

3.6. CLASSIFICATION AND USE OF MACROMOLECULAR DATA

STRUCT_NCS_ENS_GEN

STRUCT_NCS_OPER

External databases (§3.6.7.5.6)

STRUCT_REF

STRUCT_REF_SEQ

STRUCT_REF_SEQ_DIF

 β -sheets (§3.6.7.5.7)

STRUCT_SHEET

STRUCT_SHEET_TOPOLOGY

STRUCT_SHEET_ORDER

STRUCT_SHEET_RANGE

STRUCT_SHEET_HBOND

Molecular sites (§3.6.7.5.8)

STRUCT_SITE_GEN

STRUCT_SITE_KEYWORDS

STRUCT_SITE_VIEW

The results of the determination of a structure can be described in mmCIF using data items in the categories contained in the STRUCT category group. This is a very large group of categories and it has been divided into eight groups of related categories for the discussions that follow: (1) those that describe the structure at the level of biologically relevant assemblies; (2) those that describe the secondary structure of the macromolecules present; (3) those that describe the structural interactions that determine the conformation of the macromolecules; (4) those that describe properties of the structure at the monomer level; (5) those that describe ensembles of identical domains related by noncrystallographic symmetry; (6) those that provide references to related entities in external databases; (7) those that describe the β -sheets present in the structure; and (8) those that provide detailed descriptions of the structure of biologically interesting molecular sites.

3.6.7.5.1. Higher-level macromolecular structure

The data items in these categories are as follows:

(a) STRUCT

- `_struct.entry_id`
→ `_entry.id`
- `_struct.title`

(b) STRUCT_ASYM

- `_struct_asym.id`
- `_struct_asym.details`
- `_struct_asym.entity_id`
→ `_entity.id`

(c) STRUCT_BIOL

- `_struct_biol.id`
- `_struct_biol.details`

(d) STRUCT_BIOL_GEN

- `_struct_biol_gen.asym_id`
→ `_struct_asym.id`
- `_struct_biol_gen.biol_id`
→ `_struct_biol.id`
- `_struct_biol_gen.symmetry`
- `_struct_biol_gen.details`

(e) STRUCT_BIOL_KEYWORDS

- `_struct_biol_keywords.biol_id`
→ `_struct_biol.id`
- `_struct_biol_keywords.text`

(f) STRUCT_BIOL_VIEW

- `_struct_biol_view.biol_id`
→ `_struct_biol.id`
- `_struct_biol_view.id`
- `_struct_biol_view.details`
- `_struct_biol_view.rot_matrix[1][1]`
- `_struct_biol_view.rot_matrix[1][2]`
- `_struct_biol_view.rot_matrix[1][3]`

```
_struct_biol_view.rot_matrix[2][1]
_struct_biol_view.rot_matrix[2][2]
_struct_biol_view.rot_matrix[2][3]
_struct_biol_view.rot_matrix[3][1]
_struct_biol_view.rot_matrix[3][2]
_struct_biol_view.rot_matrix[3][3]
```

(g) STRUCT_KEYWORDS

- `_struct_keywords.entry_id`
→ `_entry.id`
- `_struct_keywords.text`

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.

The data items in these categories serve two related but distinct purposes.

The first purpose is to label each of the entities in the asymmetric unit, using data items in the STRUCT_ASYM category. These labels become part of the category key that identifies each coordinate record and they are used extensively throughout the STRUCT family of categories, so care must be taken to select a labelling scheme that is concise and informative.

The second function is descriptive. The categories descending from STRUCT_BIOL allow the author of the mmCIF to identify and annotate the biologically relevant structural units found by the structure determination. What constitutes a biological unit can depend on the context. Take the case of a structure with two polymers related by noncrystallographic symmetry, each of which binds a small-molecule cofactor. If the author wishes to describe the dimer interface, the biological unit could be taken to be the two protein molecules. If the author wishes to highlight the cofactor binding mode, the biological unit could be taken to be one protein molecule and its bound cofactor. In this second case, there could be an additional biological unit of the second protein molecule and its bound cofactor, which may or may not be identical in conformation to the first.

The relationships between categories used to describe higher-level structure are illustrated in Fig. 3.6.7.7.

