

## 3. CIF DATA DEFINITION AND CLASSIFICATION

Example 3.6.7.12. *Noncrystallographic symmetry in the structure of trimeric haemerythrin (PDB 1HR3) described with data items in the STRUCT\_NCS\_ENS, STRUCT\_NCS\_ENS\_GEN, STRUCT\_NCS\_DOM and STRUCT\_NCS\_DOM\_LIM categories. For brevity, the data items in the STRUCT\_NCS\_OPER category are not shown.*

```
_struct_ncs_ens.id          trimer
_struct_ncs_ens.point_group 3

loop_
_struct_ncs_ens_gen.ens_id
_struct_ncs_ens_gen.dom_id_1
_struct_ncs_ens_gen.dom_id_2
_struct_ncs_ens_gen.oper_id
trimer chain_A chain_B 1
trimer chain_A chain_C 2

loop_
_struct_ncs_dom.id
chain_A chain_B chain_C

loop_
_struct_ncs_dom_lim.dom_id
_struct_ncs_dom_lim.beg_label_asym_id
_struct_ncs_dom_lim.beg_label_comp_id
_struct_ncs_dom_lim.beg_label_seq_id
_struct_ncs_dom_lim.beg_label_alt_id
_struct_ncs_dom_lim.end_label_asym_id
_struct_ncs_dom_lim.end_label_comp_id
_struct_ncs_dom_lim.end_label_seq_id
_struct_ncs_dom_lim.end_label_alt_id
chain_A A ala 1 . A ala 118 .
chain_B B ala 1 . B ala 118 .
chain_C C ala 1 . C ala 118 .
```

specify a pair of domains and reference the appropriate operator in STRUCT\_NCS\_OPER, but now the data item `_struct_ncs_oper.code` is given the value `generate`.

Noncrystallographic symmetry in a trimeric molecule is shown in Fig. 3.6.7.13 and described in Example 3.6.7.12.

## 3.6.7.5.6. External databases

The data items in these categories are as follows:

## (a) STRUCT\_REF

- `_struct_ref.id`
- `_struct_ref.biol_id`  
→ `_struct_biol.id`
- `_struct_ref.db_code`
- `_struct_ref.db_name`
- `_struct_ref.details`
- `_struct_ref.entity_id`  
→ `_entity.id`
- `_struct_ref.seq_align`
- `_struct_ref.seq_dif`

## (b) STRUCT\_REF\_SEQ

- `_struct_ref_seq.align_id`
- `_struct_ref_seq.db_align_beg`
- `_struct_ref_seq.db_align_end`
- `_struct_ref_seq.details`
- `_struct_ref_seq.ref_id`  
→ `_struct_ref.id`
- `_struct_ref_seq.seq_align_beg`  
→ `_entity_poly_seq.num`
- `_struct_ref_seq.seq_align_end`  
→ `_entity_poly_seq.num`

## (c) STRUCT\_REF\_SEQ\_DIF

- `_struct_ref_seq_dif.align_id`  
→ `_struct_ref_seq.align_id`
- `_struct_ref_seq_dif.seq_num`  
→ `_entity_poly_seq.num`
- `_struct_ref_seq_dif.db_mon_id`  
→ `_chem_comp.id`
- `_struct_ref_seq_dif.details`

```
_struct_ref_seq_dif.mon_id
→ _chem_comp.id
```

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.

Data items in the STRUCT\_REF category allow the author of an mmCIF to provide references to information in external databases that is relevant to the entities or biological units described in the mmCIF. For example, the database entry for a protein or nucleic acid sequence could be referenced and any differences between the sequence of the macromolecule whose structure is reported in the mmCIF and the sequence of the related entry in the external database can be recorded. Alternatively, references to external database entries can be used to record the relationship of the structure reported in the mmCIF to structures already reported in the literature, for example by referring to previously determined structures of the same or a similar protein, or to a small-molecule structure determination of a bound inhibitor or cofactor. STRUCT\_REF data items are not intended to be used to reference a database entry for the structure in the mmCIF itself (this would be the role of data items in the DATABASE\_2 category), but it would not be formally incorrect to do so.

