERVATION')

aln.check()
```

```

# Color the first template structure according to gaps in alignment:
aln = alignment(env)
aln.append(file='align2d.ali', align_codes=('2ctx', '1nbt'),
           alignment_format='PIR', remove_gaps=True)
mdl = model(env)
mdl.read(aln=aln, model_segment=('2ctx', '2ctx'))
mdl.color(aln=aln)
mdl.write(file='2ctx.aln.pdb')

# Color the first template structure according to secondary structure:
mdl.write_data(file='2ctx', output='SSM')
mdl.write(file='2ctx.ssm.pdb')

# Superpose the target structure onto the first template:
mdl2 = model(env)
mdl2.read(aln=aln, model_segment=('1nbt', '1nbt'))
mdl2.pick_atoms(aln=aln, atom_types='CA')
mdl2.superpose(mdl2, aln)
mdl2.write(file='1nbt.fit.pdb')

```

4.9.16 alignment.malign() — align two or more sequences

Options:

<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>gap_penalties_1d = <float : 2></code>	900 50	gap creation and extension penalties for sequence/sequence alignment
<code>off_diagonal = <int : 1></code>	100	to speed up the alignment
<code>align_block = <int : 1></code>	0	the last sequence in the first block of sequences
<code>matrix_offset = <float : 1></code>	0.00	substitution matrix offset for local alignment
<code>overhang = <int : 1></code>	0	un-penalized overhangs in protein comparisons
<code>local_alignment = <bool : 1></code>	<code>False</code>	whether to do local as opposed to global alignment

Description: This command performs a multiple sequence alignment. The sequences to be aligned are the sequences in the current alignment arrays. The command uses the dynamic programming method for the best sequence alignment, given the gap initiation and extension penalties specified by `gap_penalties_1d`, and residue type weights read from file `rr_file`. See command `alignment.align()` for more information.

The algorithm for the multiple alignment is as follows. First, sequence 2 is aligned with sequence 1 (*i.e.*, block of sequences from 1–`align_block`). Next, sequence 3 is aligned with an average of the aligned sequences 1 and 2; *i.e.*, the weight matrix is an average of the weights 1–3 and 2–3. For this averaging, the gap–residue and gap–gap weights are obtained from the residue–residue weight matrix file, not from gap penalties. If the corresponding weights are not in the file, they are set to the worst and best residue–residue score, respectively.

See instructions for `alignment.align()` for more details.

Example: examples/commands/malign.py

```
# Example for: alignment.malign()
```

```
# This will read all sequences from a file, align them, and write
# the alignment to a new file:

env = environ()

aln = alignment(env, file='toxin.ali', align_codes='all')
aln.malign(gap_penalties_1d=(-600, -400))
aln.write(file='toxin-seq.pap', alignment_format='PAP')
```

4.9.17 alignment.consensus() — consensus sequence alignment

Options:

gap_penalties_1d = <float : 2>	900 50	gap creation and extension penalties for sequence/sequence alignment
align_block = <int : 1>	0	the last sequence in the first block of sequences
max_gap_length = <int : 1>	999999	maximal length of gap in protein comparisons
read_weights = <bool : 1>	False	whether to read the whole NxM weight matrix for ALIGN*
write_weights = <bool : 1>	False	whether to write the whole NxM weight matrix for ALIGN*
input_weights_file = <str : 1>	''	External weight matrix input to MODELLER (SALIGN/ALIGN)
output_weights_file = <str : 1>	''	File into which the weight file is written (iff WRITE_WEIGHTS = 'on')

Description: This command is similar to **alignment.align()** except that a consensus alignment of two blocks of sequences is produced. A consensus alignment is obtained from a consensus similarity matrix using the specified gap penalties and the global dynamic programming method. The consensus similarity matrix is obtained by aligning the two blocks of sequences many times with different parameters and methods and counting how many times each pair was aligned. This command is still experimental and no detailed description is given at this time.

This command also produces the alignment accuracy that can be printed out by the **alignment.write()** command in the 'PAP' format (0 inaccurate, 9 accurate). If the gap initiation penalty is 0, the gap extension penalty of say 0.4 means that only those positions will be equivalenced that were aligned in at least 80% of the individual alignments (*i.e.*, 2 times 0.40).

Example: examples/commands/align_consensus.py

```
# Example for: alignment.consensus()

# This will read 2 sequences and prepare a consensus alignment
# from many different pairwise alignments.

env = environ()

aln = alignment(env)
aln.append(file='toxin.ali', align_codes=('2ctx', '2abx'))
aln.consensus(gap_penalties_1d=(0, 0.4), align_block=1)
aln.write(file='toxin-seq.pap', alignment_format='PAP')
```

4.9.18 `alignment.compare_structures()` — compare 3D structures given alignment

Options:

<code>edat = <energy_data></code>		objective function parameters
<code>io = <io_data></code>		Options for reading atom files
<code>output = <str : 1></code>	<code>'LONG'</code>	selects output: <code>'SHORT'</code> <code>'LONG'</code> <code>'RMS'</code> <code>'DRMS'</code>
<code>matrix_file = <str : 1></code>	<code>'family.mat'</code>	the filename of the pairwise distance matrix
<code>compare_mode = <int : 1></code>	<code>3</code>	selects the type of comparison: <code>1</code> <code>2</code> <code>3</code>
<code>rms_cutoffs = <float : 11></code>	<code>3.5 3.5 60 60 15 60 60 60 60 60 60</code>	cutoffs for RMS, DRMS, Alpha Phi Psi Omega chi1 chi2 chi3 chi4 chi5
<code>fit_atoms = <str : 1></code>	<code>'CA'</code>	whether to superpose before comparison
<code>distance_atoms = <str : 2></code>	<code>'CA' 'CA'</code>	atom type used for variability calculations
<code>fit = <bool : 1></code>	<code>True</code>	whether to do pairwise least-squares fitting or ALIGN2D alignment
<code>asgl_output = <bool : 1></code>	<code>False</code>	whether to write output for ASGL
<code>refine_local = <bool : 1></code>	<code>True</code>	whether to refine superposition using locally similar substructures

Description: This command compares the structures in the given alignment. It does not make an alignment, but it calculates the RMS and DRMS deviations between atomic positions and distances, and class differences between the mainchain and sidechain dihedral angles. In contrast to the `model.superpose()` command, `alignment.compare_structures()` works with a multiple alignment and it writes more information about the pairwise comparisons.

`output` selects short (`'SHORT'`) or long (`'LONG'`) form of output to the `log` file. If it contains word `'RMS'` or `'DRMS'` it also outputs the RMS or DRMS deviation matrix to file `matrix_file`. This file can be used with the PHYLIP program or with the `environ.dendrogram()` or `environ.principal_components()` commands of MODELLER to calculate a clustering of the structures.

`compare_mode` selects the form of the positional variability calculated for each position along the sequence:

- 1, for true RMS deviation over all proteins that have a residue at the current position. This does not make any sense for periodic quantities like dihedral angles.
- 2, for the average absolute distance over all pairs of residues that have a residue at the current position.
- 3, the same as 2 except that average distance, not its absolute value is used (convenient for comparison of 2 structures to get the \pm sign of the changes for dihedral angles and distances).

`rms_cutoffs` specifies cutoff values for calculation of the position, distance, and dihedral angle RMS deviations for pairwise overall comparisons. If difference between two equivalent points is larger than cutoff it is not included in the RMS sum. The order of cutoffs in this vector is: atomic position, intra-molecular distance, α , Φ , Ψ , ω , χ_1 , χ_2 , χ_3 , χ_4 , and χ_5 (there are 5 dihedrals in a disulfide bridge), where α is the virtual C_α dihedral angle between four consecutive C_α atoms. These cutoffs do not affect positional variability calculations.

`fit_atoms` string specifies all the atom types (including possibly a generic `'ALL'`) to be fitted in the least-squares superposition. These atom types are used in the least-squares superposition, and in calculation of the position and distance RMS deviations.

`distance_atoms[1]` specifies the atom type that is used for getting the average structure and RMS deviation at each alignment position in the ASGL output file `'posdif.asgl'`. This ASGL file contains the positional variability of the selected atom type in the family of compared proteins. The ASGL output files can then be used with ASGL scripts `'posdif'` and `'dih'` to produce POSTSCRIPT plots of the corresponding variabilities at each alignment position. `asgl_output` has to be `True` to obtain the ASGL output files.

If `fit = True`, a least-squares superposition is done before the comparisons; otherwise, the orientation of the molecules in the input atom files is used.

Example: See `alignment.malign3d()` command.

4.9.19 `alignment.align3d()` — align two structures

Options:

<code>io = <io_data></code>		Options for reading atom files
<code>gap_penalties_3d = <float : 2></code>	0.0 1.75	gap creation and extension penalties for structure/structure superposition
<code>fit_atoms = <str : 1></code>	'CA'	one atom type used for superposition
<code>fit = <bool : 1></code>	True	whether to align
<code>output = <str : 1></code>	'LONG'	'SHORT' 'LONG' 'VERY_LONG'
<code>align3d_trf = <bool : 1></code>	False	whether to transform the distances before dynamic programming
<code>align3d_repeat = <bool : 1></code>	False	do several starts to maximize number of equivalent positions
<code>off_diagonal = <int : 1></code>	100	to speed up the alignment
<code>matrix_offset = <float : 1></code>	0.00	substitution matrix offset for local alignment
<code>overhang = <int : 1></code>	0	un-penalized overhangs in protein comparisons
<code>local_alignment = <bool : 1></code>	False	whether to do local as opposed to global alignment

Description: This command uses the current alignment as the starting point for an iterative least-squares superposition of two 3D structures. This results in a new pairwise structural alignment. A good initial alignment may be obtained by sequence alignment (`alignment.align()`). For superpositions, only one atom per residue is used, as specified by `fit_atoms[1]`.

The alignment algorithm is as follows. First, structure 2 is least-squares fit on structure 1 using all the equivalent residue positions in the initial alignment that have the specified atom type. Next, the residue-residue distance matrix is obtained by calculating Euclidean distances between all pairs of selected atoms from the two structures. The alignment of the two structures is then obtained by the standard dynamic programming optimization based on the residue-residue distance matrix.

`gap_penalties_3d[1]` is a gap creation penalty (usually 0), and `gap_penalties_3d[2]` is a gap extension penalty, say 1.75. This procedure identifies pairs of residues as equivalent when they have their selected atoms at most 2 times `gap_penalties_3d[2]` angstroms apart in the current orientation (this is so when the gap initiation penalty is 0). The reason is that an equivalence costs the distance between the two residues while an alternative, the gap-residue and residue-gap matches, costs twice the gap extension penalty.

From the dynamic programming run, a new alignment is obtained. Thus, structure 2 can be fitted onto structure 1 again, using this new alignment, and the whole cycle is repeated until there is no change in the number of equivalent positions and until the difference in the rotation matrices for the last two superpositions is very small. At the end, the framework, that is the alignment positions without gaps, is written to the `log` file.

If `fit` is `False`, no alignment is done.

If `output` contains 'SHORT', only the best alignment and its summary are displayed. If `output` contains 'LONG', summaries are displayed for all initial alignments in each framework cycle. If `output` contains 'VERY_LONG', all alignments are displayed.

If `align3d_trf` is `True`, the weights in the weight matrix are modified distances [Subbiah *et al.*, 1993].

If `align3d_repeat` is `True`, three additional initial alignments are tried and the one resulting in the largest number of equivalent positions is selected.

Example: `examples/commands/align3d.py`

```

# Example for: alignment.align3d(), model.superpose()

# This will align 3D structures of two proteins:

log.verbose()
env = environ()

# First example: read sequences from a sequence file:
aln = alignment(env)
aln.append(file='toxin.ali', align_codes=['1fas', '2ctx'])
aln.align(gap_penalties_1d=[-600, -400])
aln.align3d(gap_penalties_3d=[0, 4.0])
aln.write(file='toxin-str.ali')

# Second example: read sequences from PDB files to eliminate the
# need for the toxin.ali sequence file:
mdl = model(env)
aln = alignment(env)
for code in ['1fas', '2ctx']:
    mdl.read(file=code)
    aln.append_model(mdl, align_codes=code, atom_files=code)
aln.align(gap_penalties_1d=[-600, -400])
aln.align3d(gap_penalties_3d=(0, 2.0))
aln.write(file='toxin-str.ali')

# And now superpose the two structures using current alignment to get
# various RMS's:
mdl = model(env, file='1fas')
mdl.pick_atoms(aln, atom_types='CA')
mdl2 = model(env, file='2ctx')
mdl.superpose(mdl2, aln)

```

4.9.20 alignment.malign3d() — align two or more structures

Options:

io = <io_data>		Options for reading atom files
gap_penalties_3d = <float : 2>	0.0 1.75	gap creation and extension penalties for structure/structure superposition
off_diagonal = <int : 1>	100	to speed up the alignment
matrix_offset = <float : 1>	0.00	substitution matrix offset for local alignment
overhang = <int : 1>	0	un-penalized overhangs in protein comparisons
local_alignment = <bool : 1>	False	whether to do local as opposed to global alignment
fit_atoms = <str : 1>	'CA'	one atom type for superposition
fit = <bool : 1>	True	whether to align
output = <str : 1>	'LONG'	'SHORT' 'LONG' 'VERY_LONG' 'NO_ALIGNMENT'
write_fit = <bool : 1>	False	whether to write out fitted coordinates to .fit files
edit_file_ext = <str : 2>	' .pdb' ' .fit.pdb'	old and new file extensions for filename construction in MALIGN3D

<code>current_directory = <bool : 1></code>	<code>True</code>	whether to write output .fit files to current directory
<code>write_whole_pdb = <bool : 1></code>	<code>True</code>	whether to write out all lines in the input PDB file

Description: This command uses the current alignment as the starting point for an iterative least-squares superposition of two or more 3D structures. This results in a new multiple structural alignment. A good initial alignment may be obtained by sequence alignment (`alignment.malign()`). For superpositions, only one atom per residue is used, as specified by `fit_atoms`. The resulting alignment can be written to a file with the `alignment.write()` command. The multiply superposed coordinates remain in memory and can be used with such commands as `model.transfer_xyz()` if `alignment[x].atom.file` is not changed in the meantime. It is best to use the structure that overlaps most with all the other structures as the first protein in the alignment. This may prevent an error exit due to too few equivalent positions during framework construction.

The alignment algorithm is as follows. There are several cycles, each of which consists of an update of a framework and a calculation of a new alignment; the new alignment is based on the superposition of the structures onto the latest framework. The framework in each cycle is obtained as follows. The initial framework consists of the atoms in structure 1 that correspond to `fit_atoms`. If there is no specified atom types in any of the residues at a given position, the coordinates for this framework position are approximated by the neighboring coordinates. Next, all other structures are fit to this framework. The final framework for the current cycle is then obtained as an average of all the structures, in their fitted orientations, but only for residue positions that are common to all of them, given the current alignment. Another result is that all the structures are now superposed on this framework. Note that the alignment has not been changed yet. Next, the multiple alignment itself is re-derived in $N - 1$ dynamic programming runs, where N is the number of structures. This is done as follows. First, structure 2 is aligned with structure 1, using the inter-molecular atom-atom distance matrix, for all atoms of the selected type, as the weight matrix for the dynamic programming run. Next, structure 3 is aligned with an average of structures 1 and 2 using the same dynamic programming technique. Structure 4 is then aligned with an average of structures 1-3, and so on. Averages of structures $i-j$ are calculated for all alignment positions where there is at least one residue in any of the structures $i-j$ (this is different from a framework which requires that residues from all structures be present). Note that in this step, residues out of the current framework may get aligned and the current framework residues may get unaligned. Thus, after the series of $N - 1$ dynamic programming runs, a new multiple alignment is obtained. This is then used in the next cycle to obtain the next framework and the next alignment. The cycles are repeated until there is no change in the number of equivalent positions. This procedure is best viewed as a way to determine the framework regions, not the whole alignment. The results from this command are expected to be similar to the output of program MNYFIT [Sutcliffe *et al.*, 1987].

`gap_penalties_3d[1]` is a gap creation penalty (usually 0), and `gap_penalties_3d[2]` is a gap extension penalty, say 1.75. This procedure identifies pairs of positions as equivalent when they have their selected atoms at most 2 times `gap_penalties_3d[2]` angstroms apart in the current superposition (this is so when the gap initiation penalty is 0), as described for the `alignment.align3d()` command.

Argument output can contain the following values:

- 'SHORT', only the final framework is written to the `log` file.
- 'LONG', the framework after the alignment stage in each cycle is written to the `log` file.
- 'VERY_LONG', the framework from the framework stage in each cycle is also written to the `log`.

If `write_fit` is `True`, the fitted atom files are written out in their final fitted orientations. To construct the file-names, first the file extension in `edit_file_ext[1]` is removed (if present), and then the extension in `edit_file_ext[2]` is added. By default this creates files with a `_fit` extension.

If `current_directory` is `True`, the fitted atom files will go to the current directory. Otherwise, the output will be in the directory with the original files.

If `write_whole_pdb` is `True`, the whole PDB files are written out; otherwise only the parts corresponding to the aligned sequences are output.

If `fit` is `False`, the initial alignment is not changed. This is useful when all the structures have to be superimposed with the initial alignment (`fit = False` and `write_fit = True`).

Example: examples/commands/malign3d.py

```

# Example for: alignment.malign3d(), alignment.compare_structures()

# This will read all sequences from a sequence file, multiply align
# their 3D structures, and then also compare them using this alignment.

env = environ()

aln = alignment(env, file='toxin.ali', align_codes='all')
aln.malign(gap_penalties_1d=(-600, -400))
aln.malign3d(gap_penalties_3d=(0, 2.0), write_fit=True, write_whole_pdb=False)
aln.write(file='toxin-str.pap', alignment_format='PAP')

