

NCI Thesaurus Semantics

The NCI Thesaurus is built using the Ontylog dialect of description logic (DL). The semantics of Ontylog (effective May 2004) are summarized at <ftp://ftp1.nci.nih.gov/pub/cacore/EVS/ThesaurusSemantics/OntylogSemantics.doc>. Ontylog is fairly widely used in biomedical terminology construction. Notably, SNOMED/RT and SNOMED/CT are based on it, as is the US Veterans Health Administration's National Drug File/Reference Terminology.

As of May, 2004, the NCI Thesaurus employs all the semantic constructions offered by Ontylog *except* modal restriction and right identity. NCI is experimenting with modal restriction, and expects begin using it in the Thesaurus later in 2004. At this time NCI has no plans to employ right identity. The *concept* is the fundamental notion in Ontylog. In Ontylog concepts are abstract classes: there is no notion if an instance. Each concept denotes a semantic unit of meaning. Concepts are placed into *is_a* hierarchies through a process of classification. Following the usual practice of description logics, classification is accomplished by subsumption testing across the acyclic graph structure in which each concept is a node. The graph's edges are semantic relationships among the concepts.

NCI employs the Apelon, Inc. Terminology Development Environment (TDE) to build the NCI Thesaurus, and the Apelon Distributed Terminology Server (DTS) to make the Thesaurus accessible to users via programming interface and Web browser. The TDE and DTS software products implement Ontylog DL in software.

Each concept in NCI Thesaurus is either primitive (description limited to necessary conditions) or defined (description includes necessary and sufficient conditions). Subsumption of primitive concepts is established by the concepts' *defining super concepts*. Subsumption of defined concepts is established by its *direct super concepts*. Defining super concepts are manually determined, while direct super concepts are determined algorithmically during classification.

Originally all concepts in the NCI Thesaurus were primitive. As the terminology matures, the proportion of defined concepts is increasing. Top-level concepts will remain primitive as they express axiomatic knowledge used to infer the meaning of defined concepts and assure that branching at the top of the hierarchy trees in the Thesaurus is well formed. Note that NCI Thesaurus supports concept *polyhierarchy*. That is, a concept may have more than one super concept.

Many description logics make a distinction between generic concepts that describe sets and individual concepts that describe actual instances or elements of the sets. They break the collection of terms up into a T-Box for the former and an A-Box for the latter. NCI Thesaurus does not support instances for two reasons. First, it is designed to be a large and complex vocabulary that will be embedded in runtime systems supporting basic, translational and clinical research. In these systems instances are typically experimental data and research subject records and are stored in databases where transaction semantics and other core system concerns are paramount. The kind of inferencing that is performed

over the instances is readily performed over the generic concepts so there is no need for assertions about individuals in the Thesaurus. If Thesaurus supported exceptions, individuals might be required in the terminology. However, Ontylog has an enforced semantics, so there are no exceptions in the Thesaurus. Second, the distinction between generic and individual is considered a hard problem in mathematical foundations.

Semantic associations among concepts are called *roles* in NCI Thesaurus. Roles are binary associations between pairs of concepts. The roles currently available for use in Thesaurus are listed in

ftp://ftp1.nci.nih.gov/pub/cacore/EVS/ThesaurusSemantics/March04current_roles.xls.

Note that not all roles have been instantiated in the Thesaurus. They are available, but some have not yet have been used to restrict any concepts. Following DL convention, there is a *domain* and *range* value associated with each role. Ontylog uses the notion of a *kind* as values for domain and range. All concepts belong to one, and only one, kind.

Kinds may be thought of as disjoint classes or as data types. A list of the kinds in the Thesaurus, and each kind's definition, is provided at

ftp://ftp1.nci.nih.gov/pub/cacore/EVS/ThesaurusSemantics/Kind_Definitions.doc. A

graphic is also available that represents the roles and their domains and ranges. We use a distinctly colored box for each kind. Arrows with role names represent roles¹. Because many kinds and roles are required to satisfy the needs of cancer researchers, the graphic is large and visually complex. The current graphic is available in portable network graphic and Visio format at <ftp://ftp1.nci.nih.gov/pub/cacore/EVS/ThesaurusSemantics>. We have adopted the practice of including the domain and range names in the name of each role (Disease_Has_Primary_Anatomic_Site has the domain Disease and the range Anatomy).

The relationships among concepts in NCI Thesaurus are not fixed. We add and change them as we become aware of changes in the state of knowledge in the literature.

Users are welcome to suggest adding new relationships between existing concepts and/or removing or changing relationships that already link concepts. Naturally we would want to know why the requested changes are scientifically valid. We assert relationships that are well-established. We like to see three Medline or other citations, for example, reporting an association.

In creating role relationships for disease concepts, we have introduced a “may have” construct to indicate relationships that may or may not be present in a specific instance, specifically, characteristics that are not true of all individual cases of that disease. To give an example:

Accelerated Phase Chronic Myelogenous Leukemia:

Disease_May_Have_Finding: Myeloblasts 10-19 Percent of Bone Marrow Nucleated Cells

Disease_May_Have_Finding: Myeloblasts 10-19 Percent of Peripheral Blood White Cells

Disease_May_Have_Finding: Basophilia

¹ In Ontylog, roles are unidirectional, so we use one-headed arrows. In working with the more expressive logics, such as SHIQ(D) logic of RACER, for example, we would use bidirectional arrows. Arrows take the color of their domain kind, and the arrowhead touches their range kind.

To make the diagnosis, you need at least one of these findings but may not have all three. All are unquestionable and important characteristics.

Roles are not just binary assertions of semantic relationship between concepts that we put in the Thesaurus for human interpretation. As discussed earlier, they are used by our description logic engine to define the meaning of concepts in the Thesaurus and to determine tree-placement of concepts in the Thesaurus. When enough of the concepts are defined, rather than primitive, Thesaurus should be able to support inferencing by artificial intelligence, decision support and other such applications that require a "knowledge base".

The NCI Center for Bioinformatics relies on use cases to define the needs of our user communities for informatics resources. The NCI Thesaurus has adopted the discipline of listing the use case and individual need from within the use case to roles and kinds. The role and kind associations to use cases enables us to track not only which community relies on each role and kind, but also the use that they make of it. The current list of use case to kind mappings

<ftp://ftp1.nci.nih.gov/pub/cacore/EVS/ThesaurusSemantics/RoleToUCmapping040317.xls>

is provided as general information.

NCI Thesaurus is ontology-like. However unlike a typical ontology, it provides a lot of lexical content such as synonyms and preferred terms for concepts, and facts about concepts such as reference numbers used in other systems to refer to the concept. This information is represented in what Ontylog calls *properties*. Properties are labeled values associated with a concept. The labeled values' filler values are strings and other literals. In most DLs, this sort of information would restrict individuals, not abstract classes as in Ontylog. Refer to the DTS User Guide for additional discussion.