When the data items in these categories are used to provide references to external database entries describing the sequence of a polymer, data items from all three categories could be used. The value of the data item `_struct_ref.seq_align` is used to indicate whether the correspondence between the sequence of the entity or biological unit in the mmCIF and the sequence in the related external database entry is complete or partial. If the value is `partial`, the region (or regions) of the alignment may be identified using data items in the STRUCT\_REF\_SEQ category. Comments on the alignment may be given in `_struct_ref_seq.details` (Example 3.6.7.13).

The value of the data item `_struct_ref.seq_dif` is used to indicate whether the two sequences contain point differences. If the value is `yes`, the differences may be identified and annotated using data items in the STRUCT\_REF\_SEQ\_DIF category. Comments on specific point differences may be recorded in `_struct_ref_seq_dif.details`.

Example 3.6.7.13. *The relationship of the sequence of the protein PDB 5HVP to a sequence in an external database described with data items in the STRUCT\_REF and STRUCT\_REF\_SEQ categories.*

```
loop_
_struct_ref.id
_struct_ref.biol_id
_struct_ref.entity_id
_struct_ref.db_name
_struct_ref.db_code
_struct_ref.seq_align
_struct_ref.seq_dif
seq_pdb 1 . PDB 5HVP .
seq_genbank . 1 GenBank AAG30358 complete yes

loop_
_struct_ref_seq.align_id
_struct_ref_seq.ref_id
_struct_ref_seq.seq_align_beg
_struct_ref_seq.seq_align_end
_struct_ref_seq.db_align_beg
_struct_ref_seq.db_align_end
_struct_ref_seq.details
align_seq_pdb_genbank seq_genbank 1 99 24 122
; The genbank reference is to the sequence of
residues 1-376 of the viral pol 1 polypeptide;
the protease is proteolytically released from
this precursor during viral maturation.
;
```

References do not have to be to entries in databases of sequences: any external database can be referenced. For other kinds of databases, only the data items in the STRUCT\_REF category would usually be used. The element of the structure that is referenced could be either an entity or a biological unit, that is, either a building block of the structure or a structurally meaningful assembly of those building blocks. Since the identification of the part of the structure being linked to an entry in an external database can be made using either `_struct_ref.biol_id` or `_struct_ref.entity_id`, and since any part of the structure could be linked to any number of entries in external databases, the data item `_struct_ref.id` was introduced as the category key.

### 3.6.7.5.7. $\beta$ -sheets

Data items in these categories are as follows:

#### (a) STRUCT\_SHEET

- `_struct_sheet.id`
- `_struct_sheet.details`
- `_struct_sheet.number_strands`
- `_struct_sheet.type`

#### (b) STRUCT\_SHEET\_TOPOLOGY

- `_struct_sheet_topology.range_id_1`  
  → `_struct_sheet_range.id`
- `_struct_sheet_topology.range_id_2`  
  → `_struct_sheet_range.id`
- `_struct_sheet_topology.sheet_id`  
  → `_struct_sheet.id`
- `_struct_sheet_topology.offset`
- `_struct_sheet_topology.sense`

#### (c) STRUCT\_SHEET\_RANGE

- `_struct_sheet_range.id`
- `_struct_sheet_range.sheet_id`  
  → `_struct_sheet.id`
- `_struct_sheet_range.beg_label_asym_id`  
  → `_struct_asym.id`
- `_struct_sheet_range.beg_label_comp_id`  
  → `_chem_comp.id`
- `_struct_sheet_range.beg_label_seq_id`  
  → `_atom_site.label_seq_id`
- `_struct_sheet_range.end_label_asym_id`  
  → `_struct_asym.id`
- `_struct_sheet_range.end_label_comp_id`  
  → `_chem_comp.id`
- `_struct_sheet_range.end_label_seq_id`  
  → `_atom_site.label_seq_id`
- `_struct_sheet_range.beg_auth_asym_id`  
  → `_atom_site.auth_atom_id`
- `_struct_sheet_range.beg_auth_comp_id`  
  → `_atom_site.auth_comp_id`
- `_struct_sheet_range.beg_auth_seq_id`  
  → `_atom_site.auth_seq_id`
- `_struct_sheet_range.end_auth_asym_id`  
  → `_atom_site.auth_atom_id`
- `_struct_sheet_range.end_auth_comp_id`  
  → `_atom_site.auth_comp_id`
- `_struct_sheet_range.end_auth_seq_id`  
  → `_atom_site.auth_seq_id`
- `_struct_sheet_range.symmetry`