The STRUCT category serves to link the structure to the overall identifier for the data block, using `_struct.entry_id`, and to supply a title that describes the entire structure. The importance of this title as a succinct description of the structure should not be underestimated, and the author should express concisely but clearly in `_struct.title` the components of interest and the importance of this particular study. It is useful to think of this title as describing

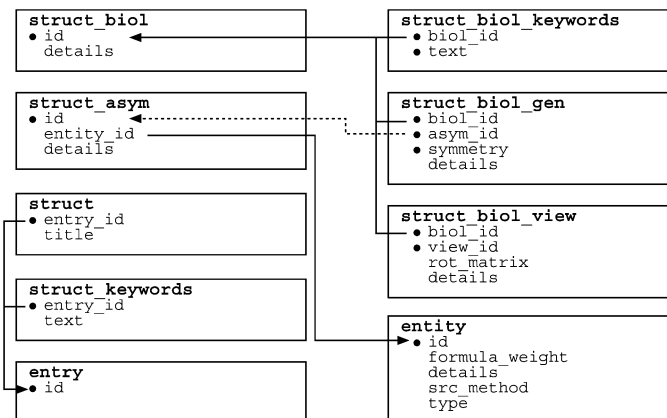


Fig. 3.6.7.7. The family of categories used to describe the higher-level macromolecular structure. Boxes surround categories of related data items. Data items that serve as category keys are preceded by a bullet (•). Lines show relationships between linked data items in different categories with arrows pointing at the parent data items.

Example 3.6.7.8. *The higher-level structure of the complex of HIV-1 protease with an inhibitor (PDB 5HVP) described with data items in the STRUCT_ASYM, STRUCT_BIOL, STRUCT_BIOL_KEYWORDS and STRUCT_BIOL_GEN categories.*

```

loop_
_struct_asy.id
_struct_asy.entity_id
_struct_asy.details
  A 1 'one monomer of the dimeric enzyme'
  B 1 'one monomer of the dimeric enzyme'
  C 2
'one partially occupied position for the inhibitor'
  D 2
'one partially occupied position for the inhibitor'

loop_
_struct_biol.id
_struct_biol.details
  1
; significant deviations from twofold symmetry exist
in this dimeric enzyme
;
  2
; The drug binds to this enzyme in two roughly
twofold symmetric modes.

Hence this biological unit (2) is roughly twofold
symmetric to biological unit (3). Disorder in the
protein chain indicated with alternative ID 1
should be used with this biological unit.
;
  3
; The drug binds to this enzyme in two roughly
twofold symmetric modes.

Hence this biological unit (3) is roughly twofold
symmetric to biological unit (2). Disorder in the
protein chain indicated with alternative ID 2
should be used with this biological unit.
;

loop_
_struct_biol_gen.biol_id
_struct_biol_gen.asy_id
_struct_biol_gen.symmetry
  1 A 1_555 1 B 1_555
  2 A 1_555 2 B 1_555 2 C 1_555
  3 A 1_555 3 B 1_555 3 D 1_555

```

the motivation for the structure determination, rather than the result. For instance, if the goal of the study was to determine the structure of enzyme A at pH 7.2 as part of a study of the mechanism of the reaction catalysed by the enzyme, an appropriate value for `_struct.title` would be 'Enzyme A at pH 7.2', even if the structure was found to contain two molecules per asymmetric unit, a bound calcium ion and a disordered loop between residues 47 and 52.

The `STRUCT_KEYWORDS` category allows an author to include keywords for the structure that has been determined. Other categories, such as `STRUCT_BIOL_KEYWORDS` and `STRUCT_SITE_KEYWORDS`, allow more specific keywords to be given, but the `STRUCT_KEYWORDS` category is the most likely category to be searched by simple information retrieval applications, so the author of an mmCIF might want to duplicate any keywords given elsewhere in the mmCIF in `STRUCT_KEYWORDS` as well.

The chemical entities that form the contents of the asymmetric unit are identified using data items in the `ENTITY` categories. The data items in the `STRUCT_ASYM` category link these entities to the structure itself. A unique identifier is attached to each occurrence of each entity in the asymmetric unit using `_struct_asy.id`. This identifier forms a part of the atom label in the `ATOM_SITE` category, which is used throughout the many categories in the `STRUCT` group