# Make two comparisons: no cutoffs, and 3.5A/60 degree cutoffs for RMS, DRMS,
# and dihedral angle comparisons:
aln.compare_structures(rms_cutoffs=[999]*11)
aln.compare_structures(rms_cutoffs=(3.5, 3.5, 60, 60, 60, 60, 60, 60, 60,
                                     60, 60))

```

4.9.21 alignment.salign() — align two or more sequences/structures of proteins**Options:**

io = <io_data>		Options for reading atom files
align3d_trf = <bool : 1>	False	whether to transform the distances before dynamic programming
align_block = <int : 1>	0	the last sequence in the first block of sequences
alignment_type = <str : 1>	'PROGRESSIVE'	'PAIRWISE' 'TREE' 'PROGRESSIVE' for SALIGN
align_what = <str : 1>	'BLOCK'	what to align in ALIGN; 'BLOCK' 'ALIGNMENT' 'LAST' 'PROFILE'
auto_overhang = <bool : 1>	False	overhang values made dependent on sequence length difference
comparison_type = <str : 1>	'MAT'	'MAT' or 'PSSM' for comparing matrices or PSSMs when profiles are compared
current_directory = <bool : 1>	True	whether to write output .fit files to current directory
dendrogram_file = <str : 1>	''	File into which the SALIGN dendrogram is written out
ext_tree = <bool : 1>	False	To read a pre computed tree for SALIGN
feature_weights = <float : 6>	1 0 0 0 0 0	feature weights for SALIGN
fit_atoms = <str : 1>	'CA'	atom type(s) being superposed
fit_on_first = <bool : 1>	False	whether or not all structures are to be fit on the first structure, given the final alignment
fit_pdbnam = <bool : 1>	True	whether or not to add _fit to the PDB file name in output alifile by SALIGN
fit = <bool : 1>	True	whether to do pairwise least-squares fitting or ALIGN2D alignment
fix_offsets = <float : 5>	0 1 2 3 4	offsets of the ALIGN2D alignment score for "fixed" positions indicated by '123456789' in line '_fix_pos'

<code>gap_function = <bool : 1></code>	False	whether or not to switch on functional gap penalty in <code>salgn</code>
<code>gap_gap_score = <float : 1></code>	0	dissimilarity score for aligning gap with gap, in <code>SALIGN</code>
<code>gap_penalties_1d = <float : 2></code>	900 50	gap creation and extension penalties for sequence/sequence alignment
<code>gap_penalties_2d = <float : 9></code>	0.35 1.2 0.9 1.2 0.6 8.6 1.2 0 0	gap penalties for sequence/structure alignment: helix, beta, accessibility, straightness, and CA-CA distance factor, <code>dst min</code> , <code>dst power</code> , <code>t</code> , <code>structure_profile</code> ; best U,V=-450,0
<code>gap_penalties_3d = <float : 2></code>	0.0 1.75	gap creation and extension penalties for structure/structure superposition
<code>gap_residue_score = type1</code>	0	dissimilarity score for aligning gap with residue, in <code>SALIGN</code>
<code>improve_alignment = <bool : 1></code>	True	whether or not to optimize alignment in <code>SALIGN</code>
<code>input_weights_file = <str : 1></code>	''	External weight matrix input to <code>MODELLER (SALIGN/ALIGN)</code>
<code>local_alignment = <bool : 1></code>	False	whether to do local as opposed to global alignment
<code>matrix_comparison = <str : 1></code>	'CC'	'CC', 'MAX', 'AVE', - kinds of matrix comparisons
<code>matrix_offset = <float : 1></code>	0.00	substitution matrix offset for local alignment
<code>max_gap_length = <int : 1></code>	999999	maximal length of gap in protein comparisons
<code>normalize_pp_scores = <bool : 1></code>	False	whether or not to normalize position-position scores in <code>SALIGN</code>
<code>no_ter = <bool : 1></code>	False	whether to not write TER into PDB
<code>n_subopt = <int : 1></code>	1	number of optimal and suboptimal alignments <code>ALIGN/ALIGN2D</code>
<code>off_diagonal = <int : 1></code>	100	to speed up the alignment
<code>output_weights_file = <str : 1></code>	''	File into which the weight file is written (iff <code>WRITE_WEIGHTS = 'on'</code>)
<code>output = <str : 1></code>	'LONG'	what and/or how to output
<code>overhang = <int : 1></code>	0	un-penalized overhangs in protein comparisons
<code>overhang_auto_limit = <int : 1></code>	60	<code>auto_overhang</code> effective if seq length diff greater than <code>overhang_auto_limit</code>
<code>overhang_factor = <float : 1></code>	0.4	Factor to multiply seq. length diff with when <code>auto_overhang</code> is on
<code>rms_cutoffs = <float : 11></code>	3.5 3.5 60 60 15 60 60 60 60 60 60	cutoffs for RMS, DRMS, Alpha Phi Psi Omega chi1 chi2 chi3 chi4 chi5
<code>rr_file = <str : 1></code>	'\$(LIB)/as1.sim.mat'	input residue-residue scoring file
<code>similarity_flag = <bool : 1></code>	False	when turned on, the <code>SALIGN</code> command does not convert numbers into a distance sense.
<code>smooth_prof_weight = <float : 1></code>	10	for smoothing the profile aa frequency with a prior
<code>subopt_offset = <float : 1></code>	0.0	offset for residue-residue score in getting suboptimals in <code>ALIGN/ALIGN2D</code>
<code>substitution = <bool : 1></code>	False	whether to use the background in <code>PSSM</code> comparison

<code>weights_type = <str : 1></code>	<code>'SIMILAR'</code>	or <code>'DISTANCE'</code> for the kind of substitution values
<code>write_fit = <bool : 1></code>	<code>False</code>	whether to write out fitted coordinates to <code>.fit</code> files
<code>write_weights = <bool : 1></code>	<code>False</code>	whether to write the whole NxM weight matrix for ALIGN*
<code>write_whole_pdb = <bool : 1></code>	<code>True</code>	whether to write out all lines in the input PDB file

Description: This command is a general dynamic programming based alignment procedure for aligning sequences, structures or a combination of the two. It is loosely based on the program COMPARER [Šali & Blundell, 1990]. SALIGN can be used to generate multiple protein structures/sequences alignments or to align two blocks of sequences/structures that are in memory.

Please note that the method is still in development, and has not yet been fully benchmarked. As with any other alignment method, generated alignments should be assessed for quality.

Broadly classifying, three different types of protein alignment categories are tackled by this command:

1. Multiple structure alignments
2. Aligning a structure block to a sequence block
3. Multiple and pair-wise protein sequence alignment

The command incorporates the functionality of several old MODELLER commands (`alignment.align()`, `alignment.align2d()`, `alignment.malign()`, `alignment.align3d()`, and `alignment.malign3d()`). Some of the examples below illustrate the equivalent script files to replace the old alignment commands with `alignment.salign()`.

In addition to these, this command has several new alignment features including profile-profile sequence alignments and a dendrogram based multiple sequence/structure alignment among others.

All pair-wise alignments make use of local or global dynamic programming. A switch from one to another can be effected by setting `local_alignment` to `True` or `False`. The dynamic programming can be carried out using affine gap penalties (as previously used in `alignment.align()`) or an environment dependent gap penalty function (as used in `alignment.align2d()`). The choice of gap penalty types can be regulated by switching the variable `gap_function` on or off. All arguments that associated to the `alignment.align()` and `alignment.align2d()` commands apply.

If at least one of the blocks in a pairwise alignment consists of structures, dynamic programming can be performed using structure dependent gap penalties.

Features of proteins used for alignment

Central to the dynamic programming algorithm is the weight matrix. In SALIGN, this matrix is constructed by weighting the contribution from six features of protein structure and sequence:

Feature 1 is the residue type. $W_{i,j}^1$ is obtained from the residue type – residue type dissimilarity matrix, specified in the file `rr_file`. $W_{i,j}^1$ dissimilarity score for positions i and j in the two compared sub-alignments is the average dissimilarity score for a comparison of all residues in one sub-alignment with all residues in the other sub-alignment (note that gaps are ignored here). Should only feature weight 1 be non-zero, the user has an option of considering residue-residue similarity scores instead of distance scores by setting `similarity_flag` to `True`.

Feature 2 is the inter-molecular distance for a pair of residues (unless `align3d_trf` is `True`: see `alignment.align3d()`). Only one atom per residue is of course selected, as specified by `fit_atoms` (e.g., C_α , although we should also allow for C_β in the future, which requires an intervention for Gly). This ‘position’ feature is complicated because it depends on the relative orientation of the structures corresponding to the two compared alignments. $W_{i,j}^2$ is the Euclidean distance between the compared positions i and j in the two compared sub-alignments that are already optimally aligned and superposed based on their coordinates alone. This optimal alignment is obtained by an iterative procedure as follows (the same as

in `alignment.align3d()`). The average structures for both sub-alignments are calculated for all sub-alignment positions with at least one defined selected atom. This calculation is straightforward because the structures within the two sub-alignments are already superposed with each other (see below). Then, the distance matrix for dynamic programming with affine gap penalties is calculated as the matrix of Euclidean distances between the two averages. The dynamic programming results into a new alignment, dependent also on the gap initiation and extension penalties `gap_penalties_3d` (a reasonable setting is (0, 3)). `gap_penalties_3d[0]` is a gap creation penalty (usually 0), and `gap_penalties_3d[1]` is a gap extension penalty, say 3. When the gap initiation penalty is 0, pairs of positions are identified as equivalent when they have their selected atoms at most 2 times `gap_penalties_3d[0]` angstroms apart in the current superposition, as described for the `alignment.align3d()` command. The new alignment is then used to generate the new superposition of the two averages, and the iteration of the distance matrix calculation, alignment and superposition is repeated until there are no changes in the number of equivalent positions and in the rotation matrix relating the two averages.

The values of both `improve_alignment` and `fit` are used in the calculation of the position feature. That is, the initial alignment and the orientation of the coordinates can be selected not to change at all during the calculation of the inter-molecular distance matrix.

When the calculation of the inter-molecular distance matrix is finished, all the structures in the second sub-alignment are rotated and translated following the optimal rotation and translation of the second average on the first average. These superpositions prepare the individual structures for the next of the $n - 1$ stages of the progressive multiple alignment, and also orient all the structures for writing out to `'alignment[i].code_fit.pdb'` atom files if `write_fit = True`. If `fit_pdbnam = False`, the PDB filenames in the output alignment file will not have the `'_fit.pdb'` extensions. Thus, feature 2 needs to be selected by `feature_weight[2] > 0` if you wish to write out the structures superposed according to the tree-following procedure; also, `fit_on_first` must be `False`, otherwise the structures are written out superposed on the first structure according to the final alignment (see also below).

The alignment produced within the routine that calculates W^2 does not generally correspond to the alignment calculated based on W . Therefore, the multiply superposed structures are not necessarily superposed based on the final multiple alignment produced by `alignment.salign()`. If you wish such a superposition, you can use `alignment.malign3d()` with `fit = False` and `write_fit = True` (the meaning of `fit` is different between `alignment.salign()` and `alignment.malign3d()`).

Unless the position feature is selected, the initial alignment does not matter. If the position feature is selected, a good starting alignment is a multiple sequence alignment, obtained either by `alignment.malign()` or by `alignment.salign()` used without the position feature (the initial alignment can also be prepared using the position feature). If the position feature is used, each pair of structures needs to have at least 3 aligned residues at all points during the alignment.

There are several possibilities as to the final orientation of the input coordinates. If `fit_on_first` is `True`, all the coordinate sets are superposed on the first structure, using the final multi-feature multiple alignment. If `fit_on_first` is `False`, and position feature was used, and `fit` was `True`, the coordinates will be superposed in the progressive manner guided by the tree, by the routine that calculates the inter-molecular distance matrices; this superposition is based only on the positions of the selected atoms (feature 2), not on other features such as residue type, secondary, structure, *etc.* If `improve_alignment` is `False`, it does not make much sense to have `fit = True` (use `fit_on_first = True`).

Feature 3 is the fractional sidechain accessibility. The pair-wise residue-residue dissimilarity is calculated by classifying residues into the buried (< 15%), semi-exposed, and exposed classes (> 30%). The dissimilarity is 0 for equal classes, 1 for neighboring classes and 2 for the buried-exposed match. The position-position dissimilarity is the average residue-residue dissimilarity for comparing all residues from one group to all residues in the other group (gaps are ignored).

Feature 4 is the secondary structure type, distinguishing between helix, strand, and other. The pair-wise residue-residue dissimilarity is 0 for equal classes, 1 for 'helix' or 'strand' matched to 'other', and 2 for 'helix' matched to 'strand'. Position-position dissimilarity is calculated in the same way as for feature 3.

Feature 5 is the local conformation. A pair-wise residue-residue score is DRMSD between the selected atoms (`fit_atoms`) from the segments of 5 residues centered on the two matched residues. Position-position dissimilarity is calculated in the same way as for feature 3.

Feature 6 is a user specified feature for which an external matrix (in MODELLER matrix format; see the substitution matrices in the `modlib` directory for examples) has to be specified using `input_weights_file`. The user can input either a similarity matrix (`weights_type = SIMILAR`) or a distance matrix (`weights_type = DISTANCE`).

Alignment of protein sequences

- **Alignment of two sequence blocks**

For all alignments that are to be performed using only sequence information all features but the first should be made 0. Pairs of sequence blocks are aligned using SALIGN the same way in which `alignment.align()` aligned sequence blocks. The example in the next section illustrates how SALIGN accomplishes this. `gap_function` should be kept `False` to align sequence blocks locally or globally using affine gap penalties.

- **Alignment of protein sequences by their profiles**

Two blocks of sequences can be aligned using the information contained within each of the multiple sequence blocks [Martí-Renom *et al.*, 2004]. `align_block` demarcates the end of the first block and `align_what` is set to `profile` indicating that the blocks will be aligned using their profiles. Also, since this kind of alignment is effected only between two blocks, `alignment_type` is set to `pairwise`. The weight matrix for dynamic programming is created by comparing the sequence information in the two blocks. Two kinds of comparisons can be performed:

1. A correlation coefficient of the variation of the the 20 amino acids at each position (`alignment_type` is set to `PSSM`).
2. Comparing the residue substitution matrices implied at each position of the two blocks (`alignment_type` is set to `MAT`).

Matrix comparisons are of three types: taking the maximum, average or correlation coefficient of residue-residue substitution at the aligned positions (`matrix_comparison` set to `MAX`, `AVE` or `CC`).

Example: `examples/salign/salign_profile_profile.py`

```
# profile-profile alignment using salign

log.level(1, 0, 1, 1, 1)
env = environ()

aln = alignment(env, file='mega_prune.faa', alignment_format='FASTA')

aln.salign(rr_file='${LIB}/blosum62.sim.mat',
           gap_penalties_1d=(-500, 0), output='',
           align_block=15, # no. of seqs. in first MSA
           align_what='PROFILE',
           alignment_type='PAIRWISE',
           comparison_type='PSSM', # or 'MAT' (Caution: Method NOT benchmarked
                                   # for 'MAT')
           similarity_flag=True, # The score matrix is not rescaled
           substitution=True, # The BLOSUM62 substitution values are
                               # multiplied to the corr. coef.
           #write_weights=True,
           #output_weights_file='test.mtx', # optional, to write weight matrix
           smooth_prof_weight=0.10) # For mixing data with priors

#write out aligned profiles (MSA)
aln.write(file='salign.ali', alignment_format='PIR')

# Make a pairwise alignment of two sequences
aln = alignment(env, file='salign.ali', alignment_format='PIR',
               align_codes=('12asA', '1b8aA'))
aln.write(file='salign_pair.ali', alignment_format='PIR')
```

```
aln.write(file='salign_pair.pap', alignment_format='PAP')
```

- **Multiple Protein sequence alignment**

Multiple alignments of sequences is similar to multiply aligning structures, the difference being that for sequence alignments only the first feature is non-zero. The example in the structure alignment section can illustrate multiple sequence alignments as well.

Alignment of protein structures with sequences

As stated earlier, all `alignment.align` and `alignment.align2d` related commands apply to `alignment.salign` too. The example below is a `alignment.salign` equivalent of `alignment.align2d` (and `alignment.align`). For a description of the `gap_penalties_2d` see the section on `alignment.align2d`.

Example: examples/salign/salign_align2d.py

```
# align2d/align using salign

# parameters to be input by the user
# 1. gap_penalties_1d
# 2. gap_penalties_2d = (4.35, 1.2, 0.9, 1.2, 0.6, 8.6, 1.2, 0., 0.) (default)
# 3. input alignment file

log.verbose()
env = environ()
env.io.atom_files_directory='./'

aln = alignment(env, file='align2d_in.ali', align_codes='all')
aln.salign(rr_file='${LIB}/as1.sim.mat', # Substitution matrix used
           output='',
           max_gap_length=20,
           gap_function=True,          # If False then align2d not done
           feature_weights=(1., 0., 0., 0., 0., 0.),
           gap_penalties_1d=(-1000, 0),
           gap_penalties_2d=(4.35, 1.2, 0.9, 1.2, 0.6, 8.6, 1.2, 0.0, 0.0),
           # d.p. score matrix
           #write_weights=True, output_weights_file='salign.mtx'
           similarity_flag=True)      # Ensuring that the dynamic programming
                                     # matrix is not scaled to a difference matrix

aln.write(file='align2d.ali', alignment_format='PIR')
aln.write(file='align2d.pap', alignment_format='PAP')
```

Alignment of protein structures

Structure alignments can make use of all the 5 structure/sequence features as well as the 6th user provided feature matrix. Pairwise alignments of structures can make use of the constant gap penalties or the environment dependent gap penalties. Multiple structure alignments are constructed from pairwise structure alignments.

- **Pairwise protein structure alignment**

For optimal pairwise alignments it is suggested to call SALIGN multiple times, typically 2-3 times. The first SALIGN call will give an initial alignment which is refined in the subsequent calls. Usually, feature 2 is made non-zero only during the 'refinement' stage as rigid body refinement is done at the position implied by the alignment in memory.

- **Tree Multiple Structure Alignments**

When `alignment_type` is set to `tree`, a dendrogram of the n proteins in memory is calculated using the selected features. The multiple alignment is then a progression of $n - 1$ pairwise alignments of the growing sub-alignments. A sub-alignment is an alignment of $< n$ proteins. The pairwise alignment of two sub-alignments is achieved using affine or environment dependent gap penalties, depending on whether `gap_function` is set to `False` or `True` (arguments to the `alignment.align()` and `alignment.align2d()` commands apply).

- **Progressive Multiple Structure Alignments**

If `alignment_type` is set to `progressive`, the multiple alignment follows only the last part of the ‘tree’ alignment where, in $n - 1$ alignments, all the structures/sequences are successively aligned to the first one on the list.

The alignment of proteins within a sub-alignment does not change when the sub-alignment is aligned with another protein or sub-alignment. The pairwise alignment of sub-alignments is guided by the dendrogram. First, the most similar pair of proteins are aligned. Second, the next most similar pair of proteins are aligned, or the third protein is aligned with the sub-alignment of the first two, as indicated by the dendrogram. This greedy, progressive procedure requires $n - 1$ steps to align all n proteins, and each step requires a pairwise alignment of two sub-alignments.

If in a multiple alignment, overhangs are to be penalized differently for the pairs of alignments that create the multiple, `auto_overhang` can be set to `True`. This will ensure that the value of `overhang` changes as `overhang_factor` times the numerical difference in the residues of the pair. Further, this is only effected if the difference is greater than `overhang_auto_limit`.

The dendrogram can be written out in a separate file by specifying the file name to `dendrogram_file`.

Example: examples/salign/salign_multiple.py

```
# Illustrates the SALIGN multiple structure/sequence alignment

log.verbose()
env = environ()
env.io.atom_files_directory = './:../atom_files/'

aln = alignment(env)
for (code, chain) in (('1is4', 'A'), ('1uld', 'D'), ('1ulf', 'B'),
                    ('1ulg', 'B'), ('1is5', 'A')):
    mdl = model(env, file=code, model_segment=('FIRST:'+chain, 'LAST:'+chain))
    aln.append_model(mdl, atom_files=code, align_codes=code+chain)

for (weights, write_fit, whole) in (((1., 0., 0., 0., 1., 0.), False, True),
                                   ((1., 0.5, 1., 1., 1., 0.), False, True),
                                   ((1., 1., 1., 1., 1., 0.), True, False)):
    aln.salign(rms_cutoffs=(3.5, 6., 60, 60, 15, 60, 60, 60, 60, 60),
              normalize_pp_scores=False,
              rr_file='${LIB}/as1.sim.mat', overhang=30,
              gap_penalties_1d=(-450, -50),
              gap_penalties_3d=(0, 3), gap_gap_score=0, gap_residue_score=0,
              dendrogram_file='1is3A.tree',
              alignment_type='tree', # If 'progressive', the tree is not
                                   # computed and all structures will be
                                   # aligned sequentially to the first
#ext_tree=True, # Tree building can be avoided if the tree
              # is input
              #input_weights_file='1is3A_exmat.mtx',
              feature_weights=weights, # For a multiple sequence alignment only
              # the first feature needs to be non-zero
              improve_alignment=True, fit=True, write_fit=write_fit,
```

```

write_whole_pdb=whole, output='ALIGNMENT QUALITY')

aln.write(file='1is3A.pap', alignment_format='PAP')
aln.write(file='1is3A.ali', alignment_format='PIR')

# The number of equivalent positions at different RMS_CUTOFF values can be
# computed by changing the RMS value and keeping all feature weights = 0
aln.salign(rms_cutoffs=(1.0, 6., 60, 60, 15, 60, 60, 60, 60, 60),
           normalize_pp_scores=False, rr_file='$(LIB)/as1.sim.mat', overhang=30,
           gap_penalties_1d=(-450, -50), gap_penalties_3d=(0, 3),
           gap_gap_score=0, gap_residue_score=0, dendrogram_file='1is3A.tree',
           alignment_type='progressive', feature_weights=[0]*6,
           improve_alignment=False, fit=False, write_fit=True,
           write_whole_pdb=False, output='QUALITY')

```

Sub-optimal alignments

The weight matrix can be offset at random, many times over, to generate several 'sub-optimal' alignments. The number of sub-optimal alignments to be output can be specified with `n_subopt`. Though the matrix positions at which these offsets are applied cannot be controlled, the user can choose by how much the matrix will be offset (`subopt_offset`). The output suboptimal alignments are written into the file 'suboptimal_alignments.out' and is of the format that specifies the serial residue numbers of the two sequences.

Example: examples/salign/salign_subopt.py

```

log.verbose()
env = environ()
aln = alignment(env, file='fm07254_test.ali', alignment_format='PIR')
aln.salign(feature_weights=(1., 0, 0, 0, 0, 0), gap_penalties_1d=(-450, -50),
           n_subopt = 5, subopt_offset = 15)

```

Alignments using external restraints

Fix positions: The user can choose to have certain alignment positions "fixed" by offsetting the appropriate matrix entries. This is achieved by setting `fix_offset` to `True`. Below the given pair of sequence blocks in the input alignment file, a new pseudo sequence is introduced with the align code `._fix_pos`. The residues of this pseudo sequence are integer values from 0 through 4 (alternatively, a blank is equivalent to 0). The placement of the non-zero integers is done so that they are below the positions that are to be offset, in an effort to ensure that the alignment implied by that position is restrained. The extent to which this alignment position is restrained depends on the offset that the numerical value (0-4) corresponds to. This can also be specified by the user by setting the matrix offset values with `fix_offsets`.

Example: examples/salign/salign_fix_positions.py

```

# Demonstrating the use of alignment restraints, only available in
# align2d and salign:

log.verbose()
env = environ()

# The special alignment entry '_fix_pos' has to be the last entry in the
# alignment array. Its sequence contains characters blank, 1, 2, 3, and 4
# at the restrained alignment positions. The residue-residue score from

```

```

# the substitution matrix for these positions will be offset by the scalar
# value FIX_OFFSETS[1..5].
aln = alignment(env, file='fix_positions.ali', align_codes=('1leh', '3btoA',
                                                         '_fix_pos'))

# fix_offsets specifies the offset corresponding to character '1234' in the
# _fix_pos entry in the alignment
# (this offsets unlabeled positions for 0, the ones indicated by 1 by
# 1000, those indicated by 2 by 2000, etc.)
aln.salign(fix_offsets=(0, 1000, 2000, 3000, 4000),
           gap_penalties_2d=(0, 0, 0, 0, 0, 0, 0, 0, 0), # Any values are
                           # possible here
           local_alignment=False, # Local alignment works, too
           gap_penalties_1d=(-600, -400)) # This is best with the default value
                                           # of gap_penalties_2d

# Write it out, the _fix_pos is erased automatically in salign:
aln.write(file='fix_positions_salign.pap', alignment_format='PAP')

```

External weight matrix: An example of using feature 6.

Example: `examples/salign/salign_external_matrix.py`