#### (d) STRUCT\_SHEET\_ORDER

- `_struct_sheet_order.range_id_1`  
  → `_struct_sheet_range.id`
- `_struct_sheet_order.range_id_2`  
  → `_struct_sheet_range.id`
- `_struct_sheet_order.sheet_id`  
  → `_struct_sheet.id`
- `_struct_sheet_order.offset`
- `_struct_sheet_order.sense`

#### (e) STRUCT\_SHEET\_HBOND

- `_struct_sheet_hbond.range_id_1`  
  → `_struct_sheet_range.id`
- `_struct_sheet_hbond.range_id_2`  
  → `_struct_sheet_range.id`

- `_struct_sheet_hbond.sheet_id`  
  → `_struct_sheet.id`
- `_struct_sheet_hbond.range_1_beg_label_atom_id`  
  → `_atom_site.label_atom_id`
- `_struct_sheet_hbond.range_1_beg_label_seq_id`  
  → `_atom_site.label_seq_id`
- `_struct_sheet_hbond.range_1_end_label_atom_id`  
  → `_atom_site.label_atom_id`
- `_struct_sheet_hbond.range_1_end_label_seq_id`  
  → `_atom_site.label_seq_id`
- `_struct_sheet_hbond.range_2_beg_label_atom_id`  
  → `_atom_site.label_atom_id`
- `_struct_sheet_hbond.range_2_beg_label_seq_id`  
  → `_atom_site.label_seq_id`
- `_struct_sheet_hbond.range_2_end_label_atom_id`  
  → `_atom_site.label_atom_id`
- `_struct_sheet_hbond.range_2_end_label_seq_id`  
  → `_atom_site.label_seq_id`
- `_struct_sheet_hbond.range_1_beg_auth_atom_id`  
  → `_atom_site.auth_atom_id`
- `_struct_sheet_hbond.range_1_beg_auth_seq_id`  
  → `_atom_site.auth_seq_id`
- `_struct_sheet_hbond.range_1_end_auth_atom_id`  
  → `_atom_site.auth_atom_id`
- `_struct_sheet_hbond.range_1_end_auth_seq_id`  
  → `_atom_site.auth_seq_id`
- `_struct_sheet_hbond.range_2_beg_auth_atom_id`  
  → `_atom_site.auth_atom_id`
- `_struct_sheet_hbond.range_2_beg_auth_seq_id`  
  → `_atom_site.auth_seq_id`
- `_struct_sheet_hbond.range_2_end_auth_atom_id`  
  → `_atom_site.auth_atom_id`
- `_struct_sheet_hbond.range_2_end_auth_seq_id`  
  → `_atom_site.auth_seq_id`

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.

Different methods of describing  $\beta$ -sheets are in widespread use. The mmCIF dictionary provides data items for two methods and it is anticipated that future versions of the dictionary could cover others. The model used in the STRUCT\_SHEET\_TOPOLOGY category is the simpler of the two. It is a convenient shorthand for describing the topology, but it does not provide details about strand registration and it is not suitable for describing sheets that contain strands from more than one polypeptide. A more general model is provided by the linked data items in the STRUCT\_SHEET\_RANGE, STRUCT\_SHEET\_ORDER and STRUCT\_SHEET\_HBOND categories. For both methods of representing  $\beta$ -sheets, data items in the parent category STRUCT\_SHEET can be used to provide an identifier for each sheet, a free-text description of its type, the number of participating strands and a free-text description of any peculiar aspects of the sheet. The relationships between categories used to describe  $\beta$ -sheets are shown in Fig. 3.6.7.14.

In the description of  $\beta$ -sheet topology based on the STRUCT\_SHEET\_TOPOLOGY category, the strand that occurs first in the polypeptide chain is numbered 1. Subsequent strands are described by their position in the sheet relative to the previous strand (+1, -3 *etc.*) and by their orientation relative to the previous strand (parallel or antiparallel).

While writing this chapter, a few errors in the mmCIF dictionary were discovered. The use of `_struct_sheet_topology.range_id_1` and `*_2` as pointers to the residues participating in  $\beta$ -sheets is one; the correct data items should be `_struct_sheet_topology.comp_id_1` and `*_2`, and these data items should be pointers to `_atom_site.label_comp_id`. This error will be corrected in future versions of the dictionary. As the data model encoded in the current version of the dictionary is incorrect, no example of its use is given.