# Partitioning a Vocabulary's IS-A Hierarchy into Trees

Huanying (Helen) Gu[1], Yehoshua Perl[1], James Geller[1],
Michael Halper[2], James J. Cimino[3], Mansnimar Singh[4]

[1]CIS Dept., NJIT, Newark, NJ 07102    [2]Math & CS Dept., Kean College, Union, NJ 07083
[3]Dept. of Medical Informatics, Columbia University, New York, NY 10032
[4]CHREF, Inc., Wallingford, CT 06492

*Controlled medical vocabularies are useful in application areas such as medical information-systems and decision-support. However, such vocabularies are large and complex, and working with them can be daunting. It is important to provide a means for orienting users to the vocabulary's contents. This paper introduces a methodology for partitioning a vocabulary into small, meaningful pieces. The partitioning is done with respect to the vocabulary's IS-A hierarchy. The methodology, based on a set of rules for refining the IS-A hierarchy, is a process carried out by a user in conjunction with the computer. The methodology is demonstrated on a complex portion of a vocabulary.*

## INTRODUCTION

Controlled medical vocabularies ("vocabularies" for short) [1,2] play an important role in many medical enterprises that employ a large number of disparate information systems (e.g., clinical databases). Often, each such system has its own inherent "language" or terminology. A vocabulary allows for the integration of the different systems and the standardization of common information handling tasks, helping to reduce the overall cost of data processing. A vocabulary can also aid in the orientation of users of the information systems.

However, a vocabulary can be of an overwhelming complexity. E.g., the Medical Entities Dictionary (MED) [1], developed at Columbia-Presbyterian Medical Center, contains over 48,000 concepts, over 61,000 IS-A links and over 71,000 other links. Obviously, comprehending such a vocabulary is extremely difficult.

To enhance vocabulary comprehension, we have mapped the MED into an object-oriented database (OODB) schema [3,4]. The schema of this OODB, called the Object-Oriented Healthcare Vocabulary Repository, captures the structure of the vocabulary in a compact form which aids comprehension. However, each class in the schema contains on average about 500 concepts. Thus, further comprehension efforts are needed.

In this paper, we present a methodology to make large, complex vocabularies easier to understand. Our approach is based on the partitioning of the vocabulary into *manageably-sized, meaningful* units. The partitioning centers around the IS-A hierarchy which relates specialized concepts (sub-concepts) to generalized concepts (super-concepts). It serves as the vocabulary's backbone and supports property inheritance.

Our methodology is based on the following premises: (1) A vocabulary's IS-A hierarchy is more comprehensible than the entire vocabulary; (2) An IS-A hierarchy consisting of a collection of trees is easier to comprehend than a lattice containing the same number of concepts. With these in mind, we present a methodology for extracting from a vocabulary (typically represented as a semantic network) a hierarchy composed of small trees, each representing a logical unit whose graphical representation can fit neatly on a computer screen. This hierarchy greatly enhances comprehension of the original vocabulary by users and system designers. The methodology is a joint effort of a user (the vocabulary designer or manager) and a computer. The process requires that a user, with the help of the computer, first refine the IS-A hierarchy according to some prescribed rules [5]. Then the computer can reduce the vocabulary to a collection of reasonably-sized trees.

In [5], we used a similar paradigm to partition large OODB subclass hierarchies. This paper presents an adaptation to a vocabulary's IS-A hierarchy and introduces a partitioning methodology. The IS-A relationship is differentiated into four kinds of relationships, three of which are unique to our approach. The paper demonstrates our methodology on a complex subnet of the MED.

## PARTITIONING FRAMEWORK

Throughout this paper, a vocabulary is a semantic network of concepts. The links are the semantic relationships (e.g., IS-A) between the concepts. As a first step, all non-IS-A relationships are removed from the network. The user, with the computer's aid, is required to make some refinements to the vocabulary's IS-A hierarchy before the computer performs the partitioning step. The refinement is carried out with respect to a set of prescribed rules based on the distinction between two major kinds of IS-A relationships [6]: *CATEGORY-OF* (abbreviated CATG-OF) and *ROLE-OF*.

**Definition 1:** CATG-OF relates a specialized concept to a more general concept when both are in the same context.

**Definition 2:** ROLE-OF relates a specialized concept to a more general concept when both are in different application contexts.

