

Using OODB Modeling to Partition a Vocabulary into Structurally and Semantically Uniform Concept Groups*

Li-min Liu¹, Michael Halper², James Geller¹, Yehoshua Perl¹

¹CIS Dept.
NJIT
Newark, NJ 07102 USA
{limin, geller, perl}@homer.njit.edu

²Dept. of Mathematics & Computer Science
Kean University
Union, NJ 07083 USA
mhalper@turbo.kean.edu

Abstract

Controlled Vocabularies (CVs) are networks of concepts that unify disparate terminologies and facilitate the process of information sharing within an application domain. We describe a general methodology for representing an existing CV as an object-oriented database (OODB), called an Object-Oriented Vocabulary Repository (OOVR). A formal description of the OOVR methodology, which is based on a structural abstraction technique, is given along with an algorithmic description and a number of theorems pertaining to some of the methodology's formal characteristics. An OOVR offers a two-level view of a CV, with the schema-level view serving as an important abstraction that can aid in orientation to the CV's contents. While an OOVR can also assist in traversals of the CV, we have identified certain special CV configurations where such traversals can be problematic. To address this, we introduce—based on the original methodology—an enhanced OOVR methodology that utilizes both structural and semantic features to partition and model a CV's constituent concepts. With its basis in the notions of area and the recursively defined articulation concept, an enhanced OOVR representation provides users with an improved CV view comprising groups of concepts uniform both in their structure and semantics. An algorithmic description of the singly-rooted OOVR methodology and theorems describing some of its formal properties are given. The results of applying it to a large existing CV are discussed.

KEYWORDS: Object-Oriented Databases, Object-Oriented Models, Object-Oriented Systems, Knowledge Representation, Database Models

1 Introduction

A controlled vocabulary (CV) is a structure that houses knowledge in the form of concepts, subsumption links, and semantic relationships. CVs have become integral components of many information processing environments particularly within the healthcare field. Among their primary benefits are their support for information sharing and integration, decision-support, and *ad hoc* querying of domain (e.g., medical) knowledge [7, 32]. Examples of such systems from the healthcare domain include MeSH [23], CPT98 [1], SNOMED [8], ICD9-CM [33] (all of which have been integrated into the Unified Medical Language System (UMLS) [16, 18]), GALEN's Core Model [29] (expressed in GRAIL [28]), and the Medical Entities Dictionary (MED)

*This research was (partially) done under a cooperative agreement between the National Institute of Standards and Technology Advanced Technology Program (under the HIIT contract #70NANB5H1011) and the Healthcare Open Systems and Trials, Inc. consortium. This research was also partially supported by two grants from the New Jersey Commission for Science and Technology, one for the New Jersey Center for Software Engineering and the other for the Multi-lifecycle Engineering Research Center.

[6]. (References to related work on semantic networks [4], knowledge representation languages [22], and ontologies [24] can be found in [19, 20].)

One major aspect of many CVs is their enormous size and scope. A CV can easily consist of many thousands of concepts with a proportional number of inter-concept relationships. Given this fact, it may be hard for potential users and even a CV's own designers to orient themselves to the vast content of a CV and exploit its many advantages.

In previous work, we have devised a novel technique for modeling a CV as an object-oriented database (OODB) [2, 3, 5, 12, 17, 21, 34], a form we call an Object-Oriented Vocabulary Repository (OOVR) [19, 20]. Using our methodology, we have constructed OOVRs based on the MED and the InterMED [25]. Both OOVRs are up and running in ONTOS DB/Explorer [26, 31], a commercial OODB management system. Access to the InterMED OOVR is available on the Web in two forms [11, 27].

We have shown that the OOVR representation aids in vocabulary orientation and comprehension by providing an abstraction of the underlying CV contents. The schematic representation also helps in uncovering errors and inconsistencies that may have been introduced into a CV during its original development and subsequent refinement and expansion [13, 14].

In its original form, the OOVR methodology was presented as a two-phase process, with an initial phase followed by a refinement phase [19, 20]. In this paper, we first give a unified presentation of the methodology and prove some formal characteristics of OOVR representations. We also present a complete algorithmic description of the methodology.

An additional benefit of an OOVR is its support for more efficient browsing and traversal of a CV. However, during our experimentation with OOVR representations, we have encountered some special cases where small portions of a CV's concept configuration hindered the traversal process. The problems stemmed primarily from the fact that the OOVR methodology groups concepts together into an abstract entity when they have the same structure but not necessarily uniform semantics.

To address these issues, we present an enhanced OOVR methodology based on a revised partitioning scheme that performs a two-step decomposition of the source CV. As with the original OOVR technique, the first step breaks down a CV into collections of concepts, called *areas*, which have members exhibiting identical structure. In the second step, a special kind of area (called a *multi-rooted intersection area*) is further partitioned into collections of concepts called *partial areas* containing concepts uniform in their structure

and their semantics. The partial areas are based on the recursively defined notion of *articulation concept*. The new kind of OOV schema that emerges from this process has classes that are all “singly rooted,” i.e., the subnetworks of the CV which are the classes’ extensions each have a unique root concept. We will present the singly-rooted OOV methodology in its algorithmic form, along with theorems pertaining to various formal characteristics of singly-rooted OOV representations.

The remainder of this paper is organized as follows. In Section 2, we discuss the general structure of a CV. Section 3 presents the original OOV methodology, including an algorithmic specification of its partitioning process, the details of the construction of an OOV schema, and theorems that capture various formal characteristics of OOV representations. In Section 4, we present a sample CV traversal using an OOV in order to demonstrate the advantages of the extra abstraction layer afforded by the OOV schema. Section 5 describes some difficulties that can arise in certain special cases of OOV traversals. Then, in Section 6, we describe the formal aspects of the singly-rooted OOV methodology, including an algorithmic description and theorems about formal aspects of singly-rooted OOV representations. Section 7 presents the results of applying the enhanced methodology to the MED. Conclusions follow in Section 8.

2 Structure of a CV

A common formalism utilized in the construction of a CV is the semantic network, where each node is used to represent a unique concept from the knowledge domain. All concepts can exhibit two kinds of properties: (1) Attributes whose values are derived from some data types (such as integer or text string), and (2) relationships which are references to other concepts in the CV. Formally, an attribute is a mapping of a concept to a data type, while a relationship is a mapping of one concept to other concepts. For a concept v , we will use $P(v)$ to denote the set of all v ’s properties.

Each concept in a CV is defined with the attribute *name* that holds the concept’s associated *term* (i.e., printable value) [10]. In order to satisfy the nonambiguity and synonymy criteria for CVs (proposed in [6, 7]), it is assumed that each concept also has the attribute *synonyms* whose value is the entire set of acceptable secondary names for a concept. The concept subsumption (IS-A) hierarchy is a fundamental aspect of a CV. Structurally, it is an acyclic collection of IS-A links, each of which connects a subconcept to a related superconcept. The multiple classification criterion [6, 7] requires that the IS-A hierarchy be a directed acyclic graph (DAG), allowing for any concept to have multiple parents.

