

3. CIF DATA DEFINITION AND CLASSIFICATION

angles associated with each phasing set, the wavelength of the radiation used in the experiment, the source of the radiation, the detector type, and the ambient temperature.

Data items in the PHASING_SET_REFLN category are used to record the values of the measured structure factors and their uncertainties. Several distinct data sets may be present in this list, with reflections in each set identified by the appropriate value of `_phasing_set_refl.n.set_id`.

3.6.6.2. Refinement

The categories describing refinement are as follows:

REFINE group

Overall description of the refinement (§3.6.6.2.1)

REFINE

REFINE_FUNCT_MINIMIZED

Analysis of the refined structure (§3.6.6.2.2)

REFINE_ANALYZE

Restraints and refinement by shells of resolution (§3.6.6.2.3)

REFINE_LS_RESTR

REFINE_LS_RESTR_NCS

REFINE_LS_RESTR_TYPE

REFINE_LS_SHELL

REFINE_LS_CLASS

Equivalent atoms in the refinement (§3.6.6.2.4)

REFINE_B_ISO

REFINE_OCCUPANCY

History of the refinement (§3.6.6.2.5)

REFINE_HIST

The macromolecular CIF dictionary contains many more data items for describing the refinement process than the core CIF dictionary does. In addition to new items in the REFINE category itself, additional categories have been introduced to describe in great detail the function minimized and the restraints applied, and the history of the refinement process, which often has many cycles. The REFINE_ANALYZE category can be used to give details of many of the quantities that may be used to assess the quality of the refinement. The REFINE_LS_SHELL category allows results to be reported by shells of resolution, and in effect replaces the more general core CIF category REFINE_LS_CLASS.

3.6.6.2.1. Overall description of the refinement

The data items in these categories are as follows:

(a) REFINE

- `_refine.entry_id`
→ `_entry.id`
- `_refine.aniso_B[1][1]`
- `_refine.aniso_B[1][2]`
- `_refine.aniso_B[1][3]`
- `_refine.aniso_B[2][2]`
- `_refine.aniso_B[2][3]`
- `_refine.aniso_B[3][3]`
- `_refine.B_iso_max`
- `_refine.B_iso_mean`
- `_refine.B_iso_min`
- `_refine.correlation_coeff_Fo_to_Fc`
- `_refine.correlation_coeff_Fo_to_Fc_free`
- `_refine.details` (~ `_refine.special_details`)
- + `_refine.diff_density_max`
- + `_refine.diff_density_min`
- + `_refine.diff_density_rms`
- `_refine.ls_abs_structure_details`
- + `_refine.ls_abs_structure_Flack`
- + `_refine.ls_abs_structure_Rogers`
- `_refine.ls_d_res_high`
- `_refine.ls_d_res_low`
- + `_refine.ls_extinction_coef`
- `_refine.ls_extinction_expression`
- `_refine.ls_extinction_method`