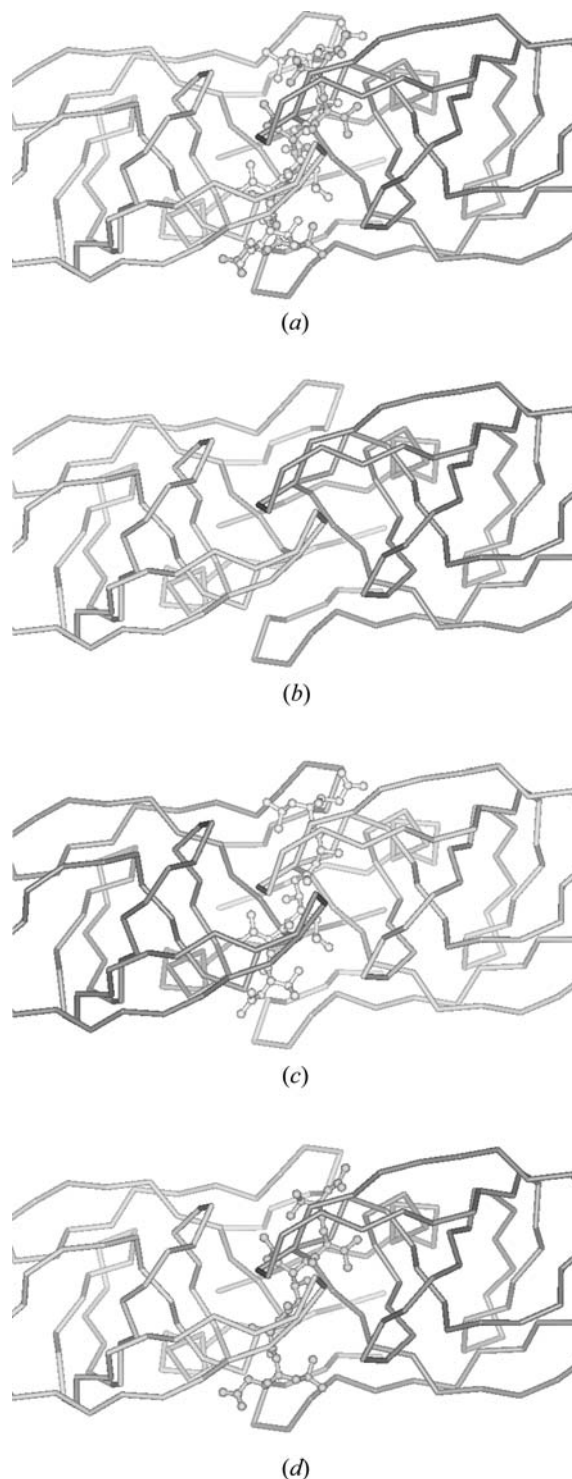


Fig. 3.6.7.8. The higher-level structure of the complex of HIV-1 protease with an inhibitor (PDB 5HVP) to be described with data items in the `STRUCT_ASYM`, `STRUCT_BIOL`, `STRUCT_BIOL_KEYWORDS` and `STRUCT_BIOL_GEN` categories. (a) Complete structure; (b), (c), (d) three different biological units.

in describing the structure. The identifier is also used in generating biological assemblies.

The usual reason for determining the structure of a biological macromolecule is to get information about the biologically relevant assemblies of the entities in the crystal structure. These assemblies take many forms and could encompass the complete contents of the asymmetric unit, a fraction of the contents of the asymmetric unit or the contents of more than one asymmetric unit. Each assembly, or 'biological unit', is given an identifier in the `STRUCT_BIOL` category and the author may annotate each biological unit using the data item `_struct_biol.details`. Key-

3.6. CLASSIFICATION AND USE OF MACROMOLECULAR DATA

words for each biological unit can be given using data items in the `STRUCT_BIOL_KEYWORD` category.

The entities that comprise the biological unit are specified using data items in the `STRUCT_BIOL_GEN` category by reference to the appropriate values of `_struct_asym.id` and by specifying any symmetry transformation that must be applied to the entities to generate the biological unit.

Data items in the `STRUCT_BIOL_VIEW` category allow the author to specify an orientation of the biological unit that provides a useful view of the structure. The comments given in `_struct_biol_view.details` may be used as a figure caption if the view is intended to be a figure in a report describing the structure.

The example of crambin in Section 3.6.3 shows the relations between the categories defining higher-level structure for the straightforward case of a single protein molecule (with a small co-crystallization molecule and solvent) in the asymmetric unit. The structure of HIV-1 protease with a bound inhibitor (PDB 5HVP), shown in Example 3.6.7.8, is considerably more complex. There are two entities: the monomeric form of the enzyme and the small-molecule inhibitor. The asymmetric unit contains two copies of the enzyme monomer (both fully occupied) and two copies of the inhibitor (each of which is partially occupied) (Fig. 3.6.7.8). Three biological assemblies are constructed for this system. One biological unit contains only the dimeric enzyme (Fig. 3.6.7.8*b*), the second contains the dimeric enzyme with one partially occupied conformation of the inhibitor (Fig. 3.6.7.8*c*) and the third contains the dimeric enzyme with the second partially occupied conformation of the inhibitor (Fig. 3.6.7.8*d*). There are alternative conformations of the side chains in the enzyme that correlate with the binding mode of the inhibitor.