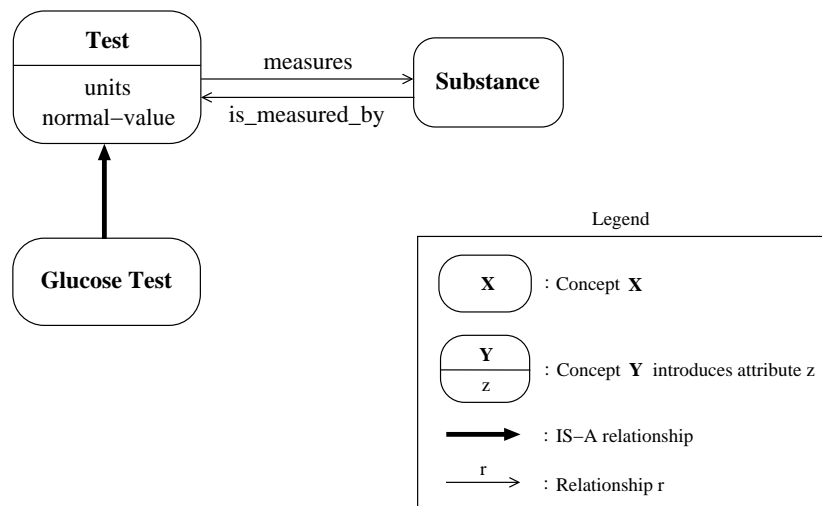


Figure 1: Concepts **Test**, **Glucose Test**, and **Glucose**

The IS-A hierarchy plays two important roles. First, it supports subsumption-based reasoning. For example, a user, wishing to know if a patient is taking antibiotics and knowing already that he is on Tetracycline, can consult the CV to learn that **Tetracycline IS-A Antibiotic**.¹ The second aspect of the IS-A hierarchy is inheritance: A subconcept inherits all the properties exhibited by its superconcepts. For example, the concept **Sodium Test IS-A Test**, and therefore the set of properties of **Sodium Test** is a superset of the properties of **Test**. If a concept has multiple parents, then it inherits properties from each of them.

Another assumption that we make, without loss of generality, is that a CV satisfies the following rule [6]:

Rule (Uniqueness of Property Introduction): A given property x can only be introduced at one concept in the CV. \square

A CV is also assumed, without loss of generality, to have a single root at the top of its IS-A hierarchy. We will refer to the root concept as **Entity**, which is defined to have the attributes *name* and *synonyms*. By inheritance, all other concepts in the CV will have these attributes, too.

We will be using the following notation when drawing a CV. A concept is a rectangle having rounded corners with its name written inside. Any attributes introduced by the concept (when shown) are listed below the name and are separated from it by a line. A relationship is a labeled arrow from the source

¹Some typographical conventions: A bold face font will be used when writing the names of concepts. Properties of concepts will appear in italics and will be written strictly in lowercase letters. Object classes will be written in italics and will start with uppercase letters.

concept to the target concept. Figure 1 shows a portion of CV with three concepts: **Test**, **Glucose Test**, and **Substance**. The concept **Test** introduces the attributes *units* and *normal-value* and the relationship *measures* to **Substance**. **Substance** introduces the relationship *is-measured-by* (the converse of *measures*) but no attributes. **Glucose Test** IS-A **Test** and therefore inherits all of **Test**'s properties.

3 OOV R Methodology

3.1 Partitioning a CV into Areas

Our OODB modeling of a CV is based on a structural abstraction of its network. The abstraction is derived from a partitioning of the network with respect to the notion of *area*. After defining area and other fundamental terminology, we prove some formal characteristics of the partition and its elements.

Definition 1: (Area) An area of a CV is an induced subgraph [9] which contains all concepts that have the exact same properties. \square

A CV is partitioned by its areas since each concept belongs to one and only one area. As we shall see, the partitioning of the CV into areas closely follows the property-introducing and inheritance patterns of the IS-A hierarchy, and this partition can be automatically identified in a top-down manner.

Definition 2: (Property set of an area) For an area A , $P(A)$ denotes the set of properties of any (and all) of its constituent concepts. \square

Definition 3: (Property-introducing concept) A concept at which one or more new properties are introduced into the CV is called a property-introducing concept. \square

An example of a property-introducing concept from the MED is **Pharmacy Items (Drugs and Non-drugs)** which, among other things, introduces the attribute *drug-trade-name*.

Definition 4: (Root of an area) A concept v residing in area A is called a root of A if A contains no parents of v . \square

The concept **Lab Diagnostic Procedure** is a root because its one parent **Diagnostic Procedure** belongs to a different area.

If an area has a single root, then the area is named after that concept. The area whose root is **Lab Diagnostic Procedure** is named "Lab Diagnostic Procedure Area."

Definition 5: (Property-introducing area) An area containing a property-introducing concept is called a property-introducing area. \square

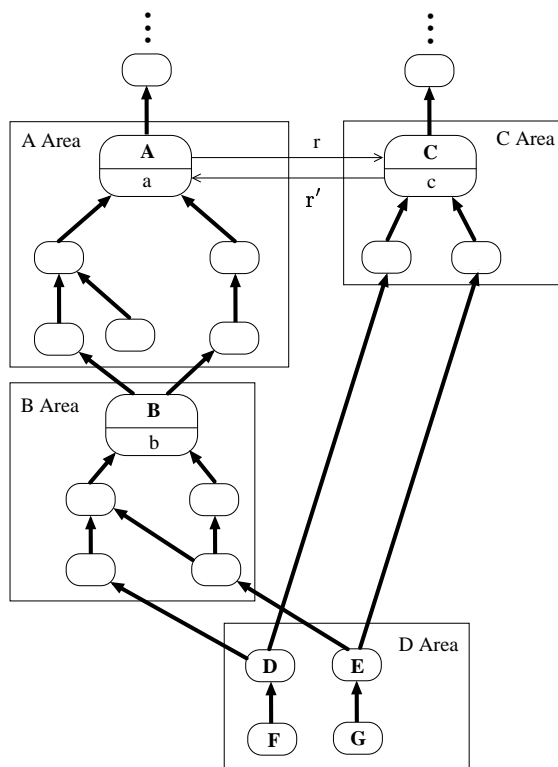


Figure 2: Four areas of a CV

An example is the Lab Diagnostic Procedure Area. In addition to property-introducing area, there is another kind of area defined in terms of *intersection concept*:

Definition 6: (Intersection concept) Let \mathbf{v} be a concept which is not a property-introducing concept and which has multiple superconcepts $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n$ ($n > 1$). The concept v is called an intersection concept if the following condition holds: $\forall i: 1 \leq i \leq n, P(\mathbf{v}) \neq P(\mathbf{w}_i)$. That is, the set of properties of v differs from all of its parents' sets of properties. Note that $P(\mathbf{v}) = \bigcup_{i=1}^n P(\mathbf{w}_i)$. \square

We use the designation “intersection concept” because v lies at the junction of (at least) two independent inheritance paths.