```

# Reads an external matrix

log.verbose()
env = environ()

aln = alignment(env, file='1dubA-1nzyA.ali', align_codes='all')

aln.salign(alignment_type='pairwise', output='',
           rr_file='$(LIB)/blosum62.sim.mat',
           #rr_file='$(LIB)/as1.sim.mat',
           #max_gap_length=20,
           gap_function=False,
           input_weights_file='external.mtx', # External weight matrix
           #weights_type='DISTANCE', # type of ext. wgt. mtx
           # ensure appropriate gap penalties for the ext. matrix
           #feature_weights=(1., 0., 0., 0., 0., 0.), gap_penalties_1d=(30, 26),
           #write_weights=True, output_weights_file='score.mtx',
           feature_weights=(1., 0., 0., 0., 0., 1.),
           gap_penalties_1d=(-500, -300))
aln.write(file='output.ali', alignment_format='PIR')
aln.write(file='output.pap', alignment_format='PAP')

```

Multiple structure alignment according to a user specified dendrogram The user has the option of inputting an $n \times n$ matrix from which a dendrogram can be inferred. The multiple tree alignment is then confined to follow this externally input dendrogram. To effect this, `ext.tree` is set to `True` and the external matrix file is specified in `input_weights.file`. A word of caution: it is not advisable to have feature 6 as non zero when an external dendrogram matrix is read.

Gap penalties and correcting for gaps

SALIGN makes use of three sets of gap penalties. `gap_penalty_1d` is for dynamic programming making use of constant gap penalties. `gap_penalty_2d` is when a variable function for gap penalty is used. `gap_penalty_3d` is used along with feature 2 only, when structures are aligned by a least squares fit of their atomic positions. All SALIGN features produce some measure of residue equivalence (similarity or distance scores). The scales of these scores differ depending on the feature used. For optimal usage, `gap_penalties_1d` should be set appropriately considering the features used. For example, if feature 1 is non zero, and the residue substitution matrix used is the BLOSUM62 similarity matrix, `gap_penalties_1d` is set to (-450, -50) and when feature 1 is zero `gap_penalties_1d` is set to (2, 3). A word of caution: gap penalties have not yet been optimized for aligning sequences by their profiles and for structure alignments.

The gap correction function is $g_{i,j} = \frac{n_{rg}}{(n_1 n_2)} r + \frac{n_{gg}}{(n_1 n_2)} g$, where n_1 and n_2 are the number of proteins in the two sub-alignments, n_{rg} is the number of gap-residue pairs, and n_{gg} is the number of gap-gap pairs when comparing protein positions from one sub-alignment with protein position from the other sub-alignment, r is `gap_residue_score` and g is `gap_gap_score`. The smaller (even negative) is `gap_gap_score`, and the larger is `gap_residue_score`, the more will the gaps be aligned with gaps.

Useful SALIGN information and commands

The `alignment.salign()` command uses position-position dissimilarity scores (except when `similarity_flag` is switched on), as opposed to similarity scores. This convention applies to all the features, including the residue-residue similarities read from the `rr_file`; however, if a residue type – residue type similarity matrix is read in, it is automatically converted into the distance matrix by $D = \max_{i,j} S_{i,j} - S$. In addition, it is also scaled linearly such that the residue-residue dissimilarity scores range from 0 to 1 (to facilitate weighting this feature with other features).

For each pairwise alignment, the weight matrix W has dimensions N and M that correspond to the lengths of the sub-alignments to be aligned based on the weight matrix W . The dissimilarity score for aligning position i with position j is calculated as $W_{i,j} = \sum_f [\frac{\omega_f}{\sum_f \omega_f} W_{i,j}^f] + g_{i,j}$, where the sum runs over all selected features f , and g is a function that may be used to correct the $W_{i,j}$ score for the presence of gaps within the sub-alignments (see below). A feature f is selected when its weight ω_f (specified in `feature_weights`) is non-zero. The matrices W^f are normalized to have the mean of 0 and standard deviation of 1 when `normalize_pp_scores` is `True`, but it is recommended not to use this option for now (*i.e.*, use `feature_weights` to scale the contributions of the different features to the final W). The weights of 1 will weigh the different features approximately evenly (the residue-residue dissimilarities of feature 1 are scaled to a range from 0 to 1, the position differences of feature 2 are in angstroms, the fractional solvent accessibility scores of feature 3 and the secondary structure scores of feature 4 range from 0 to 2, and the DRMS difference of feature 5 is expressed in angstroms).

If you enable verbose logging with `log.verbose()`, there will be more output in the 'log' file, such as the dendrogram. The dendrogram can also be written out in a separate file by specifying the file name to `dendrogram_file`.

Argument output can contain the following values:

- 'ALIGNMENT': the alignments in the first $n - 2$ stages of the pairwise alignment of sub-alignments are written out.
- 'QUALITY': the final alignment is used to obtain pairwise least-squares superpositions and the corresponding average and minimal numbers of pairs of aligned residues that are within `rms_cutoffs[0]` Å in all pairs of aligned structures. These numbers can be used as absolute quality measures for the final multiple alignment. This option requires the coordinate files for the aligned proteins.

If `write_fit` is `True`, the fitted atom files are written out in their fitted orientations. For this and other options below, also read the text above.

If `write_weights` is `True`, the dynamic programming weight matrix is written out into the file specified by `output_weights_file`.

If `current_directory` is `True`, the output `_pdb.fit` files will be written to the current directory. Otherwise, the output will be in the directory with the original files.

If `write_whole_pdb` is `True`, the whole PDB files are written out; otherwise only the parts corresponding to the aligned sequences are output.

If `fit` is `False`, the initial superposition is not changed. This is useful when all the structures have to be compared with a given alignment as is, without changing their relative orientation.

If `fit_on_first` is `True`, the structures are fit to the first structure according to the final alignment before they are written out.

If `improve_alignment` is `False`, the initial alignment is not changed, though the structures may still be superimposed if `fit = True`. This is useful when all the structures have to be superimposed with the initial alignment.

4.9.22 `alignment.to_profile()` — convert alignment to profile format

Options:

<code>clean_sequences = <bool : 1></code>	<code>True</code>	whether or not clean non-standard residues
---	-------------------	--

Description: This command will convert the alignment, currently in memory, into the profile format. For more details on the profile format, see `profile.read()`.

If `clean_sequences` is set to `'on'`, then the non-standard residues in the sequences will be cleaned before transferring into the profile format. Specifically, ASX (B) will be replaced with ASN (N), GLX (Z) will be replaced with GLN (Q) and UNK (X) will be replaced with ALA (A).

Example: `examples/commands/aln_to_prof.py`

```
env = environ()

# Read in the alignment file
aln = alignment(env)
aln.append(file='toxin.ali', alignment_format='PIR', align_codes='ALL')

# Convert the alignment to profile format
prf = aln.to_profile(clean_sequences=True)

# Write out the profile

# in text file
prf.write(file='alntoprof.prf', profile_format='TEXT')

# in binary format
prf.write(file='alntoprof.bin', profile_format='BINARY')
```

4.9.23 `alignment.expand()` — put all models into alignment

Options:

<code>root_name = <str : 1></code>	<code>'undf'</code>	root of a filename for filename construction
<code>file_id = <str : 1></code>	<code>'default'</code>	file id for filename construction
<code>expand_control = <int : 5></code>	<code>9999 9999 1 10 0</code>	for controlling EXPAND_ALIGNMENT
<code>file_ext = <str : 1></code>	<code>''</code>	file extension for filename construction

Description: `file_id`, `root_name`, `expand_control`, `file_ext`, and `file_id` are used to construct atom filenames for all the models (see `modfile.default()`). The elements of `expand_control` are, in order, the first and last values of `id1`, the first and last values of `id2`, and the number of sequences to remove from the end of the current alignment before adding the new sequences. All the models covered by `expand_control` are added to the alignment, using the last sequence in the original alignment as the guide. This allows easy multiple superposition of all the templates and models after comparative modeling.

Example: `examples/commands/expand_alignment.py`

```
# Example for: alignment.expand()

# This will add models to the alignment.

env = environ()
aln = alignment(env, file='toxin.ali', align_codes=('2ctx', '2abx'))
aln.expand(expand_control=(9999, 9999, 1, 3, 0),
           root_name='2abx', file_id='.B', file_ext = '')
aln.write(file='toxin-expand.ali')
```

4.9.24 `alignment.segment_matching()` — align segments

Options:

<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>align_block = <int : 1></code>	0	the last sequence in the first block of sequences
<code>segment_report = <int : 1></code>	1000000	for SEGMENT_MATCHING
<code>segment_cutoff = <float : 1></code>	999999	cutoff for writing out an alignment in SEGMENT_MATCHING
<code>segment_shifts = <int : 0></code>	<code>segment shifts +- in SEG-</code>	
<code>segment_growth_n = <int : 0></code>	<code>MENT_MATCHING</code>	reducing/growing segment N-termini in SEGMENT_MATCHING
<code>segment_growth_c = <int : 0></code>	<code>reducing/growing segment C-termini in</code>	
<code>min_loop_length = <int : 0></code>	<code>SEGMENT_MATCHING</code>	inter-segment minimal lengths in SEGMENT_MATCHING
<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
<code>output_directory = <str : 1></code>	<code>''</code>	output directory
<code>root_name = <str : 1></code>	<code>'undef'</code>	root of a filename for filename construction
<code>file_id = <str : 1></code>	<code>'default'</code>	file id for filename construction
<code>file_ext = <str : 1></code>	<code>''</code>	file extension for filename construction

Description: This command enumerates alignments between two blocks of sequences. More precisely, it enumerates the alignments between the segments in the first block and the sequences in the second block. The segments can be moved to the left and right as well as lengthened and shortened, relative to the initial alignment. The regions not in segments or not aligned with segments are left un-aligned, possibly to be modeled as insertions. Typically, the first block of sequences corresponds to structures, the segments to secondary structure elements in the first block, and the second block to the sequences one of which is to be modeled later on. The command is useful for generating many alignments which can then be used by another MODELLER script to generate and evaluate the corresponding 3D models.

All the sequences and segments are defined in the alignment array. The first block of sequences, the ones with segments, are the first `align_block` sequences. The regions corresponding to the segments are defined by the last entry in the alignment as contiguous blocks of non-gap residues. Any standard single character residue code may be used. The segments must be separated by gap residues, '-'. The remaining sequences from `align_block + 1` to `NSEQ - 1` are the second block of sequences. The alignment of the sequences within the two blocks does not change. A sample alignment file is

The enumeration of alignments explores all possible combinations of alignments between each segment and the 2nd block of sequences: The starting position of each segment i is varied relative to the input alignment in the interval from `segment_shift[2i - 1]` to `segment_shift[2i]`. There has to be at least `min_loop_length[i]` and `min_loop_length[i + 1]` residues that are not in any segment before and after the i -th segment, respectively. The location of the N-terminus of segment i is varied relative to the location in the input alignment in the interval from `segment_growth_n[2i - 1]` to `segment_growth_n[2i]`. Similarly, the location of the C-terminus of segment i is varied relative to the location in the input alignment in the interval from `segment_growth_c[2i - 1]` to `segment_growth_c[2i]`. The shortening and lengthening of the segments may be useful in determining the best anchor regions for modeling of a loop.

Each alignment is scored according to the similarity scoring matrix specified by filename `rr_file`. This matrix may contain residue—gap scores, the gap being residue type 21; otherwise the value is set to the smallest value in the matrix. The score for an alignment is obtained by summing scores only over all alignment positions corresponding to the segments (no gap penalty is added for loops). When there is more than one sequence in any of the two blocks, the position score is an average of all pairwise comparisons between the two blocks of sequences. In the case where the number of positions in the alignment changes (*i.e.*, the segments grow or shorten), the scores are not comparable to each other. It is feasible to enumerate on the order of 10^{10} different alignments in less than one hour of CPU time.

In general, two runs are required. In the first run, the alignments are scored and a histogram of the scores is written to file `file`. Then this file must be inspected to determine the cutoff `segment_cutoff`. In the second run, all the alignments with a score higher than `segment_cutoff` are written to files in the PIR format, using the standard filenames convention: `output_directory/root_namefile_idnnnn0000file_ext`, where `nnnn` is the alignment file counter. In addition, the alignments are also written out in the PAP format for easier inspection by eye. Thus, `segment_cutoff` has to be set to a very large value in the first run, to avoid writing alignment files. During a run, a message is written to the log every `segment_report` alignments; this is useful for knowing what is going on during very long runs.

Example: `examples/commands/segment_matching.py`

```
# Example for: alignment.segment_matching()

log.level(1, 1, 1, 1, 0)
env = environ()

aln = alignment(env, file='ednf2.pap', align_codes=('7rsa', 'edn', 'templ'),
                alignment_format='PAP', remove_gaps=True)

aln.segment_matching(output_directory='./', file='segmatch.dat',
                    align_block=1, rr_file='${LIB}/as1.sim.mat',
                    segment_shifts=(-8, 8, 0, 0),
                    segment_growth_n=(0, 0, 0, 0),
                    segment_growth_c=(0, 0, 0, 0),
                    min_loop_length=(0,2,0),
                    segment_report=100000, segment_cutoff=0,
                    root_name='segmatch', file_ext='.ali', file_id='default')
```

4.10 The profile class: using sequence profiles

The `profile` class holds a sequence profile. Sequence profiles are similar to multiple alignments, and can contain gaps, but do not contain all of the information of the `alignment` class. Profiles can be matched against each other with `profile.scan()` or enriched from a sequence database with `profile.build()`.

4.10.1 `profile()` — create a new profile

Options:

`env = <environ>` MODELLER environment

Description: This creates a new `profile` object. By default, the profile is empty. However, if you give any arguments to the `profile()` constructor, they are passed to `profile.read()`, to read in an initial profile. See the `profile.scan()` example.

4.10.2 `profile.read()` — read a profile of a sequence

Options:

`file = <str : 1>` 'default' partial or complete filename
`profile_format = <str : 1>` 'TEXT' '| 'BINARY' ; for
READ/WRITE_PROFILE

Description: This command will read a profile from a specified file. Two formats are supported: `TEXT` and `BINARY`.

The format of the profile file (text) is as follows:

```
# Number of sequences:      4
# Length of profile       : 20
# N_PROF_ITERATIONS      : 3
# GAP_PENALTIES_1D       : -900.0  -50.0
# MATRIX_OFFSET          : 0.0
# RR_FILE                 : ${MODINSTALL8v1}/modlib//as1.sim.mat
1 2ctx                    X   0   71   1   71   0   0   0   0.0  0.0  IRCFITPDITS---KDCPN-
2 2abx                    X   0   74   1   74   0   0   0   0.0  0.0  IVCHTTATIPS-SAVTCPPG
3 1nbt                    X   0   66   1   66   0   0   0   0.0  0.0  RTCLISPSS---TPQTCPNG
4 1fas                    X   0   61   1   61   0   0   0   0.0  0.0  TMCYSHITTSRAILTNCG--
```

The first six lines begin with a '#' in the first column and give a few general details of the profile.

The first line gives the number of sequences in the profile. The line should be in the following format: '(24x,i6)'.

The second line gives the number of positions in the profile. This should be in '(24x,i6)' format also.

The third line gives the value of the `n_prof_iterations` variable. The fourth line gives the value of the `gap_penalties_1d` variable. The fifth line gives the value of the `matrix_offset` variable. The sixth line gives the value of the `rr_file` variable.

The number of sequences in the profile and its length are used to allocate memory for the profile arrays. So they should provide an accurate description of the profile.

The values of the variables described in lines 3 through 6 are not used internally by MODELLER. But the command expects to find a total of six header lines. These records represent useful information when `profile.build()` was used to construct the profile.

The remaining lines consist of the alignment of the sequences in the profile. The format of these lines is of the form: '(i5,1x,a40,1x,a1,1x,7(i5,1x),f5.0,1x,g10.2,1x,32767a1)'

The various columns that precede the sequence are:

1. The index number of the sequence in the profile.
2. The code of the sequence (similar to `alignment[x].code`).

3. The type of sequence ('S' for sequence, 'X' for structure). This depends on the original source of the sequences. (See `alignment.to_profile()` and `sequence_db.read()`).
4. The iteration in which the sequence was selected as significant. (See `profile.build()`).
5. The length of the database sequence.
6. The starting position of the target sequence in the alignment.
7. The ending position of the target sequence in the alignment.
8. The starting position of the database sequence in the alignment.
9. The ending position of the database sequence in the alignment.
10. The number of equivalent positions in the alignment.
11. The sequence identity of between the target sequence and the database sequence.
12. The e-value of the alignment. (See `profile.build()`).
13. The sequence alignment.

Many of the fields described above are valid only when the profile that is written out is the result of `profile.build()`.

Example: `examples/commands/read_profile.py`

```
# Example file for: profile.read(), profile.to_alignment()

env = environ()

# Create a new, blank, profile
prf = profile(env)

# Read in the profile file
prf.read(file='toxin.prf', profile_format='TEXT')

# Convert the profile to alignment
aln = prf.to_alignment()

# Write out the alignment
aln.write(file='readprofile.pir', alignment_format='PIR')
```

4.10.3 `profile.write()` — write a profile

Options:

<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
<code>gap_penalties_ld = <float : 2></code>	<code>900 50</code>	gap creation and extension penalties for sequence/sequence alignment
<code>matrix_offset = <float : 1></code>	<code>0.00</code>	substitution matrix offset for local alignment
<code>profile_format = <str : 1></code>	<code>'TEXT'</code>	<code>'TEXT' 'BINARY' ;</code> for <code>READ/WRITE_PROFILE</code>
<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file

Description: This command will write a profile to a specified file, together with a number of variables that are associated with the profile in the memory. Two formats are supported: `TEXT` and `BINARY`.

Example: `examples/commands/aln_to_prof.py`

```

env = environ()

# Read in the alignment file
aln = alignment(env)
aln.append(file='toxin.ali', alignment_format='PIR', align_codes='ALL')

# Convert the alignment to profile format
prf = aln.to_profile(clean_sequences=True)

# Write out the profile

# in text file
prf.write(file='alntoprof.prf', profile_format='TEXT')

# in binary format
prf.write(file='alntoprof.bin', profile_format='BINARY')

```

4.10.4 profile.to_alignment() — profile to alignment

Output: alignment

Description: This command will convert a profile that is in memory into the alignment format (see Section 4.9.1). The function of this command is complimentary to **alignment.to_profile()**. The generated alignment is returned.

Note: Not all information of a 'PIR' format is encoded in a profile. (See **profile.read()**). So converting a profile to an alignment may need manual attention to ensure that the alignment is useful for other routines.

Example: examples/commands/read_profile.py

```

# Example file for: profile.read(), profile.to_alignment()

env = environ()

# Create a new, blank, profile
prf = profile(env)

# Read in the profile file
prf.read(file='toxin.prf', profile_format='TEXT')

# Convert the profile to alignment
aln = prf.to_alignment()

# Write out the alignment
aln.write(file='readprofile.pir', alignment_format='PIR')

```

4.10.5 profile.scan() — Compare a target profile against a database of profiles

Options:

```

profile_list_file = <str : 1>          ''          list of profiles for PROFILE.PROFILE_SCAN

```

<code>profile_format = <str : 1></code>	<code>'TEXT'</code>	<code>'TEXT' 'BINARY'</code> ; for READ/WRITE_PROFILE
<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>matrix_offset = <float : 1></code>	<code>0.00</code>	substitution matrix offset for local alignment
<code>gap_penalties_1d = <float : 2></code>	<code>900 50</code>	gap creation and extension penalties for sequence/sequence alignment
<code>max_aln_evalue = <float : 1></code>	<code>0.1</code>	Max. E-value of alignments to include in BUILD_PROFILE
<code>aln_base_filename = <str : 1></code>	<code>'alignment'</code>	basename for construction of alignment filenames used by PROFILE_PROFILE_SCAN
<code>score_statistics = <bool : 1></code>	<code>True</code>	PROFILE_PROFILE_SCAN: if turned off, the length-normalized z-scores are not computed
<code>output_scores = <bool : 1></code>	<code>False</code>	whether to output individual scores in a build_profile scan
<code>output_score_file = <str : 1></code>	<code>'default'</code>	output file for writing out individual scores in seqfilter
<code>pssm_weights_type = <str : 1></code>	<code>'HH1'</code>	type of weighting to calculate pssm; <code>'HH0' 'HH1' 'PSIC'</code>
<code>write_summary = <bool : 1></code>	<code>True</code>	whether to write summary information for PPSCAN
<code>summary_file = <str : 1></code>	<code>'ppscan.sum'</code>	output file for writing PPSCAN summary
<code>output_alignments = <bool : 1></code>	<code>True</code>	PROFILE_PROFILE_SCAN: if turned off, no alignments will be written out.

Description: This command scans the given target profile against a database of template profiles and reports significant alignments; the target profile should have been read previously with the `profile.read()` command.

All the profiles listed in `profile_list_file` should be in a format that is understood by `profile.read()`.

The `profile_list_file` should contain absolute or relative paths to the individual template profiles, one per line.

See documentation under `profile.read()` for help on `profile_format`.

`rr_file` is the residue-residue substitution matrix to use when calculating the position-specific scoring matrix (PSSM). The current implementation is optimized only for the BLOSUM62 matrix.

`gap_penalties_1d` are the gap penalties to use for the dynamic programming. `matrix_offset` is the value to be used to offset the substitution matrix. The most optimal values for these parameters are: `matrix_offset = -200`
`gap_penalties_1d = -1900 -95`

`max_aln_evalue` sets the threshold for the E-values. Alignments with e-values better than the threshold will be written out.

`aln_base_filename` sets the base filename for the alignments. The output alignment filenames will be of the form `ALN_BASE_FILENAME_XXXX.ali`. The `XXXX` is a 4-digit integer (prefixed with sufficient zeroes) that is incremented for each alignment. For example, `alignment_0001.ali`

`score_statistics` is a flag that triggers the calculation of e-values. If set to OFF, the significance estimates for the alignments will not be calculated. The calculation of alignment significance is similar to that used for `profile.build()`. This option can be useful when there are only a very small number of template profiles in `profile_list_file`, insufficient to calculate reliable statistics. Also see `profile.build()`.

`output_scores` is a flag to write out the raw alignment scores, zscores and e-values for all the comparisons. `output_score_file` sets the name of the file to which this output should be written to.

`write_summary` is a flag to output a summary of all the significant alignments into the file specified by `summary_file`.

If `output_alignments` is set to OFF, alignments will not be written out.

Example: `examples/commands/ppscan.py`