- + `_refine.ls_goodness_of_fit_all`
- + `_refine.ls_goodness_of_fit_gt`
- + `_refine.ls_goodness_of_fit_obs`
- `_refine.ls_goodness_of_fit_ref`
- `_refine.ls_hydrogen_treatment`
- `_refine.ls_matrix_type`
- `_refine.ls_number_constraints`
- `_refine.ls_number_parameters`
- `_refine.ls_number_reflns_all`
- `_refine.ls_number_reflns_obs`
(~ `_refine.ls_number_reflns`)
- `_refine.ls_number_reflns_R_free`
- `_refine.ls_number_reflns_R_work`
- `_refine.ls_number_restraints`
- `_refine.ls_percent_reflns_obs`
- `_refine.ls_percent_reflns_R_free`
- `_refine.ls_R_factor_all`
- `_refine.ls_R_factor_gt`
- `_refine.ls_R_factor_obs`
- `_refine.ls_R_factor_R_free`
- `_refine.ls_R_factor_R_free_error`
- `_refine.ls_R_factor_R_free_error_details`
- `_refine.ls_R_factor_R_work`
- `_refine.ls_R_Fsqd_factor_obs`
(~ `_refine.ls_R_Fsqd_factor`)
- `_refine.ls_R_I_factor_obs` (~ `_refine.ls_R_I_factor`)
- `_refine.ls_redundancy_reflns_all`
- `_refine.ls_redundancy_reflns_obs`
- `_refine.ls_restrained_S_all`
- `_refine.ls_restrained_S_obs`
- `_refine.ls_shift_over_esd_max`
(~ `_refine.ls_shift/esd_max`)
- `_refine.ls_shift_over_esd_mean`
(~ `_refine.ls_shift/esd_mean`)
- `_refine.ls_shift_over_su_max`
(~ `_refine.ls_shift/su_max`)
- `_refine.ls_shift_over_su_max_lt`
(~ `_refine.ls_shift/su_max_lt`)
- `_refine.ls_shift_over_su_mean`
(~ `_refine.ls_shift/su_mean`)
- `_refine.ls_shift_over_su_mean_lt`
(~ `_refine.ls_shift/su_mean_lt`)
- `_refine.ls_structure_factor_coef`
- `_refine.ls_weighting_details`
- `_refine.ls_weighting_scheme`
- `_refine.ls_wR_factor_all`
- `_refine.ls_wR_factor_obs`
- `_refine.ls_wR_factor_R_free`
- `_refine.ls_wR_factor_R_work`
- `_refine.occupancy_max`
- `_refine.occupancy_min`
- `_refine.overall_FOM_free_R_set`
- `_refine.overall_FOM_work_R_set`
- `_refine.overall_SU_B`
- `_refine.overall_SU_ML`
- `_refine.overall_SU_R_Cruickshank_DPI`
- `_refine.overall_SU_R_free`
- `_refine.solvent_model_details`
- `_refine.solvent_model_param_bsol`
- `_refine.solvent_model_param_ksol`

(b) REFINE_FUNCT_MINIMIZED

- `_refine_func minimized.type`
- `_refine_func minimized.number_terms`
- `_refine_func minimized.residual`
- `_refine_func minimized.weight`

The bullet (•) indicates a category key. The arrow (→) is a reference to a parent data item. Items in italics have aliases in the core CIF dictionary formed by changing the full stop (.) to an underscore (_) except where indicated by the ~ symbol. Data items marked with a plus (+) have companion data names for the standard uncertainty in the reported value, formed by appending the string `_esd` to the data name listed.

There is already an extensive set of data names in the REFINE category of the core dictionary, and Section 3.2.3.1 should be read with the present section. The only data items discussed in this section are entries in the mmCIF dictionary that do not have a counterpart in the core CIF dictionary. Analogues of a number of *R* factors in the core CIF dictionary have been added to the mmCIF dictionary to express these same *R* factors indepen-

dently for the free and working sets of reflections. The remaining new data items have more specialized roles, which are discussed below.

The data item `_refine.entry_id` has been added to the `REFINE` category to provide the formal category key required by the DDL2 data model.

Many macromolecular structure refinements now use the statistical cross-validation technique of monitoring a ‘free’ R factor (Brünger, 1997). R_{free} is calculated the same way as the conventional least-squares R factor, but using a small subset of reflections that are not used in the refinement of the structural model. Thus R_{free} tests how well the model predicts experimental observations that are not themselves used to fit the model.

The mmCIF dictionary provides data names for R_{free} and for the complementary R_{work} values for the ‘working’ set of reflections, which are the reflections that are used in the refinement. Separate data items are provided for unweighted and weighted versions of each R factor. A fixed percentage of the total number of reflections is usually assigned to the free group, and this percentage can be specified. Further details about the method used for selecting the free reflections can be given using `_reflns.R_free_details`. The estimated error in the R_{free} value may also be given, along with the method used for determining its value.

The purposes of having a set of reflections that are not used in the refinement are to monitor the progress of the refinement and to ensure that the R factor is not being artificially reduced by the introduction of too many parameters. However, as the refinement converges, the working and free R factors both approach stable values. It is common practice, particularly in structures at high resolution, to stop monitoring R_{free} at this point and to include all the reflections in the final rounds of refinement. It is thus worth noting a distinction between `_refine.ls_R_factor_obs` and `_refine.ls_R_factor_R_work`: `_refine.ls_R_factor_obs` relates to a refinement in which all reflections more intense than a specified threshold were used, while `_refine.ls_R_factor_R_work` relates to a refinement in which a subset of the observed reflections were excluded from the refinement and were used to calculate the free R factor. The dictionary allows the use of both values if a free R factor were calculated for most of the refinement, but all of the observed reflections were used in the final rounds of refinement; the protocol for this may be explained in `_refine.details`. When a full history of the refinement is provided using data items in the `REFINE_HIST` category, it is preferable to specify a change in protocol using data items in this category.