Definition 7: (Intersection area) An area containing an intersection concept is called an intersection area. \square

Figure 2 shows an example with three property-introducing areas (A Area, B Area, and C Area) and an intersection area. The only concepts in the property-introducing areas with their names displayed are **A**, **B**, and **C**, the respective roots. **A** introduces the attribute a ; **B**, the attribute b ; and **C**, the attribute c . **A** also introduces the relationship r directed to **C**, which itself introduces relationship r' , the converse of r .

D Area (Figure 2) is an intersection area. Unlike a property-introducing area, an intersection area can have more than one root. D Area has two roots **D** and **E**, both of which have two parents, one residing in B Area and the other in C Area. For a multi-rooted intersection area, the first identified root is used to name the area. The concept **D** was identified as a member of this area first, and hence the area was named D Area. The concepts **F** and **G** are members of D Area because they are children of **D** and **E**, respectively. **F** and **G** are not roots of D Area. None of the concepts in D Area has any intrinsic properties. All properties are inherited from outside the area. It is not possible for an intersection area to contain a property-introducing concept since such a concept would induce a new property-introducing area.

In the following, we present some formal characteristics of the partition of a CV in terms of areas.

Lemma 1: A property-introducing concept is a root of its area.

Proof: Let \mathbf{v} be a property-introducing concept that introduces property p . If \mathbf{v} is the concept **Entity**, then \mathbf{v} is clearly a root of its area. Otherwise, \mathbf{v} has parents in the CV. Thus, the parents of \mathbf{v} do not have the property p , and the property sets of the parents must all be different from $P(\mathbf{v})$. Hence, none of \mathbf{v} 's parents reside in \mathbf{v} 's area. ■

Lemma 2: A root of a property-introducing area is a property-introducing concept.

Proof: Assume to the contrary that a root \mathbf{r} of a property-introducing area A is not a property-introducing concept. Let \mathbf{v} be a property-introducing concept contained in A , and let p be a property that \mathbf{v} introduces. By Lemma 1, \mathbf{v} is also a root of A . Since \mathbf{r} is a root of A , it is not a descendant of \mathbf{v} , and thus it does not have the property p . But this implies that \mathbf{r} and \mathbf{v} are in different areas—a contradiction. ■

Lemma 3: All areas have at least one root.

Proof: Areas are induced subgraphs [9] of the CV and, therefore, because the overall IS-A hierarchy of a CV is a DAG, each subhierarchy contained in an area will be a DAG, too. Since a DAG must have at least one root, the same is true of an area. ■

Lemma 4: A property-introducing area has exactly one root.

Proof: By Lemma 3, an area must have at least one root. Suppose to the contrary that a property-introducing area A has at least two roots \mathbf{r}_1 and \mathbf{r}_2 . By Lemma 2, the root \mathbf{r}_1 (\mathbf{r}_2) is a property-introducing concept introducing, say, a property p_1 (p_2). By the “uniqueness of property introduction” rule (see Section 2), $p_1 \neq p_2$. However, \mathbf{r}_2 is not a descendent of \mathbf{r}_1 since \mathbf{r}_2 is also a root of A . Hence, \mathbf{r}_2 does not inherit the property p_1 introduced by \mathbf{r}_1 . Therefore, $P(\mathbf{r}_1) \neq P(\mathbf{r}_2)$, and \mathbf{r}_1 and \mathbf{r}_2 do not reside in the same

area—a contradiction. ■

An intersection area, in contrast, can have multiple roots (Figure 2). Lemmas 1, 2, and 4 together give us:

Theorem 1: There is a one-to-one correspondence between the property-introducing concepts, property-introducing areas, and the roots of these areas. ■

Corollary 1: The number of property-introducing areas is equal to the number of property-introducing concepts. ■

By Corollary 1 and the uniqueness of property introductions, there is at most one property-introducing area for each property, which gives us:

Corollary 2: The number of property-introducing areas is bounded by the overall number of different properties defined in the CV. ■

Note that when several properties are introduced at the same concept, there is only one corresponding area introducing them.

Lemma 5: There are only property-introducing areas and intersection areas.

Proof: Let A be an arbitrary area rooted at \mathbf{r}_A . Let $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n$ ($n \geq 1$) be the parents of \mathbf{r}_A . Note that $\forall i: 1 \leq i \leq n, P(\mathbf{w}_i) \neq P(\mathbf{r}_A)$. If the union of the property sets of \mathbf{r}_A 's parents is different from \mathbf{r}_A 's property set [i.e., $\bigcup_{i=1}^n P(\mathbf{w}_i) \neq P(\mathbf{r}_A)$], then there is some property introduced at \mathbf{r}_A . In that case, \mathbf{r}_A is a property-introducing concept and A is a property-introducing area. Otherwise, by Definition 6, \mathbf{r}_A is an intersection concept since $\forall i: 1 \leq i \leq n, P(\mathbf{w}_i) \neq P(\mathbf{r}_A)$, and A is thus an intersection area. ■

Theorem 2: A CV is partitioned into disjoint areas which are either property-introducing areas or intersection areas.

Proof: Areas are disjoint by definition. By Lemma 5, every area is either a property-introducing area or an intersection area. ■

Below, we present the algorithm that partitions a CV into its respective areas. The algorithm operates in a top-down manner in its processing of the concepts of a CV. Its input is a complete CV, and its output is the CV's entire set of areas. An area will be named after a property-introducing concept or a first-encountered intersection concept. We refer to these concepts as “naming concepts.”

In the algorithm, \mathcal{A}_v will denote a set of concepts, each of which has the same set of properties, with v as its naming concept. \mathcal{S} is a set which holds all naming concepts. \mathcal{A}_{ALL} is a set which will contain