```

# Example for: profile.scan()

env = environ()

# Read in the target profile
prf = profile(env, file='T3lzt-uniprot90.prf', profile_format='TEXT')

# Scan against all profiles in the 'profiles.list' file
prf.scan(profile_list_file = 'profiles.list',
         matrix_offset     = -200,
         rr_file           = '${LIB}/blosum62.sim.mat',
         gap_penalties_1d = (-1900, -95),
         score_statistics  = False,
         output_alignments = True,
         output_scores     = False,
         output_score_file = 'T3lzt-ppscan.scores',
         profile_format    = 'TEXT',
         max_aln_evalue    = 1,
         aln_base_filename = 'T3lzt-ppscan',
         pssm_weights_type = 'HH1',
         write_summary     = True,
         summary_file      = 'T3lzt-ppscan.sum')

```

4.10.6 profile.build() — Build a profile for a given sequence or alignment**Options:**

<code>sdb = <sequence_db></code>		Input database of sequences
<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>gap_penalties_1d = <float : 2></code>	900 50	gap creation and extension penalties for sequence/sequence alignment
<code>matrix_offset = <float : 1></code>	0.00	substitution matrix offset for local alignment
<code>n_prof_iterations = <int : 1></code>	3	number of iterations in PROFILE_SEARCH
<code>check_profile = <bool : 1></code>	True	whether to monitor profile degeneration
<code>output_scores = <bool : 1></code>	False	whether to output individual scores in a build_profile scan
<code>output_score_file = <str : 1></code>	<code>'default'</code>	output file for writing out individual scores in seqfilter
<code>max_aln_evalue = <float : 1></code>	0.1	Max. E-value of alignments to include in BUILD_PROFILE
<code>gaps_in_target = <bool : 1></code>	False	whether to include gaps in target when using build_profile

Description: This command iteratively scans a database of sequences to build a profile for the input sequence or alignment. The command calculates the score for a Smith-Waterman local alignment between the input sequence and each of the sequences in the database. The significance of the alignment scores (e-values) are calculated using a procedure similar to that described by Pearson (1998).

Alignments with e-values below `max_aln_evalue` are then added to the current alignment. A position-specific scoring matrix is then calculated for the current alignment and is used to search the sequence database. This

procedure is repeated for `n_prof_iterations` or until there are no significant alignments below the threshold, whichever occurs first.

The initial sequence or alignment can be read in either in the profile format, with `profile.read()`, or as an alignment using `alignment.append()`. In the latter case, the alignment has to be converted to the profile format using `alignment.to_profile()`.

The output contains a multiple sequence alignment (assembled) of all the homologues of the input sequence found in the database. The output can be formatted as a profile with `profile.write()` or converted into any of the standard alignment formats using `profile.to_alignment()`. It can then be written out to a file with `alignment.write()`.

The fit between the observed and theoretical distributions of the z-scores is calculated after each iteration and is reported in the log file. The fit is calculated using the Kolmogorov-Smirnov D-statistic. If the `check_profile` flag is set to 'on', then the command will not proceed if the fit deviates by more than 0.04 (D-statistic).

By default, regions of the alignment that introduce gaps in the target sequence are ignored (deleted) in the final multiple alignment. But if `gaps.in_target` is set to 'on', then the gaps are retained. (See below for comments).

If the `output_scores` flag is set to 'on', then the scores of each alignment between the input sequence and each database sequence, from all iterations, will be written out to the file specified in `output_score_file`.

Comments:

1. The procedure has been optimized only for the BLOSUM62 similarity matrix.
2. The dynamic programming algorithm has been optimized for performance on Intel Itanium2 architecture. Nevertheless, the calculation is sufficiently CPU intensive. It takes about 20 min for an iteration, using an input sequence of 250aa against a database containing 500,000 sequences on an Itanium2 machine. It could take much longer on any other machine.
3. It is advisable to have `gaps.in_target` set to 'off', when scanning against large databases to avoid the local-alignments inserting a huge number of gaps in the final alignments.
4. The statistics will not be accurate (or may even fail) if the database does not have sequences that represent the entire range of lengths possible.
5. The method can be used for fold-assignment by first building a profile for the target sequence by scanning against a large non-redundant sequence database (like swissprot) and then using the resulting profile to scan once against a database of sequences extracted from PDB structures. `gaps.in_target` can be set to 'on' in the second step to get the complete alignments that can then be used for modeling.

Example: `examples/commands/build_profile.py`

```
log.verbose()
env = environ()

#-- Prepare the input files

#-- Read in the sequence database
sdb = sequence_db(env)
sdb.read(seq_database_file='pdb95.fsa', seq_database_format='FASTA',
         chains_list='ALL', minmax_db_seq_len=(1, 40000), clean_sequences=True)

#-- Write the sequence database in binary form
sdb.write(seq_database_file='pdb95.bin', seq_database_format='BINARY',
         chains_list='ALL')

#-- Now, read in the binary database
sdb.read(seq_database_file='pdb95.bin', seq_database_format='BINARY',
         chains_list='ALL')

#-- Read in the target sequence/alignment
```

```
aln = alignment(env)
aln.append(file='toxin.ali', alignment_format='PIR', align_codes='ALL')

#-- Convert the input sequence/alignment into
#   profile format
prf = aln.to_profile()

#-- Scan sequence database to pick up homologous sequences
prf.build(sdb, matrix_offset=-450, rr_file='${LIB}/blosum62.sim.mat',
          gap_penalties_1d=(-500, -50), n_prof_iterations=5,
          check_profile=False, max_aln_evalue=0.01, gaps_in_target=False)

#-- Write out the profile
prf.write(file='buildprofile.prf')

#-- Convert the profile back to alignment format
aln = prf.to_alignment()

#-- Write out the alignment file
aln.write(file='buildprofile.ali', alignment_format='PIR')
```

4.11 The sequence_db class: using sequence databases

The `sequence_db` class holds a database of sequences. Such a database is similar to a multiple sequence alignment, but contains less auxiliary information (for example, no sequence may contain gaps). This requires less memory than a true alignment, and is thus more suited for large databases of sequences. Such a database can be scanned for matches to an input sequence with `sequence_db.search()`, used to build sequence profiles with `profile.build()` or filtered by given criteria with `sequence_db.filter()`.

4.11.1 sequence_db() — create a new sequence database

Options:

`env = <environ>` MODELLER environment

Description: This creates a new, empty, sequence database. If you give any arguments to the `sequence_db()` constructor, they are passed to `sequence_db.read()`, to read in an initial database. See the `sequence_db.filter()` example.

4.11.2 sequence_db.read() — read a database of sequences

Options:

<code>chains_list = <str : 1></code>	<code>'\$(LIB)/CHAINS_3.0_40_XN.cod'</code>	file with sequences
<code>seq_database_file = <str : 1></code>	<code>'\$(LIB)/CHAINS_all.seq'</code>	file with a list of sequence codes
<code>seq_database_format = <str : 1></code>	<code>'PIR'</code>	<code>'PIR'</code> <code>'FASTA'</code> <code>'BINARY'</code> ; for READ/WRITE_SEQUENCE_DB
<code>clean_sequences = <bool : 1></code>	<code>True</code>	whether or not clean non-standard residues
<code>minmax_db_seq_len = <int : 2></code>	<code>0 999999</code>	minimal/maximal database sequence length

Description: This command will read a database of sequences, either in the PIR, FASTA, or BINARY format.

If the format is PIR or FASTA:

- It is possible to clean all sequences of non-standard residue types by setting `clean_sequences` to `True`.
- Sequences shorter than `minmax_db_seq_len[1]` and longer than `minmax_db_seq_len[2]` are eliminated.
- Only sequences whose codes are listed in the `chains_list` file are read from the `seq_database_file` of sequences. If `chains_list` is `all`, all sequences in the `seq_database_file` file are read in, and there is no need for the `chains_list` file.

For the PIR and FASTA formats, make sure the order of sequences in the `chains_list` and `seq_database_file` is the same for faster access (there can of course be more sequences in the sequence file than there are sequence codes in the codes file).

Additionally, if the sequences are in 'PIR' format, then the protein type and resolution fields are stored in the database format. (see Section 4.9.1 for description of 'PIR' fields).

The protein type field is encoded in a single letter format. 'S' for sequence and 'X' for structures of any kind. This information is transferred to the profile arrays when using `profile.build()`. (See also `profile.read()`).

The resolution field is used to pick representatives from the clusters in `sequence_db.filter()`.

None of the options above apply to the BINARY format, which, in return, is very fast (*i.e.*, 3 seconds for 300 MB of 800,000 sequences in the TrEMBL database).

Example: See `profile.build()` command.

4.11.3 `sequence_db.write()` — write a database of sequences

Options:

<code>chains_list = <str : 1></code>	<code>'\$(LIB)/CHAINS_3.0_40_XN.cod'</code>	file with sequences
<code>seq_database_file = <str : 1></code>	<code>'\$(LIB)/CHAINS_all.seq'</code>	file with a list of sequence codes
<code>seq_database_format = <str : 1></code>	<code>'PIR'</code>	<code>'PIR'</code> <code>'FASTA'</code> <code>'BINARY'</code> ; for READ/WRITE_SEQUENCE_DB

Description: This command will write a database of sequences currently in memory, either in the PIR, FASTA, or BINARY format. The `chains_list` file is written only for the PIR or FASTA formats.

Example: See `profile.build()` command.

4.11.4 `sequence_db.search()` — search for similar sequences

Options:

<code>aln = <alignment></code>		Alignment containing the target sequence
<code>io = <io_data></code>		Options for reading atom files
<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>gap_penalties_ld = <float : 2></code>	900 50	gap creation and extension penalties for sequence/sequence alignment
<code>off_diagonal = <int : 1></code>	100	to speed up the alignment
<code>matrix_offset = <float : 1></code>	0.00	substitution matrix offset for local alignment
<code>overhang = <int : 1></code>	0	un-penalized overhangs in protein comparisons
<code>local_alignment = <bool : 1></code>	False	whether to do local as opposed to global alignment
<code>seq_database_file = <str : 1></code>	<code>'\$(LIB)/CHAINS_all.seq'</code>	file with a list of sequence codes
<code>search_group_list = <str : 1></code>	<code>'\$(LIB)/CHAINS_3.0_40_XN.grp'</code>	file with 40% groups of sequences
<code>alignment_features = <str : 1></code>	<code>'INDICES CONSERVATION'</code>	what alignment features to write out: 'ACCURACY' 'HELIX' 'BETA' 'ACCESSIBILITY' 'STRAIGHTNESS' 'CONSERVATION' 'INDICES' 'ALL' 'GAPS'
<code>search_top_list = <int : 1></code>	20	the length of the output hits list
<code>output = <str : 1></code>	<code>'LONG'</code>	<code>'SHORT'</code> <code>'LONG'</code>
<code>search_sort = <str : 1></code>	<code>'LONGER'</code>	which sequence to use for normalization when sorting the hit list: <code>'SHORTER'</code> <code>'LONGER'</code>
<code>search_randomizations = <int : 1></code>	0	number of randomizations for calculating the significance of a sequence/sequence similarity
<code>fast_search = <bool : 1></code>	False	whether to use fast sequence search or not
<code>fast_search_cutoff = <float : 1></code>	1.0	if FAST_SEARCH is ON only sequences with database scan significance higher than this value are considered for randomization significance
<code>data_file = <bool : 1></code>	False	whether results go to a separate file or not
<code>signif_cutoff = <float : 2></code>	4.0 5.0	cutoff for adding sequences to alignment, max difference from the best

Description: This command searches a sequence database for proteins that are similar to a given target sequence.

The target sequence should be the only sequence in the provided alignment, `aln`.

The database of sequences to be scanned against must be read previously by the `sequence_db.read()` command.

The command uses the dynamic programming method for the best sequence alignment, given the gap creation and extension penalties specified by `gap_penalties_1d` and residue type scores read from file `rr_file`. `gap_penalties_1d[1]` is a gap creation penalty and `gap_penalties_1d[2]` is a gap extension penalty.

The `search_top_list` top hits are written to the `log` file at the end. The hits are sorted according to the fractional sequence identity score obtained by dividing the number of identical residue pairs by the length of the longer sequence (`search_sort = 'LONGER'`) or the shorter sequence (`search_sort = 'SHORTER'`).

The final list of hits contains three different significance values:

1. **SIGNI.** Z-score from sequence randomizations. This is the most accurate significance score, but the slowest one to calculate. For each pairwise comparison, the two sequences are shuffled a specified number of times (`search_randomizations`) to obtain the mean and standard deviation of “random” scores from which the Z-score for an alignment score of a given pair of sequences is calculated.
2. **SIGNI2.** Z-score for sequence identity from the database scan. After comparison of the target sequence with all sequences in the database is done, the comparisons are sorted by the length of the database sequence. The pairwise sequence identities of the 20 sequences closest in length to the target sequence are used to calculate the average and standard deviation of the percentage sequence identities for subsequent calculation of the Z-score for the percentage sequence identity of a given pairwise alignment.
3. **SIGNI3.** Z-score for alignment score from the database scan. The procedure is the same as for **SIGNI2**, except that the alignment scores are used instead of the pairwise sequence identities.

The calculation of the Z-scores assumes that the random scores are distributed according to the Gaussian distribution, instead of the extreme value distribution [Karlin & Altschul, 1990], which is more correct.

`search_randomizations` specifies how many alignments of the shuffled sequences are done to calculate the significance score for the overall sequence similarity. If 0, the significance is not calculated. If more than 5 randomizations are done, the significance score, not sequence identity, is used for sorting the hit list.

When `fast_search` is `True` only those sequences that have a database-scan alignment score significance (**SIGNI3** in output) above `fast_search_cutoff` are used for the “full” randomization-based significance calculation. Since the mean and the standard deviation of the distribution obtained by randomizing the two compared sequences are much more appropriate than the corresponding quantities for the target/database comparisons, `fast_search` should be `True` only when you are in a hurry and the database is large.

If `data_file` is `True` the final results (list of PDB codes with significances, *etc.*) are also written to a separate file `'seqsearch.dat'`.

If `output` is `'LONG'`, the best alignment for each sequence in the database and its various scores are also written to the `log` file. If `output` is `'VERY_LONG'`, individual scores obtained for randomized sequences are also written to the `log` file (this is almost never needed).

If the selected significance score is larger than `signif_cutoff[1]` and not more than `signif_cutoff[2]` units worse than the best hit, all the members of the same group, as defined in `search_group_list`, are added to the alignment array. These sequences are taken from `seq_database_file`, which is often (but not always) the same file previously provided to `sequence_db.read()`. Subsequent `alignment.malign()`, `environ.dendrogram()` and `alignment.write()` can then be used to write out all related PDB chains aligned to the target sequence.

Example: `examples/commands/sequence_search.py`

```
# Example for: sequence_db.search()

# This will search the MODELLER database of representative protein chains
# for chains similar to the specified sequence.
```

```

log.verbose()
env = environ()

# Read in the standard database included in the Modeller distribution
sdb = sequence_db(env, seq_database_file='${LIB}/CHAINS_all.seq',
                  chains_list='very-short-for-test.cod')

# Read in the query sequence in alignment format
aln = alignment(env, file='toxin.ali', align_codes='1nbt')

sdb.search(aln, search_randomizations=20, # should use 100 in real life
           seq_database_file='${LIB}/CHAINS_all.seq',
           off_diagonal=9999, gap_penalties_1d=(-800, -400),
           signif_cutoff=(1.5, 5.0))

aln.malign()
aln.write(file='toxin-search.pap', alignment_format='PAP')

```

4.11.5 `sequence_db.filter()` — cluster sequences by sequence-identity

Options:

<code>rr_file = <str : 1></code>	<code>'\$(LIB)/as1.sim.mat'</code>	input residue-residue scoring file
<code>gap_penalties_1d = <float : 2></code>	900 50	gap creation and extension penalties for sequence/sequence alignment
<code>matrix_offset = <float : 1></code>	0.00	substitution matrix offset for local alignment
<code>output_grp_file = <str : 1></code>	<code>'seqfilt.grp'</code>	output file for seqfilter groups
<code>output_cod_file = <str : 1></code>	<code>'seqfilt.cod'</code>	output file for seqfilter representative groups
<code>seqid_cut = <int : 1></code>	95	Sequence Identity cut-off for SEQFILTER
<code>max_diff_res = <int : 1></code>	30	Length cut-off for SEQFILTER

Description: This command clusters a set of sequences by sequence identity. The command uses a greedy algorithm: the first sequence in the file becomes the first group representative. All other sequences are compared with this and if they are similar enough, as specified in `seqid_cut`, they are added as members of this group. These sequences are not used for further comparisons. The next non-member sequence becomes the next group representative and so on.

The initial set of sequences must be read previously by the `sequence_db.read()` command with `seq_database_format` being either 'PIR' or 'FASTA'.

`rr_file` is residue-residue substitution matrix. The command only handles similarity matrices for efficiency purposes.

The command uses the Smith-Waterman dynamic programming method for the best sequence alignment, given the gap creation and extension penalties specified by `gap_penalties_1d` and residue type scores read from file `rr_file`. `gap_penalties_1d[1]` is a gap creation penalty and `gap_penalties_1d[2]` is a gap extension penalty. The command only works with similarity matrices for efficiency reasons.

The final list of groups and their members is written out to `output_grp_file`. The codes of the representative sequences is written out to `output_cod_file`.

The clustering algorithm evaluates the following conditions in hierarchical order before adding a sequence to a group:

1. The difference in length: If the difference in the number of residues between the group representative and the sequence being compared is greater than `max_diff_res`, the sequence will not be included into that group.
2. The number of unaligned residues: After the local alignment is performed, a sequence will not be considered for addition into a group unless the difference between the smaller of the two sequences and the number of aligned positions in the alignment is less than `max_unaligned_res`.
3. Sequence Identity: Finally, if the sequence identity calculated from the alignment is greater than `seqid_cut`, the sequence is added to a group.

If the initial set of sequences read were in 'PIR' format with values in the resolution field, then the group representative is the sequence with the highest resolution. This is especially useful when clustering sequences from the PDB.

Example: `examples/commands/seqfilter.py`

