

3.6. CLASSIFICATION AND USE OF MACROMOLECULAR DATA

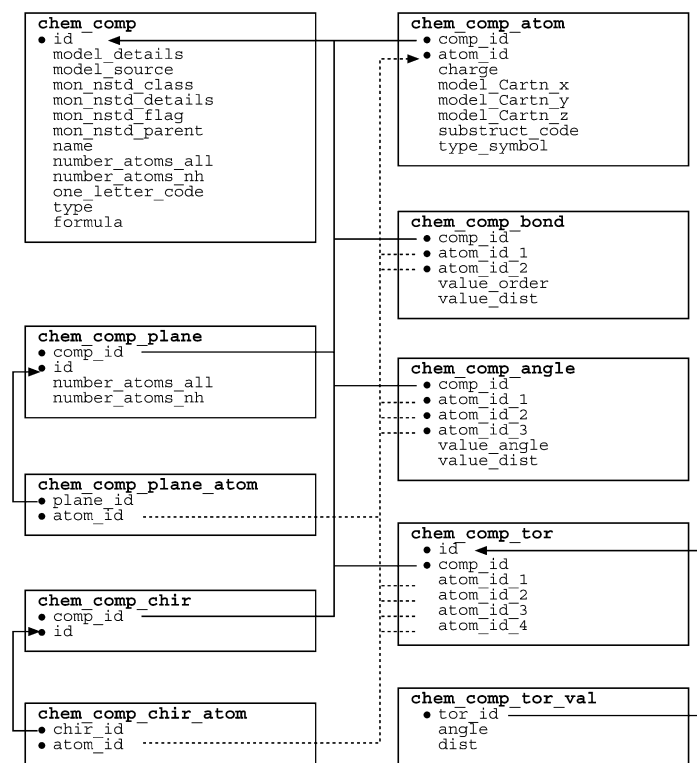


Fig. 3.6.7.3. The family of categories used to describe the chemical and structural features of the monomers and small molecules used to build a model of a structure. Boxes surround categories of related data items. Data items that serve as category keys are preceded by a bullet (•). Lines show relationships between linked data items in different categories with arrows pointing at the parent data items.

of atoms, the number of non-hydrogen atoms, and the name of the component. The name of the component will typically be a common name such as ‘alanine’ or ‘valine’; it is recommended that the IUPAC name is used for components that are not among the usual monomers that make up proteins, nucleic acids or sugars.

The one-letter or three-letter code for a standard component may be given (using `_chem_comp.one_letter_code` and `_chem_comp.three_letter_code`, respectively). Values of `x` for the one-letter code or `UNK` for the three-letter code are used to indicate components that do not have a standard abbreviation. A component that has been formed by modification of a standard component can be indicated by prefixing the code with a plus sign. A value of ‘.’, which means ‘not applicable’, should be used for components that are not monomers from which a polymeric macromolecule is built, for example co-crystallized small molecules, ions or water.

The data item `_chem_comp.type` can be used to describe the structural role of a monomer within a polymeric molecule. The types that are recognized are classified as linking monomers (for proteins, nucleic acids and sugars), monomers with an N-terminal or C-terminal cap (for proteins), and monomers with a 5′ or 3′ terminal cap (for nucleic acids). The specification of types for sugars is less complete than for proteins and nucleic acids and no types of terminal groups are currently specified for sugars. The values `non-polymer` and `other` are provided for types that have not been defined explicitly.

Information about the source of the model for the chemical component can be given using `_chem_comp.model_source` and `_chem_comp.model_details`. `_chem_comp.model_source` is a text field where the user might, for example, supply a reference to the Cambridge Structural Database or another small-molecule crystallographic database, or describe a molecular-modelling process. `_chem_comp.model_details` can be used to discuss any modification made to the model given in `_chem_comp.model_source`.

As mentioned previously, `_chem_comp.model_errf` can be used to specify the location of an external reference file if the model is not described within the current data block.

Macromolecules often contain modifications of standard monomers, such as phosphorylated serines and threonines. In the mmCIF data model, a nonstandard monomer should be treated as a separate `CHEM_COMP` entry and described in full. However, it may be useful to refer to the standard monomer from which it was derived using the `_chem_comp.mon_nstd_*` data items. There are no fixed rules for what constitutes a ‘standard’ or ‘nonstandard’ monomer in this context, but any covalent modification of a standard amino acid or nucleotide would generally be considered nonstandard. Sometimes it is difficult to decide whether a monomer is standard or nonstandard: selenomethionine is not one of the standard 20 amino acids, but it is so commonly used that geometric restraints for it are included in many standard packages for protein structure refinement.

Data items in the `CHEM_COMP_ATOM` category can be used to describe the atoms in a component. The position of each atom is given in orthogonal ångström coordinates. These coordinates correspond to the atom positions in the model of the component used in the refinement, not to the final set of refined atom positions recorded in the `ATOM_SITE` list.

Other `CHEM_COMP_ATOM` data items can be used to specify what element the atom is and its formal electronic charge, or partial charge. A code may also be assigned to the atom to indicate its role within a substructural classification of the component. The allowed codes are `main` and `side` for the main-chain and side-chain parts of amino acids, and `base`, `phos` and `sugar` for the base, phosphate and sugar parts of nucleotides. Atoms that do not belong to a substructure may be assigned the code `none`.

Data items in the `CHEM_COMP_BOND` category can be used to describe the intramolecular bonds between atoms in a component. Bond restraints may be described by the distance between the bonded atoms, the bond order, or both. The recognized bond types are the same as those for the core CIF dictionary data item `_chemical_conn_bond.type`, and they fulfil the same role: to characterize a model that could be used for database substructure searching, rather than to give a detailed description of unusual bond types.

In the `CHEM_COMP_ANGLE` category, atom 2 defines the vertex of the angle involving atoms 1, 2 and 3. The angle may be described as either an angle at the vertex atom or as a distance between atoms 1 and 3.

Data items in the `CHEM_COMP_CHIR` category can be used to describe the conformation of chiral centres within the component. The absolute configuration and the chiral volume may be specified, as well as the total number of atoms and the number of non-hydrogen atoms bonded to the chiral centre. There is also a flag to indicate whether a restrained chiral volume should match the target value in sign as well as in magnitude. Because chiral centres can involve a variable number of atoms, a separate list of the atoms should be given in `CHEM_COMP_CHIR_ATOM`.

Data items in the `CHEM_COMP_PLANE` category can be used to define planes within a component. The number of non-hydrogen atoms and the total number of atoms in each plane can be recorded. The atoms defining each plane should be listed separately in `CHEM_COMP_PLANE_ATOM`.

Data items in the `CHEM_COMP_TOR` category can be used to give details about the torsion angles in a component. A torsion angle may be described either as an angle or as a distance between the first and last atoms. (A torsion angle cannot be completely described by a distance, but sometimes a distance