all \mathcal{A}_v 's. At the end, \mathcal{A}_{ALL} will be returned. Each element v in \mathcal{S} will later be used to name an area with the format v_Area . Every element v in \mathcal{S} will have an associated set \mathcal{A}_v in \mathcal{A}_{ALL} . Every concept v has an associated counter for unprocessed parents which is denoted as “p-counter[v].” This algorithm uses two auxiliary functions: “Num_parents_of” and “Is_proper_introducing.” Num_parents_of takes a concept as input and returns its number of parents. Is_property_introducing takes a concept as input and returns “true” if it is a property-introducing concept and “false” otherwise. Statements with the same indentation are in the same block. A concept can be processed only if all its parents have been processed. Therefore, initially the root of a CV is processed.

```

set_of_areas FUNCTION AREA_Partition(CV V)
BEGIN
  // initialization
  Q := newqueue();           // Q contains concepts ready to be processed.
  A_ALL := newset();        // A_ALL will hold all areas.
  S := newset();            // S will hold property-introducing concepts and
                             // first-encountered intersection concepts.
  FOR EACH concept v in V DO // Initialize p-counter for all concepts.
    p-counter[v] := Num_parents_of(v);
  enqueue(Q, root(V));      // Insert Entity into Q.

  WHILE ( NOT emptyqueue(Q) ) DO
    v := dequeue(Q);
    // If v introduces new properties, we generate a new property introducing area.
    IF ( Is_property_introducing(v) ) THEN
      insert(S, v);
      A_v := newset();      // Create a new set A_v to keep concepts in this area.
      insert(A_v, v);       // Insert v into A_v.
      insert(A_ALL, A_v);  // Insert the new set A_v into A_ALL.
    // If v has only one parent, v is in the same area as its parent.
    ELSE IF ( Num_parents_of(v) = 1 ) THEN
      Let w be the parent of v;
      Let A_u ∈ A_ALL be the set s.t. w ∈ A_u;
      insert(A_u, v);
    ELSE
      // v has multiple parents w_1, w_2, ..., w_n (n > 1).
      IF ( ∃i: 1 ≤ i ≤ n s.t. P(v) = P(w_i) ) THEN // v is not an intersection concept
        Let A_u ∈ A_ALL be the set s.t. w_i ∈ A_u;
        insert(A_u, v);
      ELSE
        //Concept v is an intersection concept.
        Let flag found := false;
        // Determine whether v is the first-encountered intersection concept in its area.
        // It is necessary to check the elements of S to determine whether this is the case.
        FOR EACH element c ∈ S DO
          IF ( NOT found AND P(v) = P(c) ) THEN
            // Concept v is not the first-encountered intersection concept.
            insert(A_c, v);
            set flag found := true;

```

```

        IF ( NOT found ) THEN
            // Concept  $\mathbf{v}$  is the first-encountered intersection concept.
            insert( $\mathcal{S}$ ,  $\mathbf{v}$ );
             $\mathcal{A}_{\mathbf{v}} := \text{newset}()$ ;
            insert( $\mathcal{A}_{\mathbf{v}}$ ,  $\mathbf{v}$ );
            insert( $\mathcal{A}_{\text{ALL}}$ ,  $\mathcal{A}_{\mathbf{v}}$ );
        // After we process concept  $\mathbf{v}$ , we need to decrease the p-counters of all  $\mathbf{v}$ 's children by one.
        // After the decrease, if any p-counter is equal to zero, we put the associated concept
        // into the queue since it is ready to be processed.
        FOR EACH child  $\mathbf{k}$  of  $\mathbf{v}$  DO
            p-counter[ $\mathbf{k}$ ]--;
            IF ( p-counter[ $\mathbf{k}$ ] = 0 ) THEN
                enqueue( $\mathcal{Q}$ ,  $\mathbf{k}$ );
    RETURN  $\mathcal{A}_{\text{ALL}}$ ;
END  $\square$ 

```

Let us illustrate the construction of a set \mathcal{A} constituting an intersection area with multiple roots. Suppose the first concept processed in the intersection area \mathbf{D} Area (Figure 2) is \mathbf{D} . We create a set $\mathcal{A}_{\mathbf{D}}$ with \mathbf{D} as its first element. Later, we will visit the concept \mathbf{E} , the other root of this area. We compare its property set to that of the concepts \mathbf{A} , \mathbf{B} , \mathbf{C} , and \mathbf{D} in the set \mathcal{S} of naming concepts. The property sets of concepts \mathbf{A} , \mathbf{B} , and \mathbf{C} do not match that of \mathbf{E} , but \mathbf{D} 's does match. Therefore, \mathbf{E} will be inserted into the existing set $\mathcal{A}_{\mathbf{D}}$. When \mathbf{E} is processed, the p-counter of \mathbf{G} is reduced from 1 to 0, and \mathbf{G} is inserted into the queue. Later on, when \mathbf{G} is deleted from the queue, it has only one parent \mathbf{E} , and thus is added to set $\mathcal{A}_{\mathbf{D}}$.

3.2 OOV_R Schema

In the OODB-version of the CV, each concept is represented by a unique object. The OOV_R's schema is constructed automatically after the identification of all areas. There is a one-to-one correspondence between the areas in the CV and the classes in the OOV_R's schema. That is, one class is defined to represent one area. The direct extension of a given class is identical to the set of concepts in the corresponding area in the CV. Due to this, we refer to the classes in the OOV_R schema as *area classes*. If the area is a property-introducing area, then we have a *property-introducing class*. Likewise, for an intersection area, there is an *intersection class*. In an OODB schema, a class defines a set of objects whose structure and behavior are the same. In our mapping, the instances of one class are exactly all those concepts that reside in a single area which, by definition, contains all concepts exhibiting identical properties.

The intrinsic properties of a property-introducing class are defined to be exactly those introduced by the root concept of its corresponding area. In addition, all the concepts in a property-introducing area must have

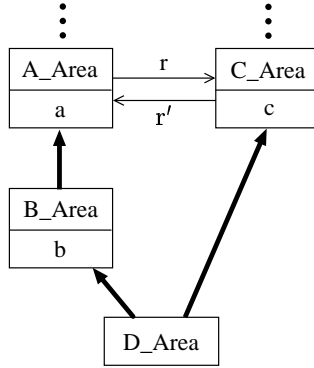


Figure 3: Area classes for the areas in Figure 2

the properties inherited by the root from its parent(s) in the CV. To capture this situation, the property-introducing class is placed in subclass relationships with those other area classes to which the parents of the root belong. In this way, the property-introducing class obtains all necessary properties: Some are defined intrinsically, while the others are inherited from other classes.

In Figure 3, we illustrate the above by showing the classes *A_Area*, *B_Area*, and *C_Area* that represent the corresponding areas in Figure 2. The classes are boxes with their names and attributes written inside. An ordinary relationship is a labeled arrow, while a subclass relationship is a bold arrow pointing from the subclass to the superclass. The ellipses indicate the omission of the subclass relationships of *A_Area* and *C_Area*. All property-introducing classes have at least one subclass relationship. The only exception is *Entity_Area*, the root of the OOVR schema.

Since an intersection area does not contain any property-introducing concepts, and, in fact, all properties of its concepts are obtained via inheritance, an intersection class does *not* introduce any properties of its own. Instead, it is defined to be a subclass of all other area classes which contain one or more parents of its root(s). An intersection class *always* exhibits multiple inheritance, i.e., it inherits from two or more superclasses.

Referring to Figure 3 again, we see the intersection class *D_Area*, representing D Area. *D_Area* is a subclass of both *B_Area* and *C_Area* because its roots (**D** and **E**) have parents residing in both those respective areas.