The decision whether a given IS-A relationship is a CATG-OF or a ROLE-OF depends on whether the super-concept and sub-concept are in the same context or not. An intuitive understanding of the

application is required to make such a decision. However, this decision is not always easy. In spite of extensive research (e.g., [7,8]), there is no widely accepted definition of context. As was shown in [8], researchers disagree on what contexts are. In our approach, we are not trying to define the notion of context. Rather, we are making the pretheoretical (axiomatic) assumption that contexts exist in human thinking, and we are requiring the designers of an application to identify them explicitly. Hence, the need for human assistance in the partitioning process.

In [5], we provided a theoretical paradigm for assigning of classes to contexts in an OODB schema, resulting in a "forest" subhierarchy (i.e., a collection of trees) of the original hierarchy. The choice between CATG-OF and ROLE-OF is supported by three prescribed rules which ensure that a forest subhierarchy can be identified. For lack of space, we state the rules modified to a semantic network without any explanations.(See [5].)

**Rule 1**: A concept must belong to exactly one context. Thus, the concepts of a hierarchy are partitioned into disjoint contexts.

**Rule 2**: Two concepts which have a CATG-OF ancestor concept can neither have a common CATG-OF descendant concept, nor can one be a CATG-OF descendant of the other.

**Rule 3**: For each context there exists one concept which is the *major* (or defining) concept for this context such that every concept in this context is a descendant of this concept.

We have proved that if these three rules are satisfied, then a concept has at most one CATG-OF super-concept. (See [5].) Due to this, we can guarantee that the CATG-OF hierarchy can be partitioned into a collection of trees. This forest structure serves to enhance comprehension of the original hierarchy. The trees denote contexts that concentrate on specific subjects and provide users with a focused view of the vocabulary.

## METHODOLOGY FOR HIERARCHY REFINEMENT

In the following description, we will explicitly note which parts are performed by the computer and which by a human expert. Extensive examples of the steps will be given in the next section.

**Step 1 (Computer):** Remove all relationships other than IS-A.

**Step 2 (User assisted by the Computer):** Refine the IS-A hierarchy of the vocabulary in order to satisfy the above three rules.

By the user's judgment, every IS-A link becomes either a CATG-OF or a ROLE-OF. Then the computer can partition the hierarchy into trees of the CATG-OF links.

Step 2 consists of a sequence of substeps carried out by a vocabulary administrator or domain expert. The domain expert makes some judgments based on his understanding of the application, while the computer performs supporting tasks which do not require complex, intuitive decisions. The judgments required of the expert are:

1. Defining some IS-A relationships as ROLE-OF while others are defined as CATG-OF so that the three rules are satisfied.
2. Identifying disjoint contexts which correspond to subtrees of CATG-OF relationships.

We will introduce an additional three-way distinction between kinds of ROLE-OF relationships.

**Step 2.1 (Computer)** *Topological sort*
Number the concepts in the hierarchy according to the order in which they are visited by a left-to-right breadth-first search.

**Step 2.2 (Human)** *Identify roots of contexts*
Scan the hierarchy top-down according to the order from Step 1. Identify concepts which should be defining concepts (roots) for contexts. The choice is made by the meaning and importance of a concept compared to its super-concepts' meanings.

After these concepts are identified, they become ROLE-OF their super-concepts to start new contexts rather than refine their super-concepts. This kind of ROLE-OF relationship is a *ROLE-OF Type I* denoting a switch of context between a super-concept and a sub-concept.

**Step 2.3 (Computer)** *Bottom up listing*
List all concepts with multiple super-concepts in bottom-up order. (Below we will explain why bottom-up processing is used here.)

**Step 2.4 (Human)** *Identify primary parent*
For each concept listed in Step 2.3, the expert needs to identify at most one super-concept in the same context. The IS-A relationship to this super-concept will then be CATG-OF, while all other super-concepts are designated ROLE-OF to denote that they belong to different contexts.

In our experience, an expert can usually easily determine which of the multiple super-concepts is the defining one, i.e., in the same context and should have a CATG-OF relationship directed to it. In a few cases, this decision is not easy. Then we try to determine which super-concept, if any, should have a CATG-OF relationship pointing to it, based on the partial context information we have already accumulated in our bottom-up processing. We distinguish three cases.

