

3. CIF DATA DEFINITION AND CLASSIFICATION

```

_refine_analyze.Luzzati_coordinate_error_free
_refine_analyze.Luzzati_coordinate_error_obs
_refine_analyze.Luzzati_d_res_low_free
_refine_analyze.Luzzati_d_res_low_obs
_refine_analyze.Luzzati_sigma_a_free
_refine_analyze.Luzzati_sigma_a_free_details
_refine_analyze.Luzzati_sigma_a_obs
_refine_analyze.Luzzati_sigma_a_obs_details
_refine_analyze.number_disordered_residues
_refine_analyze.occupancy_sum_hydrogen
_refine_analyze.occupancy_sum_non_hydrogen
_refine_analyze.RG_d_res_high
_refine_analyze.RG_d_res_low
_refine_analyze.RG_free
_refine_analyze.RG_free_work_ratio
_refine_analyze.RG_work

```

The bullet (•) indicates a category key. The arrow (→) is a reference to a parent data item.

In small-molecule crystallography, there is general agreement on the metrics that should be used to assess the quality of a structure determination, and data items in the REFINE category of the core CIF dictionary can be used to record them. For macromolecular structure determinations, no such agreement has been achieved yet and new metrics are frequently suggested as the field evolves. The REFINE_ANALYZE category can be used to record the metrics that were in common use at the time that the mmCIF dictionary was constructed; it is anticipated that new metrics will be added in future versions of the dictionary, and that some of the current metrics may fall into disuse.

Luzzati (1952) devised a method for estimating the average positional shift that would be needed in an idealized refinement to reach an *R* factor of zero by using a plot of *R* factors against resolution. For some time, macromolecular crystallographers have used a modification of this approach to assess the average positional error. Recent practice has used Luzzati plots based on the free *R* values to yield a cross-validated error estimate. Data items are provided for recording these coordinate-error estimates and the range of resolution included in the plot (Example 3.6.6.8). Related data names allow the specification of the value of σ_a used in constructing the Luzzati plot.

A general feature of introducing more parameters in the model of the structure is a reduction in the *R* factor, but the statistical significance of this is often obscured by the simultaneous reduction in the ratio of observations to parameters. Attempts to extend Hamilton's (1965) test to macromolecular structures are usually confounded by the use of restraints. Tickle *et al.* (1998) proposed the use of a Hamilton generalized *R* factor analyzed separately for reflections in the working set (those used in the refinement) and for reflections in the free set (those set aside for cross validation), and these metrics are often reported in the literature. Data items are provided for recording the Hamilton generalized *R* factor for the working and free set of reflections, and for the ratio of the two.

Other indicators of a successful refinement involve the relative order of the model. Data items are provided for recording the sum of the occupancies of the hydrogen and non-hydrogen atoms in the model. The number of disordered residues may also be recorded.

3.6.6.2.3. Restraints and refinement by shells of resolution

The data items in these categories are as follows:

```

(a) REFINE_LS_RESTR
• _refine_ls_restr.type
  _refine_ls_restr.criterion
  _refine_ls_restr.dev_ideal
  _refine_ls_restr.dev_ideal_target
  _refine_ls_restr.number
  _refine_ls_restr.rejects
  _refine_ls_restr.weight

```

Example 3.6.6.8. Aspects of the refinement of an HIV-1 protease structure (PDB 5HVP) described with data items in the REFINE_ANALYZE category.

```

loop_
_refine_analyze.entry_id                '5HVP'
_refine_analyze.Luzzati_coordinate_error_obs  0.32
_refine_analyze.Luzzati_d_res_low_obs      5.0

```

(b) REFINE_LS_RESTR_NCS

```

• _refine_ls_restr.ncs.dom_id
  → _struct.ncs.dom_id
  _refine_ls_restr.ncs.ncs_model_details
  _refine_ls_restr.ncs.rms_dev_B_iso
  _refine_ls_restr.ncs.rms_dev_position
  _refine_ls_restr.ncs.weight_B_iso
  _refine_ls_restr.ncs.weight_position

```

(c) REFINE_LS_RESTR_TYPE

```

• _refine_ls_restr.type
  → _refine_ls_restr.type
  _refine_ls_restr.type.distance_cutoff_high
  _refine_ls_restr.type.distance_cutoff_low