Our mapping technique may generate a “short-cut” of SUBCLASS_OF links. That is, it may happen that $Y \text{ SUBCLASS_OF } Z$, $X \text{ SUBCLASS_OF } Y$, and $X \text{ SUBCLASS_OF } Z$. In this case, the SUBCLASS_OF

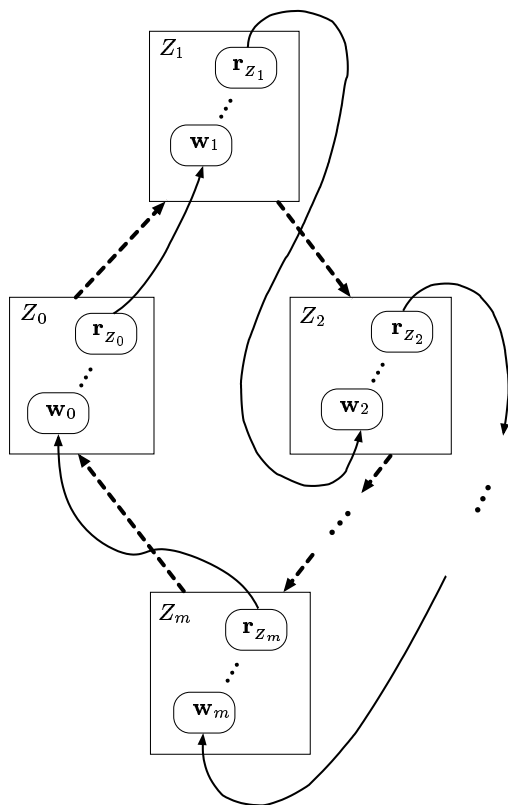


Figure 4: Invalid subclass configuration: Subclass cycles do not exist in an OOV schema. (SUBCLASS_OF drawn as a dashed arrow)

link from X to Z is a short-cut of the two links connecting X to Y and Y to Z . We have made the decision to omit this kind of subclass relationship from the OOV schema because it does not contribute to inheritance.

The final aspect of the mapping pertains to the IS-A hierarchy. All concepts have IS-A connections to other concepts (except for the root **Entity**). In the original network, **Entity** has the multivalued relationship “subconcept_of” that implements the IS-A hierarchy of concepts. In the mapping, this is translated into a multivalued, reflexive relationship *subconcept_of*, defined at the class *Entity_Area*. In this way, all concepts (objects) in the OOV representation have their required IS-A connections.

All OODB schemas must have acyclic subclass structures to avoid circular definitions of properties. Since the OOV schema is derived by an algorithm, it remains for us to prove that its induced subclass configuration is indeed acyclic. This result follows from the fact that the IS-A hierarchy of any CV is acyclic.

Theorem 3: The subclass relationships of an OOV schema are acyclic.

Proof: Assume to the contrary that an OOV schema contains a cycle of area classes Z_0, Z_1, \dots, Z_m (see

Figure 4) with respect to the SUBCLASS_OF relationship. Let the naming concepts of these classes be $\mathbf{r}_{z_0}, \mathbf{r}_{z_1}, \dots, \mathbf{r}_{z_m}$, respectively. According to the construction of the schema, for each class Z_i ($0 \leq i < m$), there is an IS-A connection from its naming concept \mathbf{r}_{z_i} to a concept \mathbf{w}_{i+1} in Z_{i+1} . (Also: \mathbf{r}_{z_m} IS-A \mathbf{w}_0 in Z_0 .) Whether \mathbf{r}_{z_i} is a property-introducing concept or an intersection concept, $P(\mathbf{w}_i) = P(\mathbf{r}_{z_i}) \supset P(\mathbf{w}_{i+1}) = P(\mathbf{r}_{z_{i+1}})$, and therefore $P(\mathbf{r}_{z_i}) \supset P(\mathbf{r}_{z_{i+1}})$. From this, we see that $P(\mathbf{r}_{z_1}) \supset P(\mathbf{r}_{z_2}) \supset \dots \supset P(\mathbf{r}_{z_m}) \supset P(\mathbf{r}_{z_1})$. In other words, $P(\mathbf{r}_{z_1}) \supset P(\mathbf{r}_{z_1})$ —a contradiction. ■

Overall, the OOVr schema provides a structural abstraction of the underlying network of the CV [19, 20]. Concepts with the same properties are grouped into areas which in turn are modeled as object classes; the concepts themselves become the objects of the OODB. This schema represents a substantial reduction in size from the original CV. In Figure 5, we show the InterMED OOVr schema. The InterMED has 2,820 concepts in its network, but its schema contains only 39 area classes, nine of them being intersection classes (below the dashed line). This schema can be used to gain an understanding of the InterMED.

4 Navigation Examples

In this section, we demonstrate how the schema helps to speed up traversals of a CV. Suppose that a user wants to search for some information, say, in the InterMED, but does not know the name of the concept for which the information is desired. For example, suppose a user is looking for a drug to treat fever and coughing in children. While the user does not remember the names of such drugs, he may recognize them when encountered. This is a natural application of an IS-A hierarchy traversal, with the user employing his knowledge about the target concept to guide the choices at the different levels of the hierarchy.

Using the InterMED OOVr representation, we enable a faster traversal involving both the schema and the underlying knowledge content. The depth of the InterMED’s IS-A hierarchy is 11, while the depth of its OOVr schema’s subclass hierarchy is just 4. Instead of traversing the InterMED hierarchy through its many levels, we traverse the OOVr schema until the proper area class (say, X_Area) is identified. This is easier because the schema presents higher-level subject areas rather than detailed concepts. The user only needs to make a very general judgment about whether a desired concept fits into a given class or not. Once that judgment is made, the user will switch to that part of the concept network belonging to X_Area . The traversal will run through this subhierarchy until the desired concept is recognized (or its absence is noted).

This traversal is shorter since the number of traversing steps is bounded by the sum of the depth of the