3.6.7.5.2. Secondary structure

The data items in these categories are as follows:

(a) `STRUCT_CONF_TYPE`

- `_struct_conf_type.id`
- `_struct_conf_type.criteria`
- `_struct_conf_type.reference`

(b) `STRUCT_CONF`

- `_struct_conf.id`
- `_struct_conf.beg_label_asym_id`
→ `_atom_site.label_asym_id`
- `_struct_conf.beg_label_comp_id`
→ `_atom_site.label_comp_id`
- `_struct_conf.beg_label_seq_id`
→ `_atom_site.label_seq_id`
- `_struct_conf.beg_auth_asym_id`
→ `_atom_site.auth_asym_id`
- `_struct_conf.beg_auth_comp_id`
→ `_atom_site.auth_comp_id`
- `_struct_conf.beg_auth_seq_id`
→ `_atom_site.auth_seq_id`
- `_struct_conf.conf_type_id`
→ `_struct_conf_type.id`
- `_struct_conf.details`
- `_struct_conf.end_label_asym_id`
→ `_atom_site.label_asym_id`
- `_struct_conf.end_label_comp_id`
→ `_atom_site.label_comp_id`
- `_struct_conf.end_label_seq_id`
→ `_atom_site.label_seq_id`
- `_struct_conf.end_auth_asym_id`
→ `_atom_site.auth_asym_id`
- `_struct_conf.end_auth_comp_id`
→ `_atom_site.auth_comp_id`
- `_struct_conf.end_auth_seq_id`
→ `_atom_site.auth_seq_id`

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

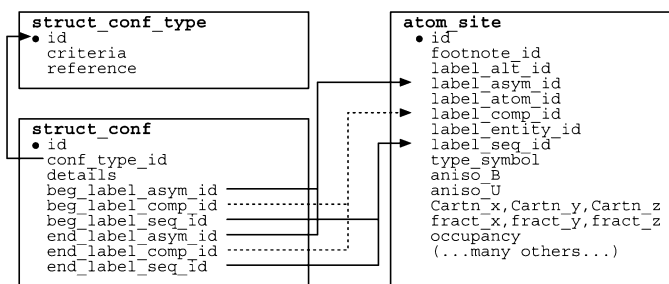


Fig. 3.6.7.9. The family of categories used to describe secondary structure. Boxes surround categories of related data items. Data items that serve as category keys are preceded by a bullet (•). Lines show relationships between linked data items in different categories with arrows pointing at the parent data items.

Example 3.6.7.9. Secondary structure in an HIV-1 protease structure (PDB 5HVP) described with data items in the `STRUCT_CONF_TYPE` and `STRUCT_CONF` categories.

```
loop_
  _struct_conf_type.id
  _struct_conf_type.criteria
  HELX_RH_AL_P 'author judgement'
  STRN         'author judgement'
  TURN_TY1_P  'author judgement'
  TURN_TY1P_P 'author judgement'
  TURN_TY2_P  'author judgement'
  TURN_TY2P_P 'author judgement'

loop_
  _struct_conf.id
  _struct_conf.conf_type_id
  _struct_conf.beg_label_comp_id
  _struct_conf.beg_label_asym_id
  _struct_conf.beg_label_seq_id
  _struct_conf.end_label_comp_id
  _struct_conf.end_label_asym_id
  _struct_conf.end_label_seq_id
  HELX1  HELX_RH_AL_P  ARG  A   87  GLN  A   92
  HELX2  HELX_RH_AL_P  ARG  B  287  GLN  B  292
  STRN1  STRN          PRO  A    1  LEU  A    5
  STRN2  STRN          CYS  B  295  PHE  B  299
  STRN3  STRN          CYS  A   95  PHE  A  299
  STRN4  STRN          PRO  B  201  LEU  B  205
  TURN1  TURN_TY1P_P  ILE  A   15  GLN  A   18
  TURN2  TURN_TY2_P   GLY  A   49  GLY  A   52
  TURN3  TURN_TY1P_P  ILE  A   55  HIS  A   69
  TURN4  TURN_TY1_P   THR  A   91  GLY  A   94
```

The primary structure of a macromolecule is defined by the sequence of the components (amino acids, nucleic acids or sugars) in the polymer chain. The polymer chains assume conformations based on the torsion angles adopted by the rotatable bonds in the polymer backbone; the resulting conformations are referred to as the secondary structure of the polymer. Several patterns of values of backbone torsion angles have been described and given names, such as α -helix, β -strand, turn and coil for proteins, and A-, B- and Z-helix for nucleic acids.

In the mmCIF dictionary, these secondary structures are described in the `STRUCT_CONF` and `STRUCT_CONF_TYPE` categories. Note that the data items in these categories describe only the secondary structure; the tertiary organization of β -strands into β -sheets is described in the `STRUCT_SHEET_*` categories. There are no data items for describing the tertiary organization of α -helices or nucleic acids in the current version of the mmCIF dictionary.

The relationships between categories used to describe secondary structure are shown in Fig. 3.6.7.9.

The type of the secondary structure is specified in the `STRUCT_CONF_TYPE` category, along with the criteria used to identify it. The range of monomers assigned to each secondary-structure element is given in the `STRUCT_CONF` category.