```

(d) REFINE_LS_SHELL

```

• _refine_ls_shell.d_res_high
• _refine_ls_shell.d_res_low
  _refine_ls_shell.number_reflns_all
  _refine_ls_shell.number_reflns_obs
  _refine_ls_shell.number_reflns_R_free
  _refine_ls_shell.number_reflns_R_work
  _refine_ls_shell.percent_reflns_obs
  _refine_ls_shell.percent_reflns_R_free
  _refine_ls_shell.R_factor_all
  _refine_ls_shell.R_factor_obs
  _refine_ls_shell.R_factor_R_free
  _refine_ls_shell.R_factor_R_free_error
  _refine_ls_shell.R_factor_R_work
  _refine_ls_shell.redundancy_reflns_all
  _refine_ls_shell.redundancy_reflns_obs
  _refine_ls_shell.wR_factor_all
  _refine_ls_shell.wR_factor_obs
  _refine_ls_shell.wR_factor_R_free
  _refine_ls_shell.wR_factor_R_work

```

(e) REFINE_LS_CLASS

```

• _refine_ls_class.code
  _refine_ls_class.d_res_high
  _refine_ls_class.d_res_low
  _refine_ls_class.R_factor_all
  _refine_ls_class.R_factor_gt
  _refine_ls_class.R_Fsqd_factor
  _refine_ls_class.R_I_factor
  _refine_ls_class.wR_factor_all

```

The bullet (•) indicates a category key. Where multiple items within a category are marked with a bullet, they must be taken together to form a compound key. The arrow (→) is a reference to a parent data item.

These categories were introduced in the mmCIF dictionary to allow a detailed description of several aspects of structure refinement to be given. Data items in the REFINE_LS_RESTR category allow geometric restraints to be specified and the deviations of restrained parameters from ideal values in the final model to be given. The type of the geometric restraints can be described in more detail using data items in the REFINE_LS_RESTR_TYPE category. Data items in the REFINE_LS_RESTR_NCS category can be used to give information about any restraints on noncrystallographic symmetry used in the refinement and the category REFINE_LS_SHELL contains data items that allow the results of refinement to be given by shells of resolution.

Data items in the REFINE_LS_RESTR category can be used to record details about the restraints applied to various classes of parameters during least-squares refinement (Example 3.6.6.9). It is clearly useful to tabulate the various classes of restraint, their deviation from ideal target values and the criteria used to reject

Example 3.6.6.9. Results of the refinement of an HIV-1 protease structure (PDB 5HVP) described with data items in the *REFINE_LS_RESTR* and *REFINE_LS_SHELL* categories.

```

loop_
_refine_ls_restr.type
_refine_ls_restr.dev_ideal_target
_refine_ls_restr.dev_ideal
_refine_ls_restr.number
_refine_ls_restr.criterion
_refine_ls_restr.rejects
'p_bond_d'      0.020  0.018  1654  '>2 sigma'  22
'p_angle_d'    0.030  0.038  2246  '>2 sigma'  139
'p_planar_d'   0.040  0.043  498   '>2 sigma'  21
'p_planar'     0.020  0.015  270   '>2 sigma'  1
'p_chiral'     0.150  0.177  278   '>2 sigma'  2
'p_singtor_nbd' 0.500  0.216  582   '>2 sigma'  0
'p_multtor_nbd' 0.500  0.207  419   '>2 sigma'  0
'p_xyhbond_nbd' 0.500  0.245  149   '>2 sigma'  0
'p_planar_tor' 3.0    2.6    203   '>2 sigma'  9
'p_staggered_tor' 15.0  17.4  298   '>2 sigma'  31
'p_orthonormal_tor' 20.0  18.1  12    '>2 sigma'  1

loop_
_refine_ls_shell.d_res_low
_refine_ls_shell.d_res_high
_refine_ls_shell.number_reflns_obs
_refine_ls_shell.R_factor_obs
  8.00  4.51  1226  0.196
  4.51  3.48  1679  0.146
  3.48  2.94  2014  0.160
  2.94  2.59  2147  0.182
  2.59  2.34  2127  0.193
  2.34  2.15  2061  0.203
  2.15  2.00  1647  0.188