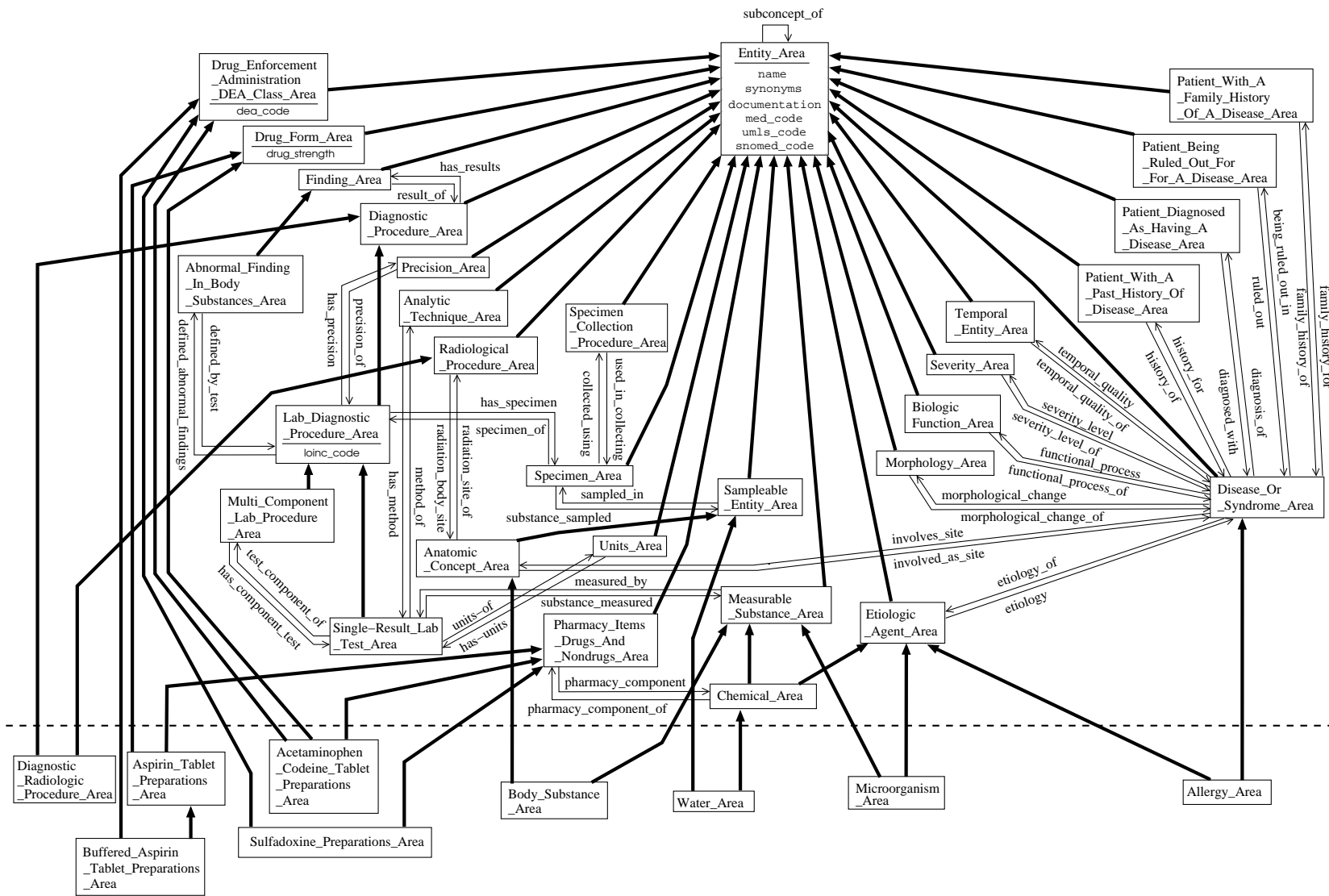


Figure 5: Schema of the InterMED OOV

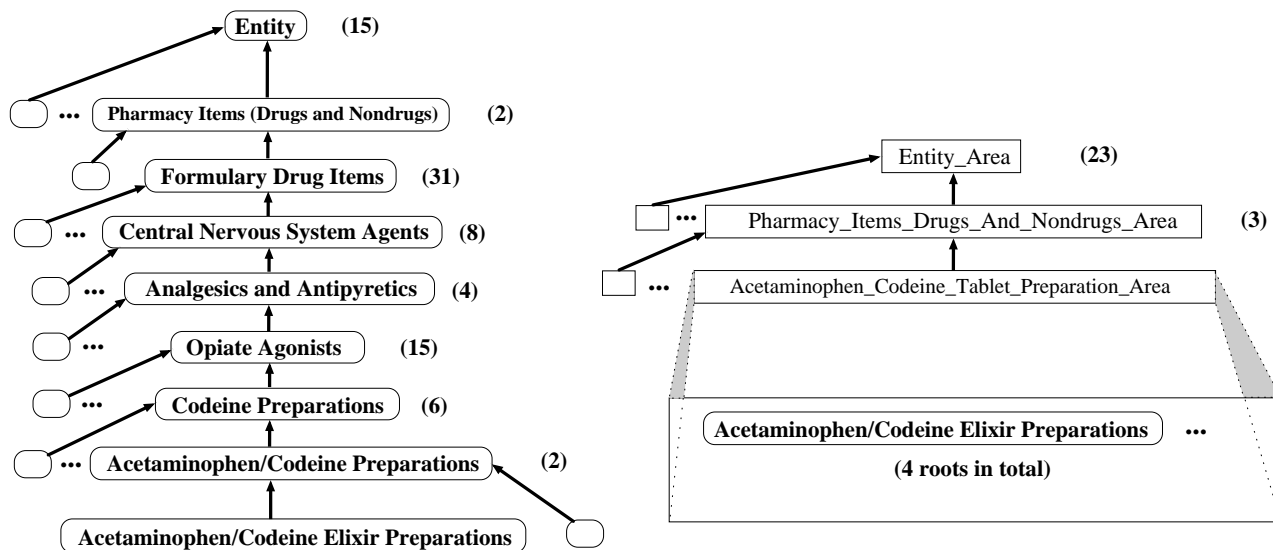


Figure 6: (a) Concept-level browsing path between **Entity** and **Acetaminophen/Codeine Elixir Preparations**; (b) Corresponding schema-level browsing path

OOVR schema and the depth of the subhierarchy of X_Area . Furthermore, the traversal is faster because the number of subclasses of a given class in the OOVR is typically much smaller than the number of children of a concept in the InterMED. As a traversal requires scanning through a list of children and choosing one of them, it will be easier and faster at the schema level.

Let us demonstrate the abovementioned traversal: looking for a drug to treat fever and coughing. First, let us perform the traversal at the concept level in the InterMED. The traversal starts at the root **Entity**, having fifteen children. Since we are looking for a medication, **Pharmacy Items (Drugs and Nondrugs)** is chosen. The process continues in this manner all the way down to **Acetaminophen/Codeine Elixir Preparations**. The entire traversal path is illustrated in Figure 6 (a). Alongside each concept, we list its number of children, indicating the range of choices encountered at that level. Overall, this traversal of a path of 9 concepts required scanning a total of 83 children.

Let us now demonstrate the same traversal in the OOVR (Figure 5). We start with $Entity_Area$ and travel through $Pharmacy_Items_Drugs_And_Nondrugs_Area$ to $Acetaminophen_Codeine_Tablet_Preparation_Area$, a leaf. At that point, the traversal switches to the concept level. Since $Acetaminophen_Codeine_Tablet_Preparation_Area$ is an intersection class with only 4 roots, we can easily find the concept **Acetaminophen/Codeine Elixir Preparations**. This is illustrated in Figure 6 (b), where the number beside a class is its respective number of subclasses. This traversal spans 3 classes, with a total of 26 subclasses, and 4 concepts.

Thus, the total number of scanned items is 30, quite a bit fewer than the 83 required before.