**Case 1:** *One of the super-concepts is definitional while the others are functional.*
We look into the context to which the sub-concept and its descendants belong. If it is definitional, we will prefer the definitional super-concept. If it is functional, then we will prefer the functional super-concept (or if there are several functional super-concepts, we will prefer the one which fits the function described in the context of the sub-concept). If the sub-concept is the only concept in its context, we will choose the definitional super-concept. In this case, one super-concept is chosen

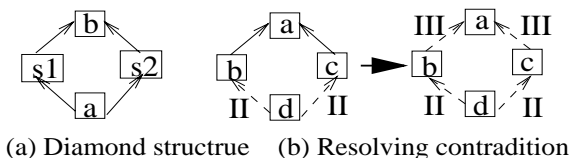(a) Diamond structrue    (b) Resolving contradition

Figure 1:

as primary super-concept. The concept becomes CATG-OF this primary super-concept and ROLE-OF the other super-concepts. This ROLE-OF is Type I since a context-switch occurs.

**Case 2:** *All super-concepts are definitional with the same or indistinguishable importance because each contributes equally to the definition of the sub-concept.*
The concept with multiple super-concepts could belong to the context of any of its super-concepts. However, by Rule 1 it cannot belong to more than one context. Also, we have no reason to prefer one over the other. Each choice will disassociate the concept from some other contexts. The solution is to require that such a concept start a new context representing an intersection of the contexts of all its super-concepts. Thus, this concept is ROLE-OF all its super-concepts. We call these relationships *ROLE-OF Type II*. ROLE-OF Type II is not an actual case of a context-switch but an artificial one required in order to satisfy the rules.

**Case 3:** *The sub-concept is a combination of the multiple super-concepts in different contexts, but one of them contributes more to the meaning than the others.*
The CATG-OF relationship should point to the preferred super-concept, as both concepts appropriately belong in the same context, while the other IS-A relationships should be ROLE-OF.

**Step 2.5 (Computer)** *Find diamond structures*
Scan the hierarchy bottom-up to find all the concepts with more than one super-concept. For each such concept $a$ and for each pair of super-concepts $s_1$ and $s_2$ of $a$, find a lowest common ancestor $b$ of both $s_1$ and $s_2$. For each pair of such concepts $a$ and $b$, output the IS-A subhierarchy containing $a$, $b$, and all the concepts which are descendants of $b$ and ancestors of $a$. We call such a subhierarchy a *diamond structure* $\langle a, b \rangle$ (Figure 1 (a)). The concept $a$ is called the *source* and $b$ is the *sink*.

**Step 2.6 (Human)** *"Cut" the diamonds*
Each diamond must contain concepts from more than one context to fulfill Rule 2. After Steps 2.1–2.5, this should be true. However, there is one scenario where we must artificially change the CATG-OFs to ROLE-OFs in order to resolve a contradiction. We call this situation the *Contradicting Diamond Case* and it occurs when the source $d$ of $\langle d, a \rangle$ is a ROLE-OF Type II of its super-concepts. All other concepts in the diamond then belong to one context. Since $d$ is the intersection of two super-concepts $b$ and $c$, both cannot belong to the same context of their super-concept $a$. Otherwise,

the intersection must belong to this common context. Thus, the concepts $b$ and $c$ also constitute separate contexts. The CATG-OFs are changed to ROLE-OFs. We call such a relationship ROLE-OF Type III. Originally, the link was a CATG-OF, but due to Rule 2 it now becomes a ROLE-OF (see Figure 1 (b)).

The steps used both top-down (2.1,2.2) and bottom-up (2.3,2.4,2.5,2.6) processing. The determination of the context of concepts is done top-down, as the context concept itself has a top-down nature with the root defining the context of its descendants. When scanning the hierarchy top-down, an expert can identify when a concept defines a new context. However, when trying to determine the concept's context from among those of its super-concepts, we need to know which descendants of the concept belong to the same context.

At this point, the computer partitions the IS-A hierarchy into a collection of trees as follows.
**Step 3 (Computer):** Remove all ROLE-OF relationships from the refined IS-A hierarchy.