```
log.verbose()
env = environ()

sdb = sequence_db(env, seq_database_file='sequences.pir',
                  seq_database_format='PIR',
                  chains_list='ALL', minmax_db_seq_len=[30, 3000],
                  clean_sequences=True)

sdb.filter(rr_file='${LIB}/id.sim.mat', gap_penalties_1d=[-3000, -1000],
           max_diff_res=30, seqid_cut=95, output_grp_file='seqfilt.grp',
           output_cod_file='seqfilt.cod')
```

4.12 The density class: handling electron microscopy density data

The `density` class stores all information from an electron microscopy density map file. Protein models can then be docked to this density (using `density.grid_search()`) to improve their quality.

4.12.1 `density()` — create a new density map

Options:

`env = <environ>` MODELLER environment

Description: This creates a new, empty, density map. If you give any arguments to the `density()` constructor, they are passed to `density.read()`, to read in an initial density. See the `density.grid_search()` example.

4.12.2 `density.read()` — read an EM (electron microscopy) density map file

Options:

<code>file = <str : 1></code>	<code>'default'</code>	input electron density map file
<code>em_density_format = <str : 1></code>	<code>'XPLOR'</code>	input density map file format for EM_GRID_SEARCH; <code>'MRC'</code> <code>'XPLOR'</code>
<code>density_type = <str : 1></code>	<code>'SPHERE'</code>	Function used to calculate density map cross-correlation in EM_GRID_SEARCH; <code>'SPHERE'</code> <code>'GAUSS'</code> <code>'HYBRID'</code> <code>'GAUSS_NORM'</code> <code>'TRACE'</code>
<code>voxel_size = <float : 1></code>	<code>0</code>	EM density map voxel size for EM_GRID_SEARCH
<code>resolution = <float : 1></code>	<code>0</code>	EM density map resolution for EM_GRID_SEARCH
<code>em_map_size = <int : 1></code>	<code>0</code>	size of the density map, for EM_GRID_SEARCH

Description: This command reads a density map from `file`, which should be provided as a cubic grid of intensities, in the X-PLOR format [Brünger, 1992]. The size of the cubic grid is given by `em_map_size`. The density map resolution and the map voxel size are given by `resolution` and `voxel_size`, respectively.

When docking the probe into the EM grid, the probe structure is converted first into probe density, by using the function indicated in the `density_type` variable. Each atom can be represented by one of several atomic density functions, including, the uniform sphere model (`'SPHERE'`), the Gaussian function (`'GAUSS'`), a normalized Gaussian function (`'NORM'`), a hybrid Gaussian/sphere model (`'HYBRID'`), and an interpolation to the closest point on the grid (`'TRACE'`). The recommended function is `'SPHERE'`.

Example: See `density.grid_search()` command.

4.12.3 `density.grid_search()` — dock a structure into an EM (electron microscopy) density map

Options:

<code>em_density_format = <str : 1></code>	<code>'XPLOR'</code>	input density map file format for EM_GRID_SEARCH; <code>'MRC'</code> <code>'XPLOR'</code>
<code>num_structures = <int : 1></code>	<code>1</code>	number of structures to dock in EM_GRID_SEARCH
<code>em_pdb_name = <str : 0></code>	PDB files to read for	EM_GRID_SEARCH

<code>chains_num = <int : 0></code>	1	number of chains in input pdb filein EM_GRID_SEARCH
<code>start_type = <str : 1></code>	'CENTER'	how to start EM_GRID_SEARCH; 'CENTER' 'ENTIRE' 'SPECIFIC'
<code>translate_type = <str : 1></code>	'NONE'	how to perform translations dur- ing EM_GRID_SEARCH; 'NONE' 'RANDOM' 'EXHAUSTIVE'
<code>number_of_steps = <int : 1></code>	1	number of steps of Monte Carlo opti- mization
<code>angular_step_size = <float : 1></code>	0	Angular search step size in degrees for EM_GRID_SEARCH
<code>temperature = <float : 1></code>	293.0	temperature for Monte Carlo simulation
<code>best_docked_models = <int : 1></code>	1	number of best docked models to keep in EM_GRID_SEARCH
<code>em_fit_output_file = <str : 1></code>	''	output file for EM_GRID_SEARCH

Requirements: PDB files

Output: 'filename.log', 'filename_init_1.pdb', 'filename_init.xplor', 'filename_best.xplor',
'filename_1_num.pdb', 'bestCC_filename'

Description: This command docks a structure of a protein/domain (probe) into a given EM density map. See original paper for the description of the method and the most recommended protocols [Topf *et al.*, 2005].

The probe is specified by the variable `em_pdb_name`. Before starting the protocol, the probe is positioned on the EM density grid based on the `start_type` variable:

- 'CENTER' will translate the center of mass of the probe to the center of mass of the grid.
- 'ENTIRE' will divide the grid into cells similar in volume to the probe, and will translate the probe to the center of each of these cells consecutively.
- 'SPECIFIC' will not perform any translation and will just use the coordinates specified by the user (the input PDB coordinates) as a starting position.

The best fit between the probe and the EM density map is obtained by changing the position of the model so as to maximize the cross-correlation between the probe density and the EM density.

The optimization of the cross-correlation score is performed by a 6D search of the probe on the EM grid, based on the variable `translate_type`. If 'NONE' is specified, only a rotational search over the three rotational Euler angles (ϕ , θ , ψ) is performed, with no translations. `angular_step_size` gives the maximal step size of searching for all combinations of Euler angles, which is recommended to be 30 degrees. The best fit from this coarse search is refined by a finer local search in all three Euler angles. For a protein of 150 residues this calculation typically takes less than 0.5 minutes on a 3.0 GHz Intel Xeon processor.

If `translate_type` is 'RANDOM', a Monte Carlo (MC) optimization is performed, and the number of MC steps has to be specified (by `number_of_steps`). A single MC step consists of (i) a random translation of the probe for one voxel on the EM grid, (ii) a search for the three Euler angles that maximize cross-correlation score, and (iii) an application of the Metropolis criterion [Metropolis *et al.*, 1953]. The temperature used for the Metropolis criterion is specified in `temperature` (typically about 5000 units). This calculation typically takes about 1-2 minutes.

When `translate_type` is 'EXHAUSTIVE', a local search is performed with the probe on the grid. The optimal orientations at its original position and all 26 (*i.e.*, $3 \times 3 \times 3 - 1$) neighboring grid points are obtained successively by enumerating all three Euler angles. A Monte Carlo criterion is applied to each one of these 27 optimal orientations (usually with temperature lower than 5000). `number_of_steps` gives the number of steps for which this process is repeated (typically 25 times). When the EM density map covers only the probe model and `start_type` is 'CENTER', this protocol can be used for a translational and rotational refinement of the initial superposition of the centers of mass. When `start_type` is 'ENTIRE' this protocol will be applied for a local

search only in those cells where the randomly oriented probe gives a positive cross-correlation score. The calculation typically takes about 10-15 minutes.

`best_docked_models` specifies how many best-fitted models should be saved by the program. This becomes more important at decreasing resolutions, as the best solution will not necessarily have the highest cross-correlation score.

Example: `examples/commands/em_grid_search.py`

```
log.verbose()
env = environ()

den = density(env, file='1cuk-a2_8A.xplor', em_density_format='XPLOR',
              voxel_size=1., resolution=8., em_map_size=32,
              density_type='SPHERE')

den.grid_search(em_density_format='XPLOR', num_structures=1,
               em_pdb_name=['1cuk-a2.pdb'], chains_num=[1],
               start_type='CENTER', number_of_steps=1, angular_step_size=30.,
               temperature=0., best_docked_models=1,
               em_fit_output_file='test-cr.log')
```

4.13 The info object: obtaining information about the MODELLER build

The `info` object holds information about the current MODELLER build (most of this can also be found in the first few lines of the log file). This information is useful when reporting bugs, or for writing scripts which require a certain MODELLER version.

4.13.1 `info.version` — return the full MODELLER version number

Description: This returns the current version, as a string. This contains all the information returned by `info.version_info` and `info.build_date`, but in a less easily machine-readable form.

4.13.2 `info.version_info` — return the version number, as a tuple

Description: This returns only the version number, as a Python tuple containing the major and minor version numbers. For example, version 8v0 would return (8, 0). (CVS builds always return 'CVS' instead.)

4.13.3 `info.build_date` — return the date this binary was built

Description: This returns the date on which this MODELLER binary was built, as a string in 'MM/DD/YYYY HH:MM:SS' format.

4.13.4 `info.exe_type` — return the executable type of this binary

Description: This returns a string identifying the machine and FORTRAN compiler type of this MODELLER binary.

4.13.5 `info.debug` — returns this binary's debug flag

Description: This returns `True` if this binary was built with compiler debugging flags, or `False` if full code optimization was turned on. Debugging builds generally run slower than optimized builds, so release binaries will invariably return `False`.

4.13.6 `info.time_mark()` — print current date, time, and CPU time

Description: This prints to the log file the total CPU time used in this run (in seconds) and the CPU time used since the last time this command was called. The date and time on which the run was started, plus the current date and time, are also printed, in 'YYYY/MM/DD HH:MM:SS.SSS' format.

4.14 The log object: controlling the amount of output

The `log` object allows you to control the amount of information output to the MODELLER log file. (It is also used internally to divert the system standard output, *e.g.* from the Python `'print'` statement, to the log file.)

4.14.1 `log.level()` — Set all log output levels

Options:

<code>output = <int : 1></code>	Whether to display output logs
<code>notes = <int : 1></code>	Whether to display note logs
<code>warnings = <int : 1></code>	Whether to display warning logs
<code>errors = <int : 1></code>	Whether to display error logs
<code>memory = <int : 1></code>	Whether to display memory logs

Description: This sets all five of MODELLER's log outputs. (This is very similar in operation to the old `OUTPUT_CONTROL` variable.) Each argument should be an integer, either 0 to display no output, or 1 to display all output. An exception is the `memory` argument, which can also be set to 2 to display additional dynamic memory information.

4.14.2 `log.minimal()` — display minimal log output

Description: This instructs MODELLER to only display important outputs, and errors.

4.14.3 `log.verbose()` — display verbose log output

Description: This instructs MODELLER to display all log output.

4.14.4 `log.very_verbose()` — display verbose log output, and dynamic memory information

Description: This instructs MODELLER to display all log output. Additionally, a breakdown of all dynamic memory used by MODELLER is displayed every time memory is allocated or freed. (Note that some information is held statically, and this is not tracked. Some routines require additional temporary storage, which is also not listed here. Finally, memory used by any Python variables is not accounted for.)

4.15 The modfile module: handling of files

The `modfile` module contains routines for dealing with files which are used by MODELLER. More complete facilities for most of these functions are available within the standard `'os'` Python module; however, the `modfile` module is provided for compatibility with old TOP scripts, and for systems with incomplete Python installations.

4.15.1 `modfile.default()` — generate an ‘automatic’ filename

Options:

<code>root_name = <str : 1></code>	filename root
<code>file_id = <str : 1></code>	file type
<code>id1 = <int : 1></code>	first numeric ID
<code>id2 = <int : 1></code>	second numeric ID
<code>file_ext = <str : 1></code>	filename extension

Description: This returns a generated file name, as `(root_name)(file_id)(id1)(id2)(file_ext)`, where `id1` and `id2` are printed as 4-digit numbers, padded with zeroes if necessary. For example, `'2ptn.B99990001.pdb'` results from `root_name = '2ptn'`, `file_id = '.B'`, `id1 = 9999`, `id2 = 1`, and `file_ext = '.pdb'`. This mimics the ‘automatic’ filename generation of previous versions of MODELLER when filenames were set to `'default'` or `'${DEFAULT}'`.

4.15.2 `modfile.delete()` — delete a file

Options:

<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
-------------------------------------	------------------------	------------------------------

Description: This command deletes the named file.

4.15.3 `modfile.inquire()` — check if file exists

Options:

<code>file = <str : 1></code>	<code>'default'</code>	partial or complete filename
-------------------------------------	------------------------	------------------------------

Output: `file_exists`

Description: This command returns 1 if the specified file exists, or 0 otherwise.

Chapter 5

Converting TOP scripts from old MODELLER versions

Previous versions of MODELLER used TOP as their scripting language. TOP is a language similar in syntax to FORTRAN, which is also used by ASGL. For increased power and flexibility, interoperability with other programs, and improved ease of use, MODELLER now uses Python 2.3 for its control language.

5.1 Running old scripts unchanged

For compatibility with old codes, MODELLER will still run most TOP scripts. By default, the program expects to read new-style Python scripts, but if the script file extension ends in `'.top'` (as in previous versions), it will be assumed to be a TOP file and will be run as such. Note, however, that this behavior is deprecated and will probably be removed entirely in a later release. Also bear in mind that some newer commands and features will not be available in TOP; thus, it is recommended that you convert your TOP scripts to Python.

If you do wish to run old scripts unchanged, please note that the `GO_TO` function is no longer present in the TOP language. If you have code which makes use of this function, you should use the `EXIT` and `CYCLE` flow control statements instead, which either terminate a `DO` loop or skip to its next iteration. Single-line `IF` statements are no longer supported either; you should use the `ELSE` and `END_IF` statements to build multi-line `IF` clauses instead. See the scripts in the `bin` directory for examples.

5.2 Converting TOP scripts to Python

5.2.1 TOP commands and variables

TOP variables come in four varieties — strings, reals, integers, logicals — as either single values or lists. The equivalents in Python are the standard `'str'`, `'float'`, `'int'` and `'bool'` types, respectively, while lists can be represented as `'list'` or `'tuple'` objects. In essence, this means that strings must be quoted, logical `on` and `off` become `bool True` and `False` respectively, and lists must be comma-separated and enclosed in parentheses. Also, the `'='` operator is not mandatory in TOP, but it is in Python. See the Python documentation for more information. For example, consider the following assignments in TOP:

```
SET STRVAR = foo           # Set a string variable (quotes not required)
SET REALVAR 3.4           # Set a real variable (= not required either)
SET INTVAR = 4            # Set an integer variable
SET LOGVAR = on          # Set a logical variable
SET INTLIST = 1 1 3 5    # Set a list of integers
SET STRLIST = 'one' 'two' 'three' # Set a list of strings
```

The equivalent Python code would be:

```

STRVAR = 'foo'           # Set a string variable
REALVAR = 3.4           # Set a real variable
INTVAR = 4              # Set an integer variable
LOGVAR = True           # Set a logical variable
INTLIST = (1, 1, 3, 5)  # Set a list of integers
STRLIST = ('one', 'two', 'three') # Set a list of strings

```

Variables in TOP are case-insensitive; that is, the names `GAP_PENALTIES_1D` and `gap_penalties_1d` refer to the same variable. (Upper case is usually used, but this is just by convention.) Python variables are case-sensitive, so these two names refer to different variables. For consistency with other codes, all Python commands and variables used by MODELLER are **lower-case**.

All variables in TOP are global; that is, once they are set, their value is kept for the rest of the program (or until changed again). This is irrespective of whether the variable is set while calling a TOP command, or whether an explicit **SET** command is used. For example, this TOP script:

```
ALIGN GAP_PENALTIES_1D = 900 50
```

will behave identically to this code:

```
SET GAP_PENALTIES_1D = 900 50
ALIGN
```

In Python, on the other hand, each command takes a number of arguments. For example, the `align()` command takes a `gap_penalties_1d` argument. The value of this argument affects only this call to `align()`, *i.e.* it is a local variable. Note, however, that for convenience, if an argument is omitted, the default value (if available) is used. This default value is taken from the `environ` class, and so behaves similarly to the global defaults (in `'${LIB}/top.ini'`) of previous MODELLER versions. To set the default for a variable `'foo'`, use:

```
env.foo = 'value'
```

where `'env'` is an `environ` object (see section 4.2). Thus, the exact equivalent to both of the TOP scripts above would be:

```
env.gap_penalties_1d = (900, 50)
aln.align()
```

where `'aln'` is an `alignment` object (see section 5.2.2). However, it is recommended that instead you use the cleaner syntax:

```
aln.align(gap_penalties_1d=(900, 50))
```

This only sets the 1D gap penalties for this invocation of `align()`, and so is less likely to cause problems later in your script. If you want to call a routine several times with the same set of arguments, it is recommended that you save the arguments in local Python variables, use subroutines or classes, or use `'for'` loops.

5.2.2 TOP models and alignments

In TOP, commands may operate implicitly on one or more of the standard models or alignments in memory. For example, **ALIGN** always operates on `ALIGNMENT1`, **READ_MODEL** always operates on `MODEL1`, and **READ_MODEL2** always operates on `MODEL2`. `MODEL2` is a 'cut-down' model, used only for some operations (such as **SUPERPOSE**) and cannot be used to build full models, for example.

In Python, the models and alignments (and sequence databases, densities, *etc.*) are explicit, and are represented by classes. You can have as many models or alignments as you like, provided you have enough memory. Commands are simply methods of these classes. For example, consider the following:

```
env = environ()
aln = env.alignment()
aln.align(gap_penalties_1d=(900, 50))
```

This creates a new instance of the `environ` class (as above, this is used to provide default variables and the like), and calls it `'env'` (you can call it whatever you like). This is then used to create a new `alignment` class object, called `'aln'`. The following `align()` command then operates on the `'aln'` alignment object. (Note, however, that new alignments are empty, so this example wouldn't do anything interesting.)

5.2.3 TOP to Python correspondence

Please use the tables below to see which Python commands and variables correspond to old TOP commands and variables. Variables which are not listed in these tables have the same names as the old TOP equivalents (albeit in lower case).

TOP command	Python equivalent
ADD_RESTRAINT	<code>model.restraints.add()</code>
ALIGN	<code>alignment.align()</code>
ALIGN2D	<code>alignment.align2d()</code>
ALIGN3D	<code>alignment.align3d()</code>
ALIGN_CONSENSUS	<code>alignment.consensus()</code>
ALN_TO_PROF	<code>alignment.to_profile()</code>
BUILD_MODEL	<code>model.build()</code>
BUILD_PROFILE	<code>profile.build()</code>
CALL	use Python subroutines
CHECK_ALIGNMENT	<code>alignment.check()</code>
CLOSE	use Python file object
COLOR_ALN_MODEL	<code>model.color()</code>
COMPARE	<code>alignment.compare_structures()</code>
COMPARE_ALIGNMENTS	<code>alignment.compare_with()</code>
CONDENSE_RESTRAINTS	<code>model.restraints.condense()</code>
DEBUG_FUNCTION	<code>model.debug_function()</code>
DEFINE_INTEGER	use Python 'int' variables
DEFINE_LOGICAL	use Python 'bool' variables
DEFINE_REAL	use Python 'float' variables
DEFINE_STRING	use Python 'str' variables
DEFINE_SYMMETRY	<code>model.symmetry.define()</code>
DELETE_ALIGNMENT	<code>del(alignment)</code>
DELETE_FILE	<code>modfile.delete()</code>
DELETE_RESTRAINT	<code>model.restraints.unpick()</code>
DENDROGRAM	<code>environ.dendrogram()</code>
DESCRIBE	<code>alignment.describe()</code>
DO	use Python while or for loops
EDIT_ALIGNMENT	<code>alignment.edit()</code>
EM_GRID_SEARCH	<code>density.grid_search()</code>
END_SUBROUTINE	use Python subroutines
ENERGY	<code>model.energy()</code>
EXIT	use Python while or for loops
EXPAND_ALIGNMENT	<code>alignment.expand()</code>
GENERATE_TOPOLOGY	<code>model.generate_topology()</code>
ID_TABLE	<code>alignment.id_table()</code>
IF	use Python if statement
INCLUDE	use Python import statement
INQUIRE	<code>modfile.inquire()</code>
IUPAC_MODEL	<code>model.to_iupac()</code>
MAKE_CHAINS	<code>model.make_chains()</code>
MAKE_REGION	<code>model.make_region()</code>
MAKE_RESTRAINTS	<code>model.restraints.make()</code>
MAKE_SCHEDULE	<code>model.schedule.make()</code>
MAKE_TOPOLOGY_MODEL	<code>topology.make()</code>

MALIGN	alignment.malign()
MALIGN3D	alignment.malign3d()
MUTATE_MODEL	model.mutate()
OPEN	use Python file object
OPERATE	use Python arithmetic
OPTIMIZE	model.optimize()
ORIENT_MODEL	model.orient()
PATCH	model.patch()
PATCH_SS_MODEL	model.patch_ss()
PATCH_SS_TEMPLATES	model.patch_ss_templates()
PICK_ATOMS	model.pick_atoms()
PICK_HOT_ATOMS	model.pick_hot_atoms()
PICK_RESTRAINTS	model.restraints.pick()
PRINCIPAL_COMPONENTS	environ.principal_components()
PROFILE_PROFILE_SCAN	profile.scan()
PROF_TO_ALN	alignment.append_profile() or profile.to_alignment()
RANDOMIZE_XYZ	model.randomize_xyz()
READ	use Python file object
READ_ALIGNMENT	alignment.append()
READ_ALIGNMENT2	alignment.append()
READ_ATOM_CLASSES	group_restraints()
READ_MODEL	model.read()
READ_MODEL2	model.read()
READ_PARAMETERS	parameters.read() or group_restraints.append()
READ_PROFILE	profile.read()
READ_RESTRAINTS	model.restraints.append()
READ_RESTYP_LIB	environ()
READ_SCHEDULE	model.schedule.read()
READ_SEQUENCE_DB	sequence_db.read()
READ_TOPOLOGY	topology.read() or topology.append()
REINDEX_RESTRAINTS	model.restraints.reindex()
RENAME_SEGMENTS	model.rename_segments()
REORDER_ATOMS	model.reorder_atoms()
RESET	do not use
RETURN	use Python subroutines
ROTATE_DIHEDRALS	model.rotate_dihedrals()
ROTATE_MODEL	model.rotate()
SALIGN	alignment.salign()
SEGMENT_MATCHING	alignment.segment_matching()
SEQFILTER	sequence_db.filter()
SEQUENCE_COMPARISON	alignment.compare_sequences()
SEQUENCE_SEARCH	sequence_db.search()
SEQUENCE_TO_ALI	alignment.append_model()
SET	use Python variables
SPLINE_RESTRAINTS	model.restraints.spline()
STOP	do not use
STRING_IF	use Python if statement
STRING_OPERATE	use Python arithmetic
SUBROUTINE	use Python subroutines
SUPERPOSE	model.superpose()
SWITCH_TRACE	model.switch_trace()
SYSTEM	environ.system()
TIME_MARK	info.time_mark()
TRANSFER_RES_NUMB	model.res_num_from()
TRANSFER_XYZ	model.transfer_xyz()
UNBUILD_MODEL	model.unbuild()

WRITE	use Python file object
WRITE_ALIGNMENT	alignment.write()
WRITE_DATA	model.write_data()
WRITE_MODEL	model.write()
WRITE_MODEL2	model.write()
WRITE_PDB_XREF	model.write_pdb_xref()
WRITE_PROFILE	profile.write()
WRITE_RESTRAINTS	model.restraints.write()
WRITE_SCHEDULE	model.schedule.write()
WRITE_SEQUENCE_DB	sequence_db.write()
WRITE_TOP	do not use
WRITE_TOPOLOGY_MODEL	topology.write()

Table 5.1: Correspondence between TOP and Python commands.

TOP variable	Python equivalent
ALIGN_CODES[x]	alignment[x-1].code
ATOM_FILES[x]	alignment[x-1].atom_file
ATOM_FILES_DIRECTORY	io_data.atom_files_directory
CONTACT_SHELL	energy_data.contact_shell
COULOMB_SWITCH	energy_data.coulomb_switch
COVALENT_CYS	energy_data.covalent_cys
DYNAMIC_ACCESS	energy_data.dynamic_access
DYNAMIC_COULOMB	energy_data.dynamic_coulomb
DYNAMIC_LENARD	energy_data.dynamic_lennard
DYNAMIC_MODELLEDER	energy_data.dynamic_modeller
DYNAMIC_PAIRS	energy_data.dynamic_pairs
DYNAMIC_SPHERE	energy_data.dynamic_sphere
EXCL_LOCAL	energy_data.excl_local
HETATM_IO	io_data.hetatm
HYDROGEN_IO	io_data.hydrogen
LENNARD_JONES_SWITCH	energy_data.lennard_jones_switch
MOLPDF	return value from model.energy()
NLOGN_USE	energy_data.nlogn_use
NONBONDED_SEL_ATOMS	energy_data.nonbonded_sel_atoms
NUMB_OF_SEQUENCES	len(alignment)
N_SCHEDULE	len(model.schedule)
OUTPUT_CONTROL	use log.level()
RADII_FACTOR	energy_data.radii_factor
RELATIVE_DIELECTRIC	energy_data.relative_dielectric
SCHEDULE_STEP	model.schedule.step
SPHERE_STDV	energy_data.sphere_stdv
TOPOLOGY_MODEL	topology.submodel
UPDATE_DYNAMIC	energy_data.update_dynamic
WATER_IO	io_data.water

Table 5.2: Correspondence between TOP and Python variables.

Chapter 6

Methods

6.1 Dynamic programming for sequence and structure comparison and searching

In this section, the basic dynamic programming method for sequence alignment is described [Šali & Blundell, 1990]. This method forms the core of the pairwise and multiple sequence and structure comparisons as well as of the sequence database searching.

6.1.1 Pairwise comparison

The residue by residue scores W_{ij} can be used directly in the sequence alignment algorithm of Needleman & Wunsch [Needleman & Wunsch, 1970] to obtain the comparison of two protein sequences or structures. The only difference between the two types of comparison is in the type of the comparison matrix. In the case of sequence, the amino acid substitution matrix is used. In the case of 3D structure, the Euclidean distance (or some function of it) between two equivalent atoms in the current optimal superposition is used [Šali & Blundell, 1990].

The problem of the optimal alignment of two sequences as addressed by the algorithm of Needleman & Wunsch is as follows. We are given two sequences of elements and an M times N score matrix \mathcal{W} where M and N are the numbers of elements in the first and second sequence. The scoring matrix is composed of scores W_{ij} describing differences between elements i and j from the first and second sequence respectively. The goal is to obtain an optimal set of equivalences that match elements of the first sequence to the elements of the second sequence. The equivalence assignments are subject to the following “progression rule”: for elements i and k from the first sequence and elements j and l from the second sequence, if element i is equivalenced to element j , if element k is equivalenced to element l and if k is greater than i , l must also be greater than j . The optimal set of equivalences is the one with the smallest alignment score. The alignment score is a sum of scores corresponding to matched elements, also increased for occurrences of non-equivalenced elements (*ie* gaps). For a detailed discussion of this and related problems see [Sankoff & Kruskal, 1983].

We summarize the dynamic programming formulae used by MODELLER to obtain the optimal alignment since they differ slightly from those already published [Sellers, 1974, Gotoh, 1982]. The recursive dynamic programming formulae that give a matrix \mathcal{D} are:

$$\begin{aligned}
D_{i,j} &= \min \begin{cases} P_{i,j} \\ D_{i-1,j-1} + W_{i,j} \\ Q_{i,j} \end{cases} \\
P_{i,j} &= \min \begin{cases} D_{i-1,j} + g(1) \\ P_{i-1,j} + v \end{cases} \\
Q_{i,j} &= \min \begin{cases} D_{i,j-1} + g(1) \\ Q_{i,j-1} + v \end{cases}
\end{aligned} \tag{6.1}$$

where $g(l)$ is a linear gap penalty function:

$$g(l) = u + v \cdot l. \tag{6.2}$$

Note that only a vector is needed for the storage of P and Q . The uppermost formula in Eq. 6.1 is calculated for $i = M$ and $j = N$. Variable l is a gap length and parameters u and v are gap-penalty constants.

The arrays \mathcal{D} , \mathcal{P} and \mathcal{Q} are initialized as follows:

$$\begin{aligned}
D_{i,0} &= \begin{cases} 0, & i \leq e \\ g(i - e), & e < i \leq N \end{cases} \\
D_{0,j} &= \begin{cases} 0, & j \leq e \\ g(j - e), & e < j \leq N \end{cases} \\
P_{i,0} &= Q_{i,0} = \infty, & i = 1, 2, \dots, M \\
P_{0,j} &= Q_{0,j} = \infty, & j = 1, 2, \dots, N
\end{aligned} \tag{6.3}$$

where parameter e is the maximal number of elements at sequence termini which are not penalized with a gap-penalty if not equivalenced. A segment at the terminus of length e is termed an “overhang”. Note a difference from [Gotoh, 1982] in the initialization of the \mathcal{P} and \mathcal{Q} arrays. Also note that only vectors Q_i and P_j need to be stored in computer, not the whole arrays.

The minimal score $d_{M,N}$ is obtained from

$$d_{M,N} = \min(D_{i,N}, D_{M,j}) \tag{6.4}$$

where $i = M, M - 1, \dots, M - e$ and $j = N, N - 1, \dots, N - e$ to allow for the overhangs. The equivalence assignments are obtained by backtracking in matrix \mathcal{D} . Backtracking starts from the element $D_{i,j} = d_{M,N}$.

6.1.2 Variable gap penalty

This work is still in progress and is not described here.

6.1.3 Local *versus* global alignment

The Kruskal and Sankoff version of the local alignment is implemented [Sankoff & Kruskal, 1983]; this is very similar to the [Smith & Waterman, 1981] method. All the routines for the local alignment are exactly the same as the routines for the global alignment except that during the construction of matrix D the alignment is restarted each time the score becomes higher than a cutoff. The second difference is that the backtracking starts from the lowest element in the matrix, wherever it is.

6.1.4 Similarity *versus* distance scores

Each scoring matrix contains a flag determining whether it is a distance or similarity matrix. An appropriate optimization is used automatically. This is achieved by using exactly the same code except that one side of comparisons is multiplied by -1 when dealing with similarities as opposed to distances.

6.1.5 Multiple comparisons

In the discussion of the previous section, we have assumed that the sequences or structures would be compared in a pairwise manner. However, such pairwise comparisons of several related proteins may not be self consistent, *ie* the following transitivity rule can be broken: If residue a from protein A is equivalent to residue b in protein B which in turn is equivalent to residue c in protein C then the residue a from protein A must also be equivalent to residue c from protein C . This property is not always attained in the set of usual pairwise comparisons relating a group of similar proteins. For this reason we proceed by simultaneously aligning all proteins. This is achieved by aligning the second sequence with the first one, the third sequence with the alignment of the first two, *etc.* A more general tree-like growth of the multiple alignment is not yet implemented.

If the number of all proteins is N , $N - 1$ alignments must be made to obtain the final multiple comparison. It is noted that once an equivalence or gap is introduced it is not changed in later stages.

6.2 Optimization of the objective function by MODELLER

This section describes the optimization methods implemented in MODELLER. The general form of the objective function and the structure of optimization are similar to molecular dynamics programs, such as CHARMM [MacKerell *et al.*, 1998].

6.2.1 Function

MODELLER minimizes the *objective function* F with respect to Cartesian coordinates of $\sim 10,000$ atoms (3D points) that form a *system* (one or more molecules):

$$F = F(\mathbf{R}) = F_{symm} + \sum_i c_i(\mathbf{f}_i, \mathbf{p}_i) \quad (6.5)$$

where F_{symm} is an optional symmetry term defined in Eq. 6.72, \mathbf{R} are Cartesian coordinates of all atoms, c is a restraint i , \mathbf{f} is a geometric feature of a molecule, and \mathbf{p} are parameters. For a 10,000 atom system there can be on the order of 200,000 restraints. The form of c is simple; it includes a quadratic function, cosine, a weighted sum of a few Gaussian functions, Coulomb law, Lennard-Jones potential, cubic splines, and some other simple functions. The geometric features presently include a distance, an angle, a dihedral angle, a pair of dihedral angles between two, three, four atoms and eight atoms, respectively, the shortest distance in the set of distances (not documented further), solvent accessibility in \AA^2 , and atom density expressed as the number of atoms around the central atom. A pair of dihedral angles can be used to restrain such strongly correlated features as the mainchain dihedral angles Φ and Ψ . Each of the restraints also depends on a few parameters \mathbf{p}_i that generally vary from a restraint to a restraint. Some restraints can restrain *pseudo-atoms* such as a gravity center of several atoms.

MODELLER allows some atoms to be *fixed* during optimization; *ie.*, only selected atoms are allowed to be moved. Similarly, MODELLER also allows only a subset of all restraints to be actually used in the calculation of the objective function. Each subset is indicated by a list of indices specifying the selected atoms or restraints.

There are two kinds of restraints, *static* and *dynamic*, that both contribute to the objective function as indicated in Eq. 6.5:

$$F = F_{symm} + F_s + F_d . \quad (6.6)$$

The static restraints and their parameters are pre-defined; *ie.*, they are given before the call to the optimizer and are not changed during optimization. The dynamic restraints are re-generated repeatedly during optimization. Usually, the CPU time is spent evenly between the two kinds of restraints, although the dynamic restraints become

more important as the size of the system increases. All dynamic restraints are always selected and they can restrain only pairs of atoms. In all other respects, the two kinds of restraints are the same.

The dynamic restraints are obtained from a *dynamic pairs list* (the non-bonded pairs list). Each dynamic pair corresponds to at least one restraint, which may or may not be violated. The dynamic pairs list includes only the pairs of atoms that satisfy the following three conditions: (1) One or both atoms in a pair are allowed to move. (2) The two atoms are not connected through one, two, or three chemical bonds. (3) The two atoms are closer than a preset cutoff distance (*e.g.*, 4 Å). There are on the order of 5000 atom pairs in the dynamic pairs list when only soft-sphere overlap restraints are used. Currently, the restraint types on the dynamic atom pairs that can be selected include the soft-sphere overlap, Lennard-Jones, Coulombic interactions, and MODELLER non-bonded spline restraints. xx atom density?

The existence of the dynamic pairs list is justified by the fact that dynamic pairs are usually a small fraction of all possible atom–atom pairs ($N \cdot (N - 1)/2$, where N is the number of atoms in a system). The use of the dynamic pairs list becomes especially beneficent as the size of the system increases.

The actual algorithm for creating the dynamic pairs list varies with the size of the system, whether or not all atoms are allowed to move, or whether or not the user wants to include the fixed environment in the calculation of non-bonded restraints involving the selected atoms. See Section 4.6.35 for more information.

The hash-function algorithm is used to determine whether or not two atoms are a dynamic atom pair. This algorithm is about 20 times slower than a lookup table but it requires much less memory and still spends a negligible fraction of the total CPU time. A hash-function table is prepared only once before the start of the optimization and any other operation involving an evaluation of the objective function (*e.g.*, `model.optimize()`, `model.energy()`, and `model.pick_hot_atoms()`).

The dynamic pairs list is not necessarily re-generated each time the objective function is evaluated, although the contribution of the restraint to the objective function is calculated in each call to the objective function routine with the current values of the Cartesian coordinates. The dynamic pairs list is re-generated only when maximal atomic shifts accumulate to a value larger than a preset cutoff. This cutoff is chosen such that there cannot be a violation of a restraint without having its atom pair on the dynamic pairs list. The dynamic pairs list is recalculated in $\sim 20\%$ and $\sim 2\%$ of the objective function calls at the beginning and the end of optimization, respectively.

Each evaluation of the objective function or of its first derivatives with respect to the Cartesian coordinates involves the following steps:

1. Calculate non-fixed pseudo-atoms from the current atomic positions (routine `objfunc:pseudo`).
2. Update the dynamic pairs list, if necessary (routine `objfunc:upddyn`).
3. Calculate the violations of selected restraints and all other quantities that are shared between the calculations of the objective function and its derivatives (routine `objfunc:getviol`).
4. Sum the contributions of all violated restraints to the objective function and the derivatives (routine `objfunc:getviol`).

6.2.2 Optimizers

MODELLER currently implements a Beale restart conjugate gradients algorithm [Shanno & Phua, 1980, Shanno & Phua, 1982] and a molecular dynamics procedure with the Verlet integrator [Verlet, 1967]. The conjugate gradients optimizer is usually used in combination with the variable target function method [Braun & Gö, 1985] which is implemented with the `automodel` class (Section 6.4). The molecular dynamics procedure can be used in a simulated annealing protocol that is also implemented with the `automodel` class.

Molecular dynamics

Force in MODELLER is obtained by equating the objective function F with internal energy in kcal/mole. The atomic masses are all set to that of C¹² (MODELLER unit is kg/mole). The initial velocities at a given temperature

are obtained from a Gaussian random number generator with a mean and standard deviation of:

$$\bar{v}_x = 0 \quad (6.7)$$

$$\sigma_x = \sqrt{\frac{k_B T}{m}} = 0.000263143\sqrt{T} \quad (6.8)$$

where k_B is the Boltzmann constant, m is the mass of one C¹² atom, and the velocity is expressed in angstroms/femtosecond.

The Newton's equations of motion are integrated by the Verlet algorithm [Verlet, 1967]:

$$v_x(i+1) = v_x(i) + \frac{\partial F}{\partial x} A \quad (6.9)$$

$$x(i+1) = x(i) + v_x(i+1)\Delta t \quad (6.10)$$

$$A = c \frac{\Delta t}{m} = 4.1868 \cdot 10^{-7} \frac{\Delta t}{m} \quad (6.11)$$

where velocities $v(i+1)$ are for $t + \Delta t/2$ and positions $x(i+1)$ for $t + \Delta t$. Parameter c is a scaling factor so that positions are expressed in angstroms, time in femtoseconds, and velocities in angstroms/femtosecond, given that the objective function is in kcal/mole and atomic mass in kg/mole. In addition, velocity is capped at a maximum value, before calculating the shift, such that the maximal shift along one axis can only be `cap_atom_shift`. The velocities can be equilibrated every `equilibrate` steps to stabilize temperature. This is achieved by scaling the velocities with a factor f :