To be formal in our comparison of the two traversal methods, we need to define the notion of *browsing path* on both the concept level and the area (class) level.

Definition 8: (Concept-level browsing path) A concept-level browsing path is a sequence of concepts $(\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n)$ such that \mathbf{c}_{i+1} IS-A \mathbf{c}_i , $1 \leq i < n$. \square

Definition 9: (Schema-level browsing path) A schema-level browsing path is a sequence of area classes (A_1, A_2, \dots, A_k) such that A_{i+1} SUBCLASS_OF A_i , $1 \leq i < k$. \square

We can now properly state how the SUBCLASS_OF relationships at the OOV schema level properly reflect the IS-A relationships in the CV. For every concept-level browsing path $(\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n)$, there exists a corresponding schema-level browsing path (A_1, A_2, \dots, A_k) which satisfies the following conditions.

1. Concept \mathbf{c}_1 is an instance of area class A_1 .
2. Concept \mathbf{c}_n is an instance of A_k .
3. There exists a partition of $(\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n)$ into disjoint subpaths of consecutive concepts, say, $\mathbf{c}_{i_1}, \dots, \mathbf{c}_{i_e}$, which are paths in the induced subnetwork of an area A_j ($1 \leq j \leq k$).

5 Inadequacy of the Multi-Rooted OODB Modeling

5.1 Browsing Multi-Rooted Intersection Areas

The traversal at the schema level is very effective when all areas are singly-rooted. In such a case, the root concept subsumes all other concepts in the area and conveys the area's general semantics. However, only property-introducing areas are guaranteed to be singly-rooted.

Traversals in the context of multi-rooted intersection classes may not proceed so smoothly. This is because the name of the class is chosen arbitrarily from among the roots. Instead of conveying the general semantics for the whole area, the chosen root may capture only the essence of the concepts which are its descendants. But some concepts in the area—aside from the other roots—may not even be descendants of that root. In fact, the roots may be very dissimilar from an interpretive viewpoint; grouping them together was strictly the result of structural similarity. It is therefore sensible to reexamine whether those concepts should have been grouped together in the first place.

As an example, let us look at the multi-rooted intersection area shown in Figure 7 (a), which was gleaned from the MED. Overall, Figure 7 (a) contains six concepts. **ICD9 Disease** belongs to *ICD9_Element_Area* and **Disease or Syndrome** belongs to *Disease_or_Syndrome_Area* [see Figure 7 (b)]. The concepts

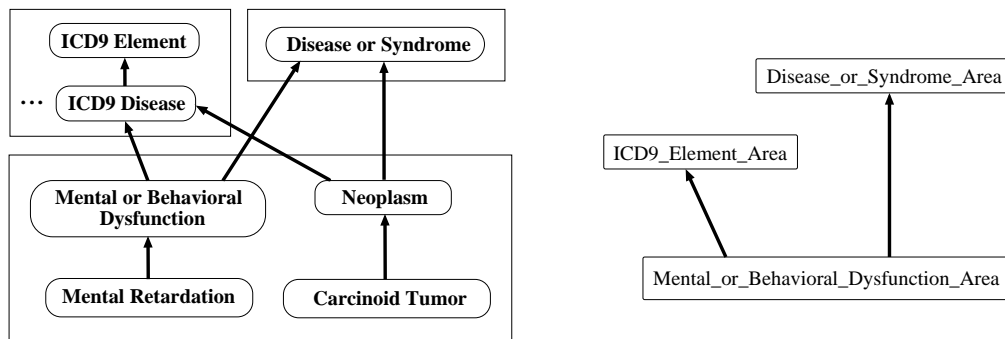


Figure 7: (a) Three areas including an intersection area (in the bottom box); (b) Their OOV schema

Mental or Behavioral Dysfunction and **Neoplasm** are children of both **ICD9 Disease** and **Disease or Syndrome**. They have the same structure and thus are roots of the same intersection area, *Mental-or-Behavioral-Dysfunction-Area*. Actually, in the MED, there are 29 roots for this intersection area! These include **Infectious Disease**, **Disorder of Circulatory System**, **Disorders of Nervous System**, etc. Our mapping methodology randomly chooses one of them to be the area’s name. In this example, **Mental or Behavioral Dysfunction** has been selected [Figure 7 (b)]. However, any of the other 28 roots could have been selected.

One will note that there is almost no similarity between **Mental or Behavioral Dysfunction** and **Neoplasm**, even though they are in the same area. With **Mental or Behavioral Dysfunction** as the area’s name, it is hard to imagine that **Neoplasm** also belongs there. In this case, the schema diagram does not provide a useful abstraction for assisting users in browsing an area of the CV.

Let us express the problem in terms of a browsing path. Consider the concept-level browsing path (**ICD9 Element**, **ICD9 Disease**, **Neoplasm**, **Carcinoid Tumor**). The corresponding schema-level browsing path is (*ICD9_Element_Area*, *Mental-or-Behavioral-Dysfunction_Area*). However, since **Carcinoid Tumor** is a descendent of **Neoplasm** and is not a descendent of **Mental or Behavioral Dysfunction**, a user will not know to choose this schema-level browsing path while searching for **Carcinoid Tumor**.

5.2 Establishing Subclass Relationships for Intersection Area Classes

We have discovered an additional problem in the modeling of multi-rooted intersection areas. Our mapping methodology yields a configuration of subclass relationships that does not reflect the pattern of IS-A links that cross the boundaries of such areas (the “cross-area IS-A relationships”). There are several equivalent

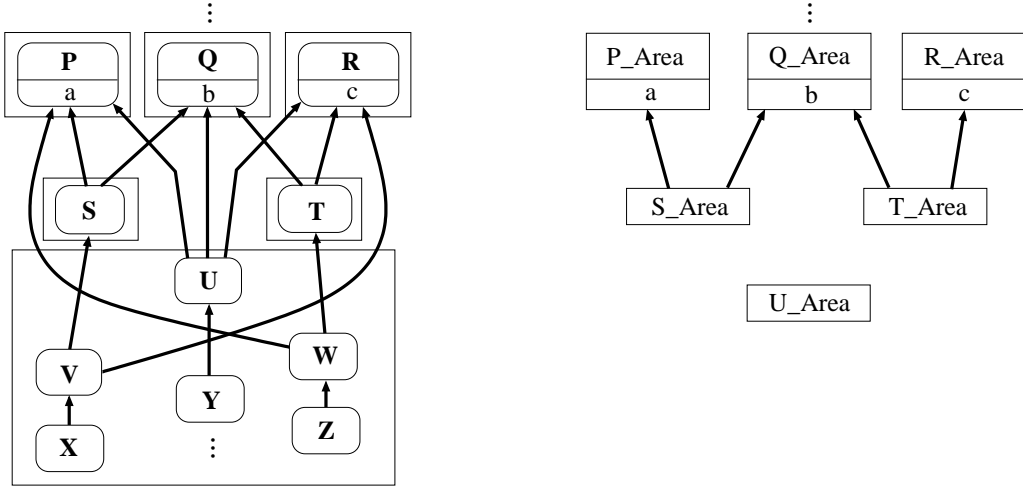


Figure 8: (a) CV excerpt; (b) Its OOV schema with subclass relationships from U_Area omitted

OODB modeling alternatives that properly capture the structure of the areas, but none of these is sufficient to convey the full extent of the cross-area IS-A relationships. This makes it more difficult to effectively utilize the OOV schema.