We will have trees since no concept has multiple CATG-OF parents.

## APPLYING THE METHODOLOGY

We will apply the methodology to a subnetwork of the MED with a very complex hierarchy. We believe that the successful application of the method to the complex subnetwork will demonstrate the probable applicability to the MED as a whole. In the MED, the concept **Cortisporin Opth Oint (40)** has the most ancestors: 39. We will focus on the subnetwork containing this concept and all its ancestors. We will call it CN. Overall, CN contains 62 IS-A relationships and 157 other links. Figure 2 shows CN after Step 1.

Next, the computer performs a topological sort in Step 2.1, yielding the numbering 1–40 (Figure 3) used by the expert to scan the hierarchy top-down to find all concepts which define new contexts (Step 2.2). All IS-A relationships from these concepts are refined to ROLE-OF. Denoting an IS-A relationship as CATG-OF or ROLE-OF, the designer makes a modeling decision which may differ from one designer to another, influencing the emerging contexts. Decisions of a pharmacist will differ from those of a surgeon. Thus, each will create a different partitioning which represents his perspective on the vocabulary's overall interpretation. Here, we try to take the view of a general vocabulary administrator.

Let us see some examples of Step 2.2. The concept **Drug Allergy Class (9)** has one super-concept **Pharmacy Concepts**, a broad term that refers to various ways of grouping drug concepts. But the concept **(9)** is a group of drug concepts with allergic or antiallergic effects. Thus, it represents one drug classification according to new dimensions and is the defining concept for a context.
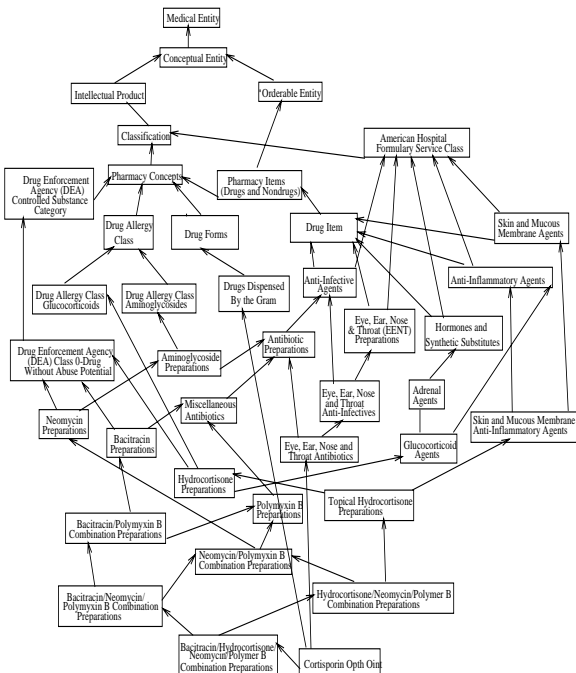
Figure 2: Complex subnetwork (CN) from MED



Figure 3: CN with CATG-OF and ROLE-OF

It is a Type I ROLE-OF of its super-concept. See Figure 3; note that we have used dashed lines for ROLE-OF to distinguish it from CATG-OF.

Consider **Glucocorticoid Agents (29)**; it has two super-concepts, **Adrenal Agents (24)** and **Anti-Inflammatory Agents (20)**. Glucocorticoid Agents are secreted by Adrenal glands. Thus **Adrenal Agents** indicates the physiological source for the glucocorticoid group of agents. Another child of **(24)** is mineralocorticoids (not shown) which are functionally distinct from glucocorticoids. **(20)** is a heterogeneous set that includes steroidal anti-inflammatory drugs like glucocorticoids and non-steroidal anti-inflammatory agents like Aspirin and Phenylbutazone Preparations. Thus, **(29)** starts a new context and is a Type I ROLE-OF its two super-concepts.

After applying Step 2.2, we have 11 defining concepts and contexts. These concepts, except **Medical Entity**, are ROLE-OF their super-concepts (see Figure 3).