Other data items help to provide an assessment of the quality of the refinement. The scale-independent correlation coefficient between the observed and calculated structure factors may be recorded for the reflections included in the refinement using the data item `_refine.correlation_coeff_Fo_to_Fc`. There is a similar data item for the reflections that were not included in the refinement.

Overall standard uncertainties for positional and displacement parameters can be recorded according to a number of conventions. A maximum-likelihood residual for the positional parameters can be given using `_refine.overall_SU_ML` and the corresponding value for the displacement parameters can be given using `_refine.overall_SU_B`. Diffraction-component precision indexes for the displacement parameters based on the crystallographic R factor (the Cruickshank DPI; Cruickshank, 1999) can be given using `_refine.overall_SU_R_Cruickshank_DPI`. The corresponding value for R_{free} can be given using `_refine.overall_SU_R_free`.

Example 3.6.6.7. Results of the overall refinement of an HIV-1 protease structure (PDB 5HVP) described using data items in the `REFINE` and `REFINE_FUNCT_MINIMIZED` categories.

```

_refine.entry_id           '5HVP'
_refine.ls_number_reflns_obs 12901
_refine.ls_number_restraints 6609
_refine.ls_number_parameters 7032
_refine.ls_R_factor_obs     0.176
_refine.ls_weighting_scheme  calc
_refine.ls_weighting_details
; Sigdel model of Konnert-Hendrickson:
  Sigdel: Afsig + Bfsig*(sin(theta)/lambda-1/6)
  Afsig = 22.0, Bfsig = -150.0 at the beginning
    of refinement.
  Afsig = 15.5, Bfsig = -50.0 at the end of
    refinement.
;
loop_
  _refine_func_t_minimized.type
  _refine_func_t_minimized.number_terms
  _refine_func_t_minimized.residual
  'sum(W*Delta(Amplitude)^2^'      3009   1621.3
  'sum(W*Delta(Plane+Rigid)^2^'    85     56.68
  'sum(W*Delta(Distance)^2^'      1219   163.59
  'sum(W*Delta(U-tempfactors)^2^'  1192   69.338

```

The quality of a data set used for the refinement of a macromolecular structure is often given not only in terms of the scaling residuals, but also in terms of the data redundancy (the ratio of the number of reflections measured to the number of crystallographically unique reflections). Data items are provided to express the redundancy of all reflections, as well as those that have been marked as ‘observed’ (*i.e.* exceeding the threshold for inclusion in the refinement). The percentage of the total number of reflections that are considered observed is another metric of the quality of the data set, and a data item is provided for this (`_refine.ls_percent_reflns_obs`).

The limited resolution of many macromolecular data sets makes it inappropriate to refine anisotropic displacement factors for each atom. For these low- to medium-resolution studies, an overall anisotropic displacement model may be refined. The data items `_refine.aniso_B*` are provided for recording the unique elements of the matrix that describes the refined anisotropy.

The two-parameter method for modelling the contribution of the bulk solvent to the scattering proposed by Tronrud is used in several refinement programs. The data items `_refine.solvent_model_*` can be used to record the scale and displacement factors of this model, and any special aspects of its application to the refinement.

The average phasing figure of merit can be given for the working and free reflections. Unusually high or low values of displacement factors or occupancies can be a sign of problems with the refinement, so data items are provided to record the high, low and mean values of each. Further indicators of the quality of the refinement are found in the `REFINE_ANALYZE` category (Section 3.6.6.2.2).

The data items in the `REFINE_FUNCT_MINIMIZED` category allow a brief description of the function minimized during refinement to be given (Example 3.6.6.7). It is not possible to reconstruct the function minimized during the refinement by automatic parsing of the values of these data items, but the details given in them may still be helpful to someone reading the mmCIF.

3.6.6.2.2. Analysis of the refined structure

The data items in this category are as follows:

`REFINE_ANALYZE`

- `_refine_analyze.entry_id`
→ `_entry.id`