$$f = \sqrt{T/E_{kin}} \quad (6.12)$$

$$E_{kin} = \frac{m}{2} \sum_i^{N_{atoms}} (v_x^2 + v_y^2 + v_z^2) \quad (6.13)$$

where E_{kin} is the current kinetic energy of the system.

6.3 Equations used in the derivation of the molecular pdf

6.3.1 Features and their derivatives

Distance

Distance is defined by points i and j :

$$d = \sqrt{\vec{r}_{ij} \cdot \vec{r}_{ij}} = |\vec{r}_{ij}| = r_{ij} \quad (6.14)$$

where

$$\vec{r}_{ij} = \vec{r}_i - \vec{r}_j. \quad (6.15)$$

The first derivatives of d with respect to Cartesian coordinates are:

$$\frac{\partial d}{\partial \vec{r}_i} = \frac{\vec{r}_{ij}}{|\vec{r}_{ij}|} \quad (6.16)$$

$$\frac{\partial d}{\partial \vec{r}_j} = -\frac{\partial d}{\partial \vec{r}_i} \quad (6.17)$$

Angle

Angle is defined by points i , j , and k , and spanned by vectors ij and kj :

$$\alpha = \arccos \frac{\vec{r}_{ij} \cdot \vec{r}_{kj}}{|\vec{r}_{ij}| |\vec{r}_{kj}|}. \quad (6.18)$$

It lies in the interval from 0 to 180°. Internal MODELLER units are radians.

The first derivatives of α with respect to Cartesian coordinates are:

$$\frac{\partial \alpha}{\partial \vec{r}_i} = \frac{\partial \alpha}{\partial \cos \alpha} \frac{\partial \cos \alpha}{\partial \vec{r}_i} = \frac{1}{\sqrt{1 - \cos^2 \alpha}} \frac{1}{r_{ij}} \left(\frac{\vec{r}_{ij}}{r_{ij}} \cos \alpha - \frac{\vec{r}_{kj}}{r_{kj}} \right) \quad (6.19)$$

$$\frac{\partial \alpha}{\partial \vec{r}_k} = \frac{\partial \alpha}{\partial \cos \alpha} \frac{\partial \cos \alpha}{\partial \vec{r}_k} = \frac{1}{\sqrt{1 - \cos^2 \alpha}} \frac{1}{r_{kj}} \left(\frac{\vec{r}_{kj}}{r_{kj}} \cos \alpha - \frac{\vec{r}_{ij}}{r_{ij}} \right) \quad (6.20)$$

$$\frac{\partial \alpha}{\partial \vec{r}_j} = -\frac{\partial d}{\partial \vec{r}_i} - \frac{\partial d}{\partial \vec{r}_k} \quad (6.21)$$

These equations for the derivatives have a numerical instability when the angle goes to 0 or to 180°. Presently, the problem is ‘solved’ by testing for the size of the angle; if it is too small, the derivatives are set to 0 in the hope that other restraints will eventually pull the angle towards well behaved regions. Thus, angle restraints of 0 or 180° should not be used in the conjugate gradients or molecular dynamics optimizations.

Dihedral angle

Dihedral angle is defined by points i, j, k , and l ($ijkl$):

$$\chi = \text{sign}(\chi) \arccos \frac{(\vec{r}_{ij} \times \vec{r}_{kj}) \cdot (\vec{r}_{kj} \times \vec{r}_{kl})}{|\vec{r}_{ij} \times \vec{r}_{kj}| |\vec{r}_{kj} \times \vec{r}_{kl}|} \quad (6.22)$$

where

$$\text{sign}(\chi) = \text{sign}[\vec{r}_{kj} \cdot (\vec{r}_{ij} \times \vec{r}_{kj}) \times (\vec{r}_{kj} \times \vec{r}_{kl})]. \quad (6.23)$$

The first derivatives of χ with respect to Cartesian coordinates are:

$$\frac{d\chi}{d\vec{r}} = \frac{d\chi}{d \cos \chi} \frac{d \cos \chi}{d\vec{r}} \quad (6.24)$$

where

$$\frac{d\chi}{d \cos \chi} = \left(\frac{d \cos \chi}{d\chi} \right)^{-1} = -\frac{1}{\sin \chi} \quad (6.25)$$

and

$$\frac{\partial \cos \chi}{\partial \vec{r}_i} = \vec{r}_{kj} \times \vec{a} \quad (6.26)$$

$$\frac{\partial \cos \chi}{\partial \vec{r}_j} = \vec{r}_{ik} \times \vec{a} - \vec{r}_{kl} \times \vec{b} \quad (6.27)$$

$$\frac{\partial \cos \chi}{\partial \vec{r}_k} = \vec{r}_{jl} \times \vec{b} - \vec{r}_{ij} \times \vec{a} \quad (6.28)$$

$$\frac{\partial \cos \chi}{\partial \vec{r}_l} = \vec{r}_{ij} \times \vec{b} \quad (6.29)$$

$$\vec{a} = \frac{1}{|\vec{r}_{ij} \times \vec{r}_{kj}|} \left(\frac{\vec{r}_{kj} \times \vec{r}_{kl}}{|\vec{r}_{kj} \times \vec{r}_{kl}|} - \cos \chi \frac{\vec{r}_{ij} \times \vec{r}_{kj}}{|\vec{r}_{ij} \times \vec{r}_{kj}|} \right) \quad (6.30)$$

$$\vec{b} = \frac{1}{|\vec{r}_{kj} \times \vec{r}_{kl}|} \left(\frac{\vec{r}_{ij} \times \vec{r}_{kj}}{|\vec{r}_{ij} \times \vec{r}_{kj}|} - \cos \chi \frac{\vec{r}_{kj} \times \vec{r}_{kl}}{|\vec{r}_{kj} \times \vec{r}_{kl}|} \right). \quad (6.31)$$

These equations for the derivatives have a numerical instability when the angle goes to 0. Thus, the following set of equations is used instead [van Schaik *et al.*, 1993]:

$$\vec{r}_{mj} = \vec{r}_{ij} \times \vec{r}_{kj} \quad (6.32)$$

$$\vec{r}_{nk} = \vec{r}_{kj} \times \vec{r}_{kl} \quad (6.33)$$

$$\frac{\partial \chi}{\partial \vec{r}_i} = \frac{r_{kj}}{r_{mj}^2} \vec{r}_{mj} \quad (6.34)$$

$$\frac{\partial \chi}{\partial \vec{r}_l} = -\frac{r_{kj}}{r_{nk}^2} \vec{r}_{nk} \quad (6.35)$$

$$\frac{\partial \chi}{\partial \vec{r}_j} = \left(\frac{\vec{r}_{ij} \cdot \vec{r}_{kj}}{r_{kj}^2} - 1 \right) \frac{\partial \chi}{\partial \vec{r}_i} - \frac{\vec{r}_{kl} \cdot \vec{r}_{kj}}{r_{kj}^2} \frac{\partial \chi}{\partial \vec{r}_l} \quad (6.36)$$

$$\frac{\partial \chi}{\partial \vec{r}_k} = \left(\frac{\vec{r}_{kl} \cdot \vec{r}_{kj}}{r_{kj}^2} - 1 \right) \frac{\partial \chi}{\partial \vec{r}_l} - \frac{\vec{r}_{ij} \cdot \vec{r}_{kj}}{r_{kj}^2} \frac{\partial \chi}{\partial \vec{r}_i} \quad (6.37)$$

The only possible instability in these equations is when the length of the central bond of the dihedral, r_{kj} , goes to 0. In such a case, which should not happen, the derivatives are set to 0. The expressions for an improper dihedral angle, as opposed to a dihedral or dihedral angle, are the same, except that indices $ijkl$ are permuted to $ikjl$. In both cases, covalent bonds ij , jk , and kl are defining the angle.

Atomic solvent accessibility

xx

Atomic density

Atomic density for a given atom is simply calculated as the number of atoms within a distance `energy_data.contact_shell` of that atom. First derivatives are not calculated, and are always returned as 0.

Atomic coordinates

The absolute atomic coordinates x_i , y_i and z_i are available for every point i , primarily for use in anchoring points to planes, lines or points. Their first derivatives with respect to Cartesian coordinates are of course simply 0 or 1.

6.3.2 Restraints and their derivatives

The chain rule is used to find the partial derivatives of the feature pdf with respect to the atomic coordinates. Thus, only the derivatives of the pdf with respect to the features are listed here.

Single Gaussian restraint

The pdf for a geometric feature f (e.g., distance, angle, dihedral angle) is

$$p = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{1}{2} \left(\frac{f - \bar{f}}{\sigma} \right)^2 \right]. \quad (6.38)$$

A corresponding restraint c in the sum that defines the objective function F is

$$c = -\ln p = \frac{1}{2} \left(\frac{f - \bar{f}}{\sigma} \right)^2 - \ln \frac{1}{\sigma\sqrt{2\pi}} \quad (6.39)$$

The first derivatives with respect to feature f are:

$$\frac{dc}{df} = \frac{f - \bar{f}}{\sigma} \frac{1}{\sigma}. \quad (6.40)$$

Multiple Gaussian restraint

The polymodal pdf for a geometric feature f (*e.g.*, distance, angle, dihedral angle) is

$$p = \sum_{i=1}^n \omega_i p_i = \sum_{i=1}^n \omega_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp \left[-\frac{1}{2} \left(\frac{f - \bar{f}_i}{\sigma_i} \right)^2 \right]. \quad (6.41)$$

A corresponding restraint c in the sum that defines the objective function F is

$$c = -\ln p = -\ln \sum_{i=1}^n \omega_i p_i \quad (6.42)$$

The first derivatives with respect to feature f are:

$$\frac{dc}{df} = \frac{1}{p} \sum_{i=1}^n \omega_i p_i \cdot \left[\frac{f - \bar{f}_i}{\sigma_i} \frac{1}{\sigma_i} \right]. \quad (6.43)$$

When any of the normalized deviations $v_i = (f - \bar{f}_i)/\sigma_i$ is large, there are numerical instabilities in calculating the derivatives because v_i are arguments to the *exp* function. Robustness is ensured as follows. The ‘effective’ normalized deviation is used in all the equations above when the magnitude of normalized violation v is larger than cutoff **rgauss1** (10 for double precision). This scheme works up to **rgauss2** (200 for double precision); violations larger than that are ignored. This trick is equivalent to increasing the standard deviation σ_i . A slight disadvantage is that there is a discontinuity in the first derivatives at **rgauss1**. However, if continuity were imposed, the range would not be extended (this is equivalent to linearizing the Gaussian, but since it is already linear for large deviations, a linearization with derivatives smoothness would not introduce much change at all).

$$M = 37 \quad ; \quad M^2/2 \text{ has to be smaller than the largest argument to exp} \quad (6.44)$$

$$A = \frac{1}{M} \frac{\text{rgauss2} - M}{\text{rgauss2} - \text{rgauss1}} \quad (6.45)$$

$$B = \frac{\text{rgauss2}}{M} \frac{M - \text{rgauss1}}{\text{rgauss2} - \text{rgauss1}} \quad (6.46)$$

$$v = \frac{f - \bar{f}_i}{\sigma_i} \quad (6.47)$$

$$F = A |v| + B \quad (6.48)$$

$$v' = v/F \quad (6.49)$$

Now, Eqs. 6.41–6.43 are used with v' instead of v . For single precision, $M = 12$, **rgauss1** = 4, **rgauss2** = 100.

Multiple binormal restraint

The polymodal pdf for a geometric feature (f_1, f_2) (*e.g.*, a pair of dihedral angles) is

$$p = \sum_{i=1}^n \omega_i p_i = \sum_{i=1}^n \omega_i \frac{1}{2\pi \sigma_{1i} \sigma_{2i} \sqrt{(1 - \rho_i^2)}} \cdot \exp \left\{ -\frac{1}{2(1 - \rho_i^2)} \left[\left(\frac{f_1 - \bar{f}_{1i}}{\sigma_{1i}} \right)^2 - 2\rho_i \frac{f_1 - \bar{f}_{1i}}{\sigma_{1i}} \frac{f_2 - \bar{f}_{2i}}{\sigma_{2i}} + \left(\frac{f_2 - \bar{f}_{2i}}{\sigma_{2i}} \right)^2 \right] \right\}. \quad (6.50)$$

where $\rho < 1$. ρ is the correlation coefficient between f_1 and f_2 . MODELLER actually uses the following series expansion to calculate p :

$$p = \sum_{i=1}^n \frac{1}{2\pi\sigma_{1i}\sigma_{2i}\sqrt{(1-\rho_i^2)}} \cdot \exp \left\{ -\frac{1}{1-\rho_i^2} \left[\frac{1-\cos(f_1-\bar{f}_{1i})}{\sigma_{1i}^2} - \rho_i \frac{\sin(f_1-\bar{f}_{1i})}{\sigma_{1i}} \frac{\sin(f_2-\bar{f}_{2i})}{\sigma_{2i}} + \frac{1-\cos(f_2-\bar{f}_{2i})}{\sigma_{2i}^2} \right] \right\}. \quad (6.51)$$

A corresponding restraint c in the sum that defines the objective function F is

$$c = -\ln p = -\ln \sum_{i=1}^n \omega_i p_i \quad (6.52)$$

The first derivatives with respect to features f_1 and f_2 are:

$$\frac{\partial c}{\partial f_1} = \frac{1}{p} \sum_{i=1}^n \left[\omega_i p_i \cdot \frac{1}{\sigma_{1i}(1-\rho_i^2)} \left(\frac{\sin(f_1-\bar{f}_{1i})}{\sigma_{1i}} - \rho_i \frac{\cos(f_1-\bar{f}_{1i})\sin(f_2-\bar{f}_{2i})}{\sigma_{2i}} \right) \right] \quad (6.53)$$

$$\frac{\partial c}{\partial f_2} = \frac{1}{p} \sum_{i=1}^n \left[\omega_i p_i \cdot \frac{1}{\sigma_{2i}(1-\rho_i^2)} \left(\frac{\sin(f_2-\bar{f}_{2i})}{\sigma_{2i}} - \rho_i \frac{\cos(f_2-\bar{f}_{2i})\sin(f_1-\bar{f}_{1i})}{\sigma_{1i}} \right) \right]. \quad (6.54)$$

Lower bound

This is like the left half of a single Gaussian restraint:

$$p = \begin{cases} p_{gauss} & ; \quad f < \bar{f} \\ 0 & ; \quad f \geq \bar{f} \end{cases} \quad (6.55)$$

where \bar{f} is a lower bound and p_{gauss} is given in Eq. 6.38. A similar equation relying on the first derivatives of a Gaussian p holds for the first derivatives of a lower bound.

Upper bound

This is like the right half of a single Gaussian restraint:

$$p = \begin{cases} p_{gauss} & ; \quad f > \bar{f} \\ 0 & ; \quad f \leq \bar{f} \end{cases} \quad (6.56)$$

where \bar{f} is an upper bound and p_{gauss} is given in Eq. 6.38. A similar equation relying on the first derivatives of a Gaussian p holds for the first derivatives of an upper bound.

Cosine restraint

This is usually used for dihedral angles f :

$$c = |b| - b \cos(nf + a) \quad (6.57)$$

where b is CHARMM force constant, a is phase shift (tested for 0 and 180°), and n is periodicity (tested for 1, 2, 3, 4, 5, and 6). The CHARMM phase value from the CHARMM parameter library corresponds to $a - 180^\circ$. The force constant b can be negative, in effect offsetting the phase a for 180° compared to the same but positive force constant.

$$\frac{dc}{df} = bn \sin(nf + a) \quad (6.58)$$

Coulomb restraint

$$c = \frac{1}{\epsilon_r} \frac{q_i q_j}{f} s(f, f_1, f_2) \quad (6.59)$$

$$s(f, f_1, f_2) = \begin{cases} 1; & f \leq f_1 \\ \frac{(f_2-f)^2(f_2+2f-3f_1)}{(f_2-f_1)^3}; & f_1 < f \leq f_2 \\ 0; & f > f_2 \end{cases} \quad (6.60)$$

where q_i and q_j are the atomic charges of atoms i and j , obtained from the CHARMM topology file, that are at a distance f . ϵ_r is the relative dielectric, controlled by the `energy.data.relative.dielectric` variable. Function $s(f, f_1, f_2)$ is a switching function that smoothes the potential down to zero in the interval from f_1 to f_2 ($f_2 > f_1$). The total Coulomb energy of a molecule is a sum over all pairs of atoms that are not in the same bonds or bond angles. 1-4 energy for the 1-4 atom pairs in the same dihedral angle corresponds to the ELEC14 MODELLER term; the remaining longer-range contribution corresponds to the ELEC term.

The first derivatives are:

$$\frac{dc}{df} = -\frac{c}{f} + \frac{c}{s} \frac{ds}{df} \quad (6.61)$$

$$\frac{ds}{df} = \begin{cases} 0; & f \leq f_1 \\ \frac{6(f_2-f)(f_1-f)}{(f_2-f_1)^3}; & f_1 < f \leq f_2 \\ 0; & f > f_2 \end{cases} \quad (6.62)$$

Lennard-Jones restraint

Usually used for non-bonded distances:

$$c = \left[\left(\frac{A}{f} \right)^{12} - \left(\frac{B}{f} \right)^6 \right] s(f, f_1, f_2) \quad (6.63)$$

The parameters f_1 and f_2 of the switching function can be different from those in Eq. 6.60. The parameters A and B are obtained from the CHARMM parameter file (NONBOND section) where they are given as E_i and r_j such that $E_{ij}(f) = -4\sqrt{E_i E_j}[(\rho_{ij}/f)^{12} - (\rho_{ij}/f)^6]$ in kcal/mole for f in angstroms and $\rho = (r_i + r_j)/2^{1/6}$; the minimum of E is $-\sqrt{E_i E_j}$ at $f = (r_i + r_j)$, and its zero is at $f = \rho$. The total Lennard-Jones energy should be evaluated over all pairs of atoms that are not in the same bonds or bond angles. The parameters A and B for 1-4 pairs in dihedral angles can be different from those for the other pairs; they are obtained from the second set of E_i and r_i in the CHARMM parameter file, if it exists. 1-4 energy corresponds to the LJ14 MODELLER term; the remaining longer-range contribution corresponds to the LJ term.

The first derivatives are:

$$\frac{dc}{df} = \frac{Cs}{f} - C \frac{ds}{df} \quad (6.64)$$

$$C = -12 \left(\frac{A}{f} \right)^{12} + 6 \left(\frac{B}{f} \right)^6 \quad (6.65)$$

Spline restraint

Any restraint form can be represented by a cubic spline [Press *et al.*, 1992]:

$$c = Ac_j + Bc_{j+1} + Cc_j'' + Dc_{j+1}'' \quad (6.66)$$

$$A = \frac{f_{j+1} - f}{f_{j+1} - f_j} \quad (6.67)$$

$$B = 1 - A \quad (6.68)$$

$$C = \frac{1}{6}(A^3 - A)(f_{j+1} - f_j)^2 \quad (6.69)$$

$$D = \frac{1}{6}(B^3 - B)(f_{j+1} - f_j)^2 \quad (6.70)$$

where $f_j \leq f \leq f_{j+1}$.

The first derivatives are:

$$\frac{dc}{df} = \frac{c_{j+1} - c_j}{f_{j+1} - f_j} - \frac{3A^2 - 1}{6}(f_{j+1} - f_j)c_j'' + \frac{3B^2 - 1}{6}(f_{j+1} - f_j)c_{j+1}'' \quad (6.71)$$

The values of c and c' beyond f_1 and f_n are obtained by linear interpolation from the termini. A violation of the restraint is calculated by finding the global minimum. A relative violation is estimated by using a standard deviation (*e.g.*, force constant) obtained by fitting a parabola to the global minimum.

Variable spacing of spline points could be used to save on memory. However, this would increase the execution time, so it is not used.

Symmetry restraint

The asymmetry penalty added to the objective function is defined as

$$F_{symm} = \sum_{i < j} \omega_i \omega_j (d_{ij} - d'_{ij})^2 \quad (6.72)$$

where the sum runs over all pairs of equivalent atoms ij , ω_i is an atom weight for atom i , d_{ij} is an intra-molecular distance between atoms ij in the first segment, and d'_{ij} is the equivalent distance in the second segment.

For each $i < j$, the first derivatives are:

$$\frac{\partial c}{\partial d_{ij}} = 2\omega_i \omega_j (d_{ij} - d'_{ij}) \frac{\vec{d}_{ij}}{d_{ij}} \quad (6.73)$$

$$\frac{\partial c}{\partial d'_{ij}} = -2\omega_i \omega_j (d_{ij} - d'_{ij}) \frac{\vec{d}'_{ij}}{d'_{ij}} \quad (6.74)$$

Thus, the total first derivatives are obtained by summing the two expressions above for all i and $j > i$ distances.

6.4 Flowchart of comparative modeling by MODELLER

This section describes a flowchart of comparative modeling by MODELLER, as implemented in the `automodel` class (see chapter 2).

Input: script file, alignment file, PDB file(s) for template(s).

Output:

<code>job.log</code>	log file
<code>job.ini</code>	initial conformation for optimization
<code>job.rsr</code>	restraints file
<code>job.sch</code>	VTFM schedule file
<code>job.B9999????</code>	PDB atom file(s) for the model(s) of the target sequence
<code>job.V9999????</code>	violation profiles for the model(s)
<code>job.D9999????</code>	progress of optimization
<code>job.BL9999????</code>	optional loop model(s)
<code>job.DL9999????</code>	progress of optimization for loop model(s)
<code>job.IL9999????</code>	initial structures for loop model(s)

The main MODELLER routines used in each step are given in parentheses.

1. Read and check the alignment between the target sequence and the template structures (**alignment.append()** and **alignment.check()**).
2. Calculate restraints on the target from its alignment with the templates:
 - (a) Generate molecular topology for the target sequence (**model.generate_topology()**). Disulfides in the target are assigned here from the equivalent disulfides in the templates (**model.patch_ss_templates()**). Any user defined patches are also done here (as defined in the **automodel.special_patches()** routine).
 - (b) Calculate coordinates for atoms that have equivalent atoms in the templates as an average over all templates (**model.transfer_xyz()**) (alternatively, read the initial coordinates from a file).
 - (c) Build the remaining unknown coordinates using internal coordinates from the CHARMM topology library (**model.build()**).
 - (d) Write the initial model to a file with extension **.ini** (**model.write()**).
 - (e) Generate stereochemical, homology-derived, and special restraints (**model.restraints.make()**) (alternatively, skip this and assume the restraints file already exists):

stereochemical	restraint_type = 'bond angle dihedral improper'
mainchain dihedrals Φ, Ψ	restraint_type = 'phi-psi-binormal'
mainchain dihedral ω	restraint_type = 'omega_dihedral'
sidechain dihedral χ_1	restraint_type = 'chi1_dihedral'
sidechain dihedral χ_2	restraint_type = 'chi2_dihedral'
sidechain dihedral χ_3	restraint_type = 'chi3_dihedral'
sidechain dihedral χ_4	restraint_type = 'chi4_dihedral'
mainchain CA-CA distance	restraint_type = 'distance'
mainchain N-O distance	restraint_type = 'distance'
sidechain-mainchain distance	restraint_type = 'distance'
sidechain-sidechain distance	restraint_type = 'distance'
block distance restraints	restraint_type = 'distance'
user defined	automodel.special_restraints() routine
non-bonded pairs distance	restraint_type = 'sphere'; calculated on the fly
 - (f) Write all restraints to a file with extension **.rsr** (**model.restraints.write()**).
3. Calculate model(s) that satisfy the restraints as well as possible. For each model:
 - (a) Generate the optimization schedule for the variable target function method (VTFM) (**model.schedule.make()**).
 - (b) Read the initial model (usually from the **.ini** file from 2.d) (**model.read()**).
 - (c) Randomize the initial structure by adding a random number between \pm automodel.deviation angstroms to all atomic positions (**model.randomize_xyz()**).
 - (d) optimize the model:
 - Partially optimize the model by VTFM; Repeat the following steps as many times as specified by the optimization schedule:
 - Read all the restraints by 'rd_restraints' (**model.restraints.append()**).
 - Select only the restraints that operate on the atoms that are close enough in sequence, as specified by the current step of VTFM (**model.restraints.pick()**).
 - Optimize the model by conjugate gradients, using only currently selected restraints (**model.optimize()**).
 - Refine the model by simulated annealing with molecular dynamics, if so selected:
 - do a short conjugate gradients optimization (**model.optimize()**).
 - increase temperature in several steps and do molecular dynamics optimization at each temperature (**model.optimize()**).
 - decrease temperature in several steps and do molecular dynamics optimization at each temperature (**model.optimize()**).

- do a short conjugate gradients optimization (**model.optimize()**).