In step 2.3, the computer finds the concepts with more than one super-concept in bottom-up order. The expert needs to identify one primary super-concept for these concepts, in Step 2.4. E.g., **(40)** has three super-concepts, **Bacitracin/Hydrocortisone/Neomycin/Polymyxin B Combination Preparations (39)**, **Drug Dispensed by Gram (15)** and **EENT Antibiotics (28)**. **(39)** defines the chemicals that form **(40)**. It uniquely defines the structural components of the ointment, and thus by Case 1 it is the primary super-concept of **(40)**. The super-concept **(15)** specifies the mode of dispensation, and the super-concept **(28)** specifies the site and action, so both
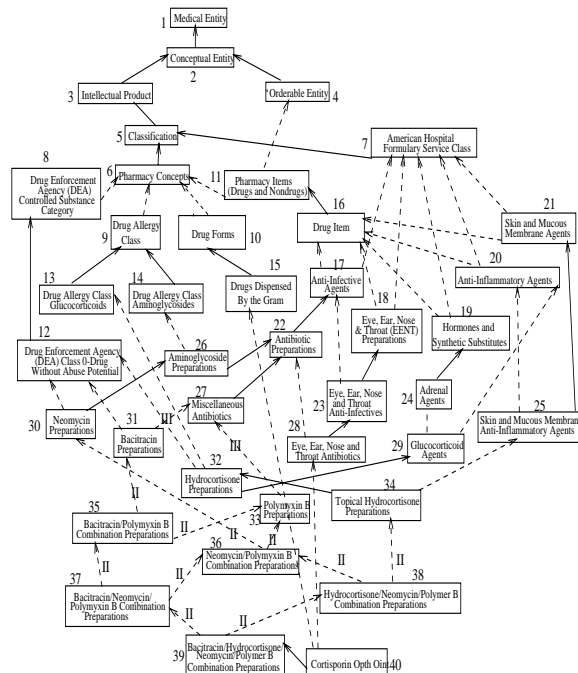
do not define the context of the concept. Thus, by Case 1, **(40)** is CATG-OF **(39)**, and ROLE-OF both **(15)** and **(28)**.

Let us check another example. **(39)** has two super-concepts, **Bacitracin/Neomycin/-Polymyxin B Combination Preparations (37)** and **Hydrocortisone/Neomycin/Polymyxin B Combination Preparations (38)**. Both contribute Neomycin and Polymyxin B to the concept. In addition, **(37)** contributes Bacitracin, and **(38)** contributes Hydrocortisone. All these chemicals together define **(39)**. According to Case 2 of Step 2.4, it is not possible to identify the primary super-concept. Hence, it is ROLE-OF of them both. This kind of ROLE-OF is Type II (marked with "II" in the figures).

After Step 2.4 is completed, none of the concepts has more than one CATG-OF super-concept.

After identifying all the diamond structures in bottom-up order in step 2.5, we need to check if any is a "contradicting diamond case" (Step 2.6). If so, then we need to change the appropriate CATG-OFs to ROLE-OFs.

Let us examine the diamond with source **Bacitracin/Polymyxin B Combination Preparations (35)** and sink **Miscellaneous Antibiotics (27)**. After Steps 2.2 and 2.4, both concepts **Bacitracin Preparations(31)** and **Polymyxin B Preparations (33)** are CATG-OF **(27)**. The source concept **(35)** is Type II ROLE-OF both **(31)** and **(33)** according to Step 2.4. This diamond structure is a contradicting case. Thus, at least one of **(31)** and **(33)** must have a Type III ROLE-OF of its super-concept **(27)**. Since **(31)** and **(33)** are in the same situation as we analyzed

Figure 4: Forest hierarchy of CN

before, both are given this designation ROLE-OF III. ("III" in Figure 3). This is the only contradicting diamond structure in CN. This way, we convey that both concepts **(31)** and **(33)** were separated from their parent's contexts just to fulfill Rule 2. But for other purposes, they may be part of the context to which **(27)** belongs.

Now, Step 3 is carried out by the computer which removes all ROLE-OF relationships to obtain the forest subhierarchy of the original subnetwork. Figure 4 shows the different contexts as trees in the forest. The relationship between concepts of different contexts (trees) is ROLE-OF.