```

parameters that lie too far from a target, as these data are often published as part of a description of the refinement and are often deposited with the coordinates in an archive. However, the types of restraints applied depend strongly on the software package used, and as new refinement packages regularly become available, it was clearly not advisable to provide program-specific data items in the mmCIF dictionary. The approach taken in the mmCIF dictionary has been to allow the value of *_refine_ls_restr.type* to be a free-text field, so that arbitrary labels can be given to restraints that are particular to a software package, but to recommend the use of specific labels for restraints applied by particular programs. The dictionary provides examples for labels specific to the programs *PROTIN/PROLSQ* (Hendrickson & Konnert, 1979) and *RESTRAIN* (Driessen *et al.*, 1989). These program-specific representations have particular prefixes; thus the value *p_bond_d* is a bond-distance restraint as applied by *PROTIN/PROLSQ*. Values for *_refine_ls_restr.type* appropriate for other refinement programs may be suggested in future versions of the mmCIF dictionary.

Data items in the *REFINE_LS_RESTR_TYPE* category can be used to specify the ranges within which quantities are allowed to vary for each type of restraint. The special value indicated by a full stop (.) represents a restraint unbounded on the high or low side.

Data items in the *REFINE_LS_RESTR_NCS* category can be used to record details about the restraints applied to atom positions in domains related by noncrystallographic symmetry during least-squares refinement, and also to record the deviation of the restrained atomic parameters at the end of the refinement. The domains related by noncrystallographic symmetry are defined in the *STRUCT_NCS_DOM* and related categories (see Section 3.6.7.5.5). The quantities that can be recorded for each restrained domain are the root-mean-square deviations of the displacement and positional parameters, and the weighting coefficients used in

the noncrystallographic restraint of each type of parameter. Any special aspects of the way the restraints were applied may be described using *_refine_ls_restr_ncs.ncs_model_details*.

Data items in the *REFINE_LS_SHELL* category are used to summarize details of the results of the least-squares refinement by shells of resolution (Example 3.6.6.9). The resolution range, in ångströms, forms the category key; for each shell the quantities reported, such as the number of reflections above the threshold for counting as significantly intense, are all defined in the same way as the corresponding data items used to describe the results of the overall refinement in the *REFINE* category.

The core dictionary category *REFINE_LS_CLASS* was introduced after the release of the first version of the mmCIF dictionary. It provides a more general way of describing the treatment of particular subsets of the observations, but it is not expected to be used in macromolecular structural studies, where partition by shells of resolution is traditional.

3.6.6.2.4. Equivalent atoms in the refinement

The data items in these categories are as follows:

(a) *REFINE_B_ISO*

- *_refine_b_iso.class*
- *_refine_b_iso.details*
- *_refine_b_iso.treatment*
- *_refine_b_iso.value*

(b) *REFINE_OCCUPANCY*

- *_refine_occupancy.class*
- *_refine_occupancy.details*
- *_refine_occupancy.treatment*
- *_refine_occupancy.value*

The bullet (•) indicates a category key.

In macromolecular structure refinement, displacement factors or occupancies are often treated as equivalent for groups of atoms. An example would be the case where most of the atoms in the structure are refined with isotropic displacement factors, but a bound metal atom is allowed to refine anisotropically. Another example would be where the occupancies for all of the atoms in the protein part of a macromolecular complex are fixed at 1.0, but the occupancies of atoms in a bound inhibitor are refined. The *REFINE_B_ISO* and *REFINE_OCCUPANCY* categories can be used to record this information (Example 3.6.6.10).

Example 3.6.6.10. The handling of displacement factors and occupancies during the refinement of an HIV-1 protease structure (PDB 5HVP) described with data items in the *REFINE_B_ISO* and *REFINE_OCCUPANCY* categories.

```

loop_
_refine_b_iso.class
_refine_b_iso.treatment
'protein'      isotropic
'solvent'      isotropic
'inhibitor'    isotropic

loop_
_refine_occupancy.class
_refine_occupancy.treatment
_refine_occupancy.value
_refine_occupancy.details
'protein'      fix 1.00 .
'solvent'      fix 1.00 .
'inhibitor orientation 1' fix 0.65 .
'inhibitor orientation 2' fix 0.35
; The inhibitor binds to the enzyme in two
alternative conformations. The occupancy of
each conformation was adjusted so as to result
in approximately equal mean thermal factors
for the atoms in each conformation.
;

```