- (e) calculate the remaining restraint violations and write them out (**model.energy()**).
- (f) write out the final model to a file with extension **.B9999????** where **????** indicates the model number (**model.write()**). Also write out the violations profile. Also write superposed templates and models if so selected by **automodel.final_malign3d = True**.
- (g) superpose the models and the templates, if so selected, and write them out (**alignment.expand()**, **alignment.malign3d()**).
- (h) do loop modeling if so selected using the **loopmodel** class.

6.5 Loop modeling method

The loop modeling method first takes the generated model, and selects all standard residues around gaps in the alignment for additional loop modeling. (To select a different region for modeling, simply redefine the **loopmodel.select_loop_atoms()** routine to select the relevant atoms.) An initial loop conformation is then generated by simply positioning the atoms of the loop with uniform spacing on the line that connects the main-chain carbonyl oxygen and amide nitrogen atoms of the N- and C-terminal anchor regions respectively, and this model is written out to a file with the **.IL** extension.

Next, a number of loop models are generated from **loopmodel.loop.starting_model** to **loopmodel.loop.ending_model**. Each takes the initial loop conformation and randomizes it by $\pm 5\text{\AA}$ in each of the Cartesian directions. The model is then optimized thoroughly twice, firstly considering only the loop atoms and secondly with these atoms “feeling” the rest of the system. The loop optimization relies on an atomistic distance-dependent statistical potential of mean force for nonbond interactions [Melo & Feytmans, 1997]. This classifies all amino acid atoms into one of 40 atom classes (as defined in **\$LIB/atmcls-melo.lib**) and applies a potential as MODELLER cubic spline restraints (as defined in **\$LIB/melo-dist1.lib**). No homology-derived restraints are used during this procedure. Each loop model is written out with the **.BL** extension.

For more information, please consult the loop modeling paper [Fiser *et al.*, 2000] or look at the loop modeling class itself in **modlib/modeller/automodel/loopmodel.py**.

Bibliography

- Braun, W. & Gö, N. (1985). *J. Mol. Biol.* **186**, 611–626.
- Brünger, A. T. (1992). *X-PLOR Manual Version 3.1*. Yale University New Haven, Connecticut.
- Felsenstein, J. (1985). *Evolution*, **39**, 783–791.
- Fiser, A., Do, R. K. G., & Šali, A. (2000). *Protein Sci.* **9**, 1–21. (Also available online).
- Gotoh, O. (1982). *J. Mol. Biol.* **162**, 705–708.
- Hubbard, T. J. P. & Blundell, T. L. (1987). *Protein Eng.* **1**, 159–171.
- IUPAC-IUB (1970). *Biochem.* **9**, 3471–3479.
- Kabsch, W. & Sander, C. (1983). *Biopolymers*, **22**, 2577–2637.
- Karlin, S. & Altschul, S. F. (1990). *Proc. Natl. Acad. Sci. USA*, **87**, 2264–2268.
- Kendrew, J. C., Klyne, W., Lifson, S., Miyazawa, T., Némethy, G., Phillips, D. C., Ramachandran, G. N., & Scheraga, H. (1970). *J. Mol. Biol.* **52**, 1–17.
- MacKerell, Jr., A. D., Bashford, D., Bellott, M., Dunbrack Jr., R. L., Evanseck, J. D., Field, M. J., Fischer, S., Gao, J., Guo, H., Ha, S., Joseph-McCarthy, D., Kuchnir, L., Kuczera, K., Lau, F. T. K., Mattos, C., Michnick, S., Ngo, T., Nguyen, D. T., Prodhom, B., Reiher, III, W. E., Roux, B., Schlenkrich, M., Smith, J. C., Stote, R., Straub, J., Watanabe, M., Wiorkiewicz-Kuczera, J., Yin, D., & Karplus, M. (1998). *J. Phys. Chem. B*, **102**, 3586–3616.
- Martí-Renom, M. A., Madhusudhan, M. S., & Šali, A. (2004). *Protein Sci.* **13**, 1071–1087. (Also available online).
- Melo, F. & Feytmans, E. (1997). *J. Mol. Biol.* **267**, 207–222.
- Melo, F., Sánchez, R., & Šali, A. (2002). *Protein Sci.* **11**, 430–448. (Also available online).
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., & Teller, E. (1953). *J. Chem. Phys.* **21**, 1087–1092.
- Needleman, S. B. & Wunsch, C. D. (1970). *J. Mol. Biol.* **48**, 443–453.
- Nicholls, A., Sharp, K. A., & Honig, B. (1991). *Proteins*, **11**, 281–296.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., & Flannery, B. P. (1992). *Numerical Recipes, 2nd edition*. Cambridge: Cambridge University Press.
- Richards, F. M. & Kundrot, C. E. (1988). *Proteins*, **3**, 71–84.
- Richmond, T. J. & Richards, F. M. (1978). *J. Mol. Biol.* **119**, 537–555.
- Šali, A. & Blundell, T. L. (1990). *J. Mol. Biol.* **212**, 403–428. (Also available online).
- Šali, A. & Blundell, T. L. (1993). *J. Mol. Biol.* **234**, 779–815. (Also available online).
- Šali, A. & Overington, J. (1994). *Protein Sci.* **3**, 1582–1596. (Also available online).

- Sankoff, D. & Kruskal, J. B. (1983). *Time warps, string edits, and macromolecules: The theory and practice of sequence comparison*. Reading, MA: Addison-Wesley Publishing Company.
- Sellers, P. H. (1974). *J. Comb. Theor.* **A16**, 253–258.
- Shanno, D. F. & Phua, K. H. (1980). *ACM Trans. Math. Soft.* **6**, 618–622.
- Shanno, D. F. & Phua, K. H. (1982). In: *Collected algorithms from ACM. Trans. Math. Software* volume 2(1).
- Smith, T. F. & Waterman, M. S. (1981). *J. Mol. Biol.* **147**, 195–197.
- Subbiah, S., Laurents, D. V., & Levitt, M. (1993). *Curr. Biol.* **3**, 141–148.
- Sutcliffe, M. J., Haneef, I., Carney, D., & Blundell, T. L. (1987). *Protein Eng.* **1**, 377–384.
- Topf, M., Baker, M. L., John, B., Chiu, W., & Šali, A. (2005). *J. Struct. Biol.* **149**, 191–203. (Also available online).
- van Schaik, R. C., Berendsen, H. J., & Torda, A. E. (1993). *J. Mol. Biol.* **234**, 751–762.
- Verlet, J. (1967). *Phys. Rev.* **159**, 98–103.

Index

- accessibility_type, 74, 98
- add_segment, 56, 57
- add_symmetry, 102
- align3d_repeat, 128
- align3d_trf, 128, 131, 133
- align_alignment, 117, 118
- align_block, 117, 118, 121–123, 125, 126, 131, 135, 142, 143
- align_codes, 114, 117
- align_what, 121–123, 131, 135
- alignment, 10, 111
- alignment()**, 113
- alignment.align**, 136
- alignment.align()**, 118, 121, 123, 125, 126, 128, 133, 135, 137, 165
- alignment.align2d**, 136
- alignment.align2d()**, 42, 117, 118, 122, 123, 133, 137, 165
- alignment.align3d()**, 79, 122, 128, 130, 133, 134, 165
- alignment.append()**, 114, 116, 118, 119, 149, 166, 180
- alignment.append_model()**, 60, 116, 166
- alignment.append_profile()**, 114, 117, 166
- alignment.check()**, 11, 20, 30, 115, 165, 180
- alignment.clear()**, 115
- alignment.compare_sequences()**, 46, 120, 166
- alignment.compare_structures()**, 22, 42, 46, 127, 165
- alignment.compare_with()**, 116, 117, 122, 165
- alignment.consensus()**, 117, 118, 126, 165
- alignment.describe()**, 119, 165
- alignment.edit()**, 118, 165
- alignment.expand()**, 141, 165, 181
- alignment.id_table()**, 46, 119, 121, 165
- alignment.malign()**, 122, 125, 130, 133, 134, 153, 166
- alignment.malign3d()**, 42, 63, 122, 128, 129, 133, 134, 166, 181
- alignment.salign**, 136
- alignment.salign()**, 21, 131, 133, 134, 140, 166
- alignment.segment_matching()**, 21, 142, 166
- alignment.to_profile()**, 141, 145, 146, 149, 165
- alignment.write()**, 60, 117, 122, 124, 126, 130, 149, 153, 167
- alignment[x].atom_file, 57, 111, 117, 130
- alignment[x].code, 111, 117, 144
- alignment_features, 117, 152
- alignment_format, 81, 114, 117
- alignment_type, 131, 135, 137
- allhmodel()**, 38
- aln, 38, 55–58, 63–65, 76–79, 98, 99, 101, 116, 152, 153
- aln_base_filename, 147
- alnfile, 35, 38, 39
- angular_step_size, 157
- asgl_output, 84, 127
- assess_methods, 26, 35, 36, 39
- atom classes, 181
- atom_accessibility, 76
- atom_files, 114, 117
- atom_ids, 100, 106, 107
- atom_types, 66, 67
- auto_overhang, 131, 137
- automodel()**, 35, 39
- automodel.assess_methods, 39
- automodel.auto_align()**, 30, 37
- automodel.cluster()**, 38
- automodel.deviation, 180
- automodel.ending_model, 26, 37, 40, 91
- automodel.final_malign3d, 37, 181
- automodel.generate_method, 29, 36
- automodel.initial_malign3d, 37
- automodel.make()**, 26, 30, 36, 37
- automodel.max_molpdf, 37
- automodel.max_var_iterations, 36
- automodel.md_level, 36, 39
- automodel.outputs, 36, 40
- automodel.rand_method, 36
- automodel.repeat_optimization, 36
- automodel.select_atoms()**, 37
- automodel.special_patches()**, 14, 38, 180
- automodel.special_restraints()**, 38, 180
- automodel.starting_model, 26, 37, 40
- automodel.user_after_single_model()**, 38
- automodel.very_fast()**, 26, 37
- automodel.write_intermediates, 37
- base_align_codes, 118
- basis_pdf_weight, 98, 101
- basis_relative_weight, 98, 101
- best_docked_models, 157, 158
- bibliography, 2
- BLOCK residues, 51

- breakpoint, 41
- bug reports, 7
- build_method, 61

- cap_atom_shift, 86, 87, 173
- chains_list, 151, 152
- chains_num, 157
- change, 71
- check_profile, 148, 149
- chop_nonstd_terminii, 81
- classes, 110
- clean_sequences, 141, 151
- close_file, 114
- cluster_cut, 38, 46, 63
- cluster_method, 63
- command, 46
- compare_mode, 127
- comparison_type, 131
- contact_shell, 96
- csrfile, 35, 36, 39
- current_directory, 130, 131, 141

- data_file, 152, 153
- debug_function_cutoff, 89, 90
- dendrogram_file, 131, 137, 140
- density()**, 156
- density.grid_search()**, 21, 156, 165
- density.read()**, 156
- density_type, 156
- detailed_debugging, 89
- deviation, 35, 38, 39, 70–72
- dih_lib_only, 98, 101
- dihedrals, 71, 72
- distance_atoms, 127
- distance_rsr_model, 98, 101

- edat, 45, 68, 74, 84, 86, 87, 89, 98, 108, 127
- edit_align_codes, 118
- edit_file_ext, 129, 130
- em_density_format, 156
- em_fit_output_file, 157
- em_map_size, 156
- em_pdb_name, 156, 157
- end_of_file, 114
- energy_data()**, 47
- energy_data.contact_shell, 47, 87, 100, 175
- energy_data.coulomb_switch, 47, 48
- energy_data.covalent_cys, 48
- energy_data.dynamic_coulomb, 47, 48, 100
- energy_data.dynamic_lennard, 47, 48, 100
- energy_data.dynamic_modeller, 48, 110
- energy_data.dynamic_sphere, 47, 100
- energy_data.excl_local, 48, 101
- energy_data.lennard_jones_switch, 47, 48
- energy_data.nlogn_use, 49, 88
- energy_data.nonbonded_sel_atoms, 48, 49, 99–101, 104
- energy_data.radii_factor, 47, 48, 74
- energy_data.relative_dielectric, 48, 178
- energy_data.update_dynamic, 47, 87
- env, 35, 38, 39, 55, 110, 113, 144, 151, 156
- env.io.hetatm, 26
- env.io.hydrogen, 26, 27
- env.io.water, 26
- environ()**, 18, 26, 45, 166
- environ.dendrogram()**, 46, 63, 120, 127, 153, 165
- environ.edat, 45, 47
- environ.io, 45, 50
- environ.libs, 46, 51
- environ.principal_components()**, 46, 120, 127, 166
- environ.system()**, 46, 166
- equilibrate, 86, 173
- equilibration, 87
- errors, 160
- exit_stage, 37, 38
- expand_control, 141, 142
- ext_tree, 131, 139
- extend_hot_spot, 69

- fast_search, 152, 153
- fast_search_cutoff, 152, 153
- feature_weight, 134
- feature_weights, 131, 140
- file, 46, 52–56, 74, 76, 81, 83, 84, 88, 109, 110, 114, 117, 142–145, 156, 161
- file naming, 41
- file types, 42
- file_exists, 161
- file_ext, 141–143, 161
- file_id, 141–143, 161
- fit, 78, 79, 118, 123, 124, 127–131, 134, 141
- fit_atoms, 127–131, 133, 134
- fit_on_first, 131, 134, 141
- fit_pdbnam, 131, 134
- fix_offset, 138
- fix_offsets, 131, 138
- frequently asked questions, 13

- gap_extension, 66
- gap_function, 132, 133, 135, 137
- gap_gap_score, 132, 140
- gap_penalties_1d, 121–123, 125, 126, 132, 140, 144, 145, 147, 148, 152–154
- gap_penalties_2d, 123, 132, 136
- gap_penalties_3d, 128–130, 132, 134
- gap_penalty_1d, 140
- gap_penalty_2d, 140
- gap_penalty_3d, 140
- gap_residue_score, 132, 140
- gaps_in_target, 148, 149
- grid_unit, 74
- group_restraints()**, 110, 166

- `group_restraints.append()`, 110, 166
- `id1`, 142, 161
- `id2`, 142, 161
- `improve_alignment`, 132, 134, 141
- `info.build_date`, 159
- `info.debug`, 159
- `info.exe_type`, 159
- `info.time_mark()`, 159, 166
- `info.version`, 159
- `info.version_info`, 159
- `infile`, 35, 36, 39
- `inimodel`, 39
- `init_velocities`, 86, 87
- `initialize_xyz`, 61
- `input_profile_file`, 123
- `input_weights_file`, 121, 123, 126, 132, 135, 139
- `installation`, 4
- `intersegment`, 98, 100
- `io`, 45, 55–58, 63, 98, 114, 115, 118, 119, 122, 127–129, 131, 152
- `io_data()`, 50
- `io_data.atom_files_directory`, 26, 42, 50
- `io_data.hetatm`, 50, 52
- `io_data.hydrogen`, 50
- `io_data.water`, 50
- `knowns`, 35, 38, 39
- Lennard-Jones, 18, 47, 48, 51, 52, 57, 87, 91, 95, 97, 99, 100, 105, 171, 172, 178
- `libraries.parameters`, 52
- `libraries.topology`, 52
- `library_schedule`, 35, 38, 39, 82
- `local_alignment`, 121, 123, 125, 128, 129, 132, 133, 152
- `log.level()`, 160, 167
- `log.minimal()`, 160
- `log.verbose()`, 45, 140, 160
- `log.very_verbose()`, 160
- `loop modeling`, 31, 181
- `loop_assess_methods`, 39
- `loopmodel()`, 39
- `loopmodel.loop.assess_methods`, 36
- `loopmodel.loop.ending_model`, 32, 40, 181
- `loopmodel.loop.md_level`, 39
- `loopmodel.loop.outputs`, 40
- `loopmodel.loop.starting_model`, 32, 40, 181
- `loopmodel.select_loop_atoms()`, 40, 181
- `matrix_comparison`, 132, 135
- `matrix_file`, 46, 119, 120, 127
- `matrix_offset`, 121–123, 125, 128, 129, 132, 144, 145, 147, 148, 152, 154
- `max_aln_evalue`, 147, 148
- `max_diff_res`, 154, 155
- `max_gap_length`, 121, 123, 126, 132
- `max_gaps_match`, 120, 121
- `max_iterations`, 36, 86, 87
- `max_nonstdres`, 81
- `max_unaligned_res`, 155
- `maximal_distance`, 98, 101
- `md_return`, 86, 87
- `md_time_step`, 86, 87
- `mdl`, 64, 78, 79, 107, 116, 117, 120
- `memory`, 160
- `method`, 8
- `min_atom_shift`, 86, 87
- `min_base_entries`, 118
- `min_loop_length`, 142, 143
- `minimal_chain_length`, 81
- `minimal_resolution`, 81
- `minimal_stdres`, 81
- `minmax_db_seq_len`, 151
- `minmax_loop_length`, 66
- `missing parameters`, 18, 51
- `mnch_lib`, 98
- `model()`, 55
- `model.assess_dope()`, 21, 36, 91
- `model.assess_ga341()`, 21, 36, 55, 90, 91
- `model.build()`, 51, 57, 61, 64, 165, 180
- `model.color()`, 77, 165
- `model.debug_function()`, 89, 165
- `model.energy()`, 14, 20, 57, 84, 87, 89, 91, 92, 102, 165, 167, 172, 181
- `model.generate_topology()`, 19, 55–58, 72, 99, 165, 180
- `model.make_chains()`, 81, 165
- `model.make_region()`, 76, 77, 165
- `model.mutate()`, 60, 66, 166
- `model.optimize()`, 36, 57, 69, 85–87, 89, 100, 101, 104, 166, 172, 180, 181
- `model.orient()`, 72, 74, 166
- `model.patch()`, 13, 14, 19, 52, 53, 57, 166
- `model.patch_ss()`, 13, 59, 166
- `model.patch_ss_templates()`, 13, 14, 58, 60, 166, 180
- `model.pick_atoms()`, 15, 48, 51, 60, 65, 67, 79, 101, 166
- `model.pick_hot_atoms()`, 48, 57, 68, 87, 95, 166, 172
- `model.randomize_xyz()`, 66, 70, 82, 166, 180
- `model.read()`, 51, 55–57, 59, 61, 72, 81, 105, 166, 180
- `model.rename_segments()`, 65, 166
- `model.reorder_atoms()`, 71, 166
- `model.res_num_from()`, 64, 166
- `model.restraints`, 55
- `model.restraints.add()`, 14, 15, 100, 106, 107, 165
- `model.restraints.append()`, 14, 109, 166, 180
- `model.restraints.clear()`, 104, 109
- `model.restraints.condense()`, 15, 101, 104, 105, 107, 165

- model.restraints.make()**, 14, 15, 18, 19, 42, 51, 58, 66, 87, 92, 98, 104, 109, 165, 180
- model.restraints.pick()**, 14, 21, 66, 69, 82, 83, 102, 104, 166, 180
- model.restraints.reindex()**, 107, 166
- model.restraints.spline()**, 108, 166
- model.restraints.unpick()**, 15, 107, 165
- model.restraints.unpick_all()**, 104
- model.restraints.write()**, 109, 167, 180
- model.rotate()**, 73, 166
- model.rotate_dihedrals()**, 66, 71, 166
- model.schedule.make()**, 35, 82, 83, 104, 165, 180
- model.schedule.read()**, 82, 83, 166
- model.schedule.step, 83, 84, 87
- model.schedule.write()**, 82, 83, 167
- model.seq_id, 55, 90
- model.superpose()**, 21, 78, 127, 166
- model.switch_trace()**, 87, 88, 166
- model.symmetry.define()**, 102, 165
- model.to_iupac()**, 55, 70, 165
- model.transfer_xyz()**, 38, 51, 55, 57, 61, 63, 130, 166, 180
- model.unbuild()**, 63, 166
- model.write()**, 56, 167, 180, 181
- model.write_data()**, 48, 74, 77, 167
- model.write_pdb_xref()**, 76, 167
- model.format, 55, 56
- model.segment, 55, 56, 76, 112
- modfile.default()**, 21, 142, 161
- modfile.delete()**, 161, 165
- modfile.inquire()**, 161, 165
- molpdf, 84, 87, 88, 90, 91

- n_prof_iterations, 144, 148, 149
- n_subopt, 132, 138
- neighbor_cutoff, 74
- no_ter, 56, 132
- normalize_pp_scores, 132, 140
- normalize_profile, 84, 85
- notes, 160
- num_structures, 156
- number_of_steps, 74, 157

- off_diagonal, 121, 123, 125, 128, 129, 132, 152
- optimization_method, 86, 87
- orient, 74
- output, 74, 84–87, 102, 108, 127–130, 132, 140, 152, 153, 160
- output_alignments, 147
- output_cod_file, 154
- output_directory, 42, 54, 56, 76, 83, 89, 109, 117, 119, 120, 142, 143
- output_grp_file, 154
- output_profile_file, 123
- output_score_file, 147–149
- output_scores, 147–149

- output_weights_file, 121, 123, 126, 132, 140
- overhang, 118, 121, 123, 125, 128, 129, 132, 137, 152
- overhang_auto_limit, 132, 137
- overhang_factor, 132, 137

- parameters, 110
- parameters.append()**, 52, 53
- parameters.clear()**, 53
- parameters.read()**, 53, 166
- parlib, 35, 39
- patch_default, 56, 58
- pick_atoms_set, 65, 66
- pick_hot_cutoff, 69
- prf, 117
- probe_radius, 74
- profile()**, 144
- profile.build()**, 21, 144, 145, 147, 148, 151, 152, 165
- profile.read()**, 141, 144, 146, 147, 149, 151, 166
- profile.scan()**, 21, 144, 146, 166
- profile.to_alignment()**, 117, 146, 149, 166
- profile.write()**, 145, 149, 167
- profile_format, 144, 145, 147
- profile_list_file, 146, 147
- program distribution, 3
- program updates, 21
- psa_integration_step, 74
- pssm_weights_type, 147

- rand_seed, 45
- read_profile, 123
- read_weights, 121, 123, 126
- reference_atom, 78, 79
- reference_distance, 78, 79
- refine_local, 21, 78, 127
- region_size, 76
- remove_gaps, 114
- renumber_residues, 65
- res_types, 66, 67
- residue_grouping, 98
- residue_ids, 19, 57, 98, 100
- residue_span_range, 21, 48, 69, 84, 86, 87, 98–101, 104, 108
- residue_span_sign, 87, 98–101
- residue_type, 57, 60
- resolution, 156
- restraint_group, 98
- restraint_parameters, 98, 101, 106
- restraint_sel_atoms, 98, 99, 104
- restraint_stddev, 98, 101
- restraint_stddev2, 98, 101
- restraint_type, 51, 98, 99, 105
- restraints_filter, 104, 105
- restraints_format, 109
- restyp_lib_file, 21, 45
- rewind_file, 114
- rms_cutoffs, 78, 79, 127, 132, 140

- root_name, 141–143, 161
- rotation_angle, 73
- rotation_axis, 73
- rotation_matrix, 73
- rr_file, 120–123, 125, 132, 133, 140, 142–145, 147, 148, 152–154

- schedule_scale, 82–84, 91
- score_statistics, 147
- script file, 11
- sdb, 148
- search_group_list, 152, 153
- search_randomizations, 152, 153
- search_sort, 152, 153
- search_top_list, 152, 153
- segment_cutoff, 142, 143
- segment_growth_c, 142, 143
- segment_growth_n, 142, 143
- segment_ids, 65
- segment_report, 142, 143
- segment_shift, 143
- segment_shifts, 142
- selection_from, 66, 67
- selection_mode, 66, 67, 69
- selection_search, 66
- selection_segment, 66, 67, 112
- selection_slab, 66, 67
- selection_status, 66
- seq_database_file, 151–153
- seq_database_format, 151, 152, 154
- seqid_cut, 154, 155
- sequence, 35, 38, 39, 56, 57, 72
- sequence_db()**, 151
- sequence_db.filter()**, 151, 154, 166
- sequence_db.read()**, 145, 151, 153, 154, 166
- sequence_db.search()**, 22, 151, 152, 166
- sequence_db.write()**, 152, 167
- sheet_h_bonds, 99, 100
- signif_cutoff, 152, 153
- similarity_flag, 132, 133, 140
- smooth_prof_weight, 121, 123, 132
- smoothing_window, 84, 85
- sphere_center, 66, 67
- sphere_radius, 66
- spline_dx, 99, 108, 109
- spline_min_points, 99, 108, 109
- spline_on_site, 98
- spline_range, 99, 108
- spline_select, 108
- start_type, 157
- statistical potential, 181
- structure_types, 81
- submodel, 53
- subopt_offset, 123, 132, 138
- substitution, 132
- summary_file, 147

- superpose_refine, 78, 79
- surftyp, 74
- swap_atoms_in_res, 78, 79
- symmetry_weight, 102

- temperature, 86, 87, 157
- toplib, 35, 39
- topology.append()**, 52, 72, 166
- topology.clear()**, 52
- topology.make()**, 22, 52–54, 165
- topology.read()**, 52, 54, 166
- topology.submodel, 17, 21, 35, 48, 52–54, 57, 60, 74
- topology.write()**, 54, 167
- topology_model, 35, 39
- trace_output, 86, 87, 89
- translate_type, 157
- translation, 73
- tutorial, 10

- variability_file, 120, 121
- viol_report_cut, 69, 84, 85
- viol_report_cut2, 84, 85
- voxel_size, 156

- warnings, 160
- weigh_sequences, 121, 123
- weights_type, 133, 135
- write_all_atoms, 56
- write_fit, 129, 130, 133, 134, 140
- write_profile, 123
- write_summary, 147
- write_weights, 121, 123, 126, 133, 140
- write_whole_pdb, 130, 133, 141