The CN hierarchy is partitioned into 18 contexts, many of which are very small and seem to be too detailed. But this is not a typical subnetwork of the MED. By choosing CN, we picked a network with many interrelated subjects. Also, the contexts in Figure 4 are not complete since some context members that are not ancestors of **(40)** are not shown. E.g., **(35)** has other children Polysporin Opth. Oint. 3.5 Gm, Polysporin Topical Oint. 30 Gm, UD Polysporin Oint., etc. which are in the same context as **(35)**. We applied our methodology to the InterMED (an offshoot of the MED) of 3,000 concepts. It is partitioned into 545 contexts, 394 of them consist of single concept due to InterMED's incompleteness. (I.e., if more concepts of the MED are added to the InterMED, then some of these singleton concepts will get descendants and turn into actual contexts.) Thus, the InterMED is practically partitioned into 151 actual contexts with average size of 16.

## CONCLUSIONS

Vocabularies are important tools for many medical information processing tasks. Unfortunately, understanding a large complex vocabulary is difficult and time-consuming. In this paper, we present a methodology for partitioning a vocabulary into manageably sized meaningful units called contexts to aid comprehension. The partitioning centers around the IS-A hierarchy. The user, assisted by the computer, refines the IS-A relationships into 4 kinds of relationships. The subhierarchy consisting of one kind of the refined IS-A relationships results in a partition of the original vocabulary into a collection of trees representing contexts. As a result, the user can focus each time on a single context or the interaction between pairs of contexts. We believe that, for some applications, this may present a major improvement over studying "the part of the vocabulary that happens to be displayed on the screen."

### Acknowledgments

### References

1. J.J. Cimino, P.D. Clayton, G. Hripcsak, and S. Johnson. Knowledge-based approaches to the maintenance of a large controlled medical terminology. *JAMIA*, 1(1):35–50, 1994.
2. U.S. Dept. of Health and Human Services, NIH, National Library of Medicine. *Unified Medical Language System*, 1996.
3. H. Gu, J. Cimino, M. Halper, J. Geller, and Y. Perl. Utilizing OODB schema modeling for vocabulary management. In *Proc. '96 AMIA Annual Fall Symposium*, pp. 274–278, Washington,DC, 1996.
4. L. Liu, M. Halper, H. Gu, J. Geller, and Y. Perl. Modeling a vocabulary in an object-oriented database. In *CIKM'96*, pp. 179–188, Rockville, Maryland, 1996.
5. Y. Perl, J. Geller, and H. Gu. Identifying a forest hierarchy in an OODB specialization hierarchy satisfying disciplined modeling. In *Proc. COOPIS'96*, pp. 182–195, Brussels, Belgium, 1996.
6. J. Geller, Y. Perl, and E. Neuhold. Structure and semantics in OODB class specifications. *SIGMOD Record*, 20(4):40–43, 1991.
7. Saša Buvač and Richard Fikes. A declarative formalization of knowledge translation. In *CIKM'95*, pp. 340–347, Baltimore, MD, 1995.
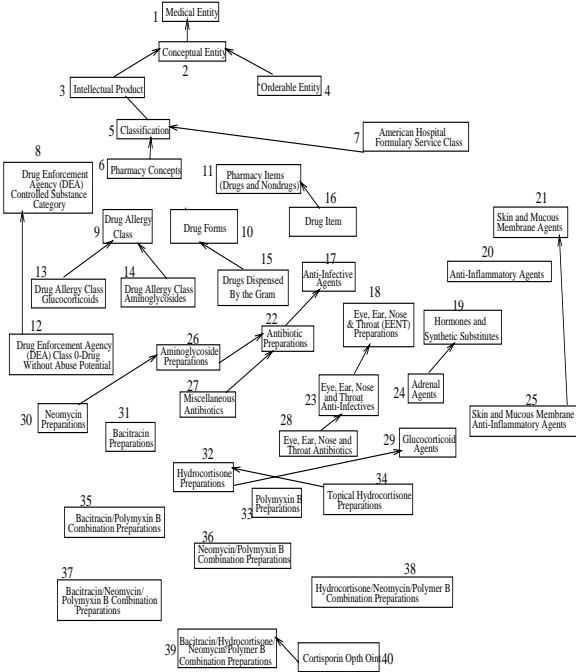8. L. Iwanska. Context in natural language processing. In *Working Notes of Workshop W13, IJCAI*. Montreal, Canada, 1995.