Consider Figure 8 (a) which contains eleven concepts, with only the top three, **P**, **Q**, and **R**, introducing new properties a , b , and c , respectively. **S**, **T**, and **U** have several parents which reside in different property-introducing areas and thus should belong to intersection areas. In fact, they should be roots of three different intersection areas since their property sets are different. **V** and **W** differ from **S**, **T**, and **U** in that one of their parents resides in an intersection area while the other resides in a property-introducing area. However, **V** and **W** have the same property set as **U**. Our mapping methodology groups **U**, **V**, **W**, **X**, **Y**, and **Z** into the same intersection area. Its name is arbitrarily chosen to be U_Area [Figure 8 (b)].

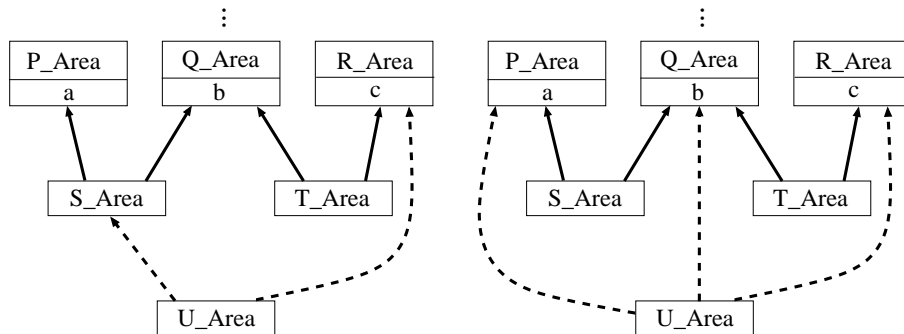


Figure 9: Alternative schemas for the areas in Figure 8 (a)

There is a problem with defining U_Area 's subclass relationships (omitted from Figure 8 (b)). Consider

the concept **V**: it has **S** and **R** as parents. Thus, one may set *U_Area*'s subclass relationships to point to *S_Area* and *R_Area* [Figure 9 (a)]. However, the absence of a relationship between *U_Area* and *T_Area* may mislead the user into thinking that there is no IS-A link between any concepts of these two areas. Actually, **W** IS-A **T**. The schema of Figure 9 (a) has no schema-level browsing path corresponding to the concept-level browsing path (**R**, **T**, **W**). A similar problem arises if **W**'s IS-A links are used to establish *U_Area*'s subclass relationships. Consider the concept **U**: it has parents **P**, **Q**, and **R**. Thus, one may set *U_Area* to be a subclass of *P_Area*, *Q_Area*, and *R_Area* [Figure 9 (b)]. This schema might lead users to believe that there are no IS-A links between concepts in *U_Area* and those in *S_Area* or *T_Area*. Here, though, **V** IS-A **S** and **W** IS-A **T**. There is no schema-level browsing path corresponding to the concept-level browsing paths (**R**, **T**, **W**) and (**P**, **S**, **V**) in Figure 9 (b).

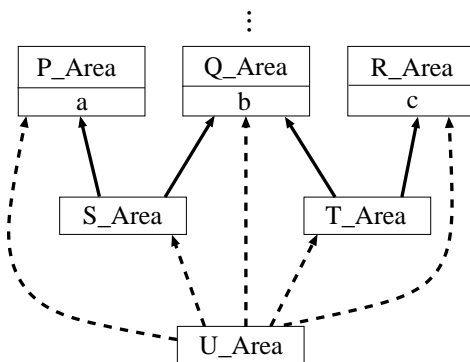


Figure 10: Another alternative schema for Figure 8 (a)

Another alternative is to set *U_Area*'s subclass relationships to mirror all IS-A relationships of its roots. For each root r , we could define subclass relationships from *U_Area* to all area classes containing a parent of r . Using this approach, we obtain a set of parent classes for *U_Area* which is the union of the parent classes from the alternatives considered above. In Figure 10, *U_Area* has five parents. The five subclass relationships from *U_Area* can, once again, be misleading. We do not know which relationship originated from which root. Furthermore, that same schema would be generated if there existed a single root in *U_Area* with five superconcepts in the areas *P_Area*, ..., *T_Area*. Thus, this is not desirable, either. All these choices are structurally equivalent since the resulting property sets for *U_Area* are the same. However, what is lost is some of the OOV schema's effectiveness in reflecting aspects of the IS-A hierarchy of the original CV.

6 Singly-Rooted OOV Methodology

The two problems that were presented in Section 5 arise from placing concepts of potentially widely varying semantics in the same intersection area and its corresponding area class. Recall that, in general, a class is a construct that gathers together objects with the same structure *and* semantics. In our mapping methodology, most classes satisfy this condition. Certainly, the structural aspect is satisfied by all area classes. Property-introducing classes are semantically cohesive due to their unique roots, which provide the areas' names. The same can be said for an intersection class having a single root. However, the synchronization of structure and semantics breaks down for multi-rooted intersection areas. All concepts of such an area have the same structure but not necessarily similar semantics because some concepts may be descendants of one root and not directly related at all to another root. It is unlikely that any single root provides appropriate “root semantics” for the entire area.

To preserve the ordinary interpretation of OODB classes as having objects with the same structure and semantics, and indeed to support effective CV access via the OOV schema, we need to further partition a multi-rooted intersection area into separate singly-rooted groupings, which we call *partial areas*. Once this is accomplished, the intersection area class can be replaced by a number of classes that have these partial areas as their respective extensions. This will ordinarily lead to the situation where several classes in the schema have the same structure, but that is not forbidden by the OODB paradigm. In the following subsection, we present a technique for carrying out this additional partitioning task.

6.1 Partial Areas of a Multi-rooted Intersection Area

It is natural to place roots of a multi-rooted intersection area into different partial areas, since each root represents a distinct semantics. However, concepts may be descendants of more than one root. In such a case, they also represent distinct semantics. Thus, we create new partial areas for these kinds of concepts. Note that these newly created partial areas are considered distinct semantic groups as well. Therefore, concepts which are descendants of the roots of more than one distinct semantic group are also considered to exhibit new semantics in a recursive process.

In order to describe partial areas, we will need some new definitions. We will be using the term “path” to exclusively denote an *upward* path of IS-A links from some concept in the CV to one of its ancestors. The predicate “Desc” will be employed to indicate a descendant/ancestor relationship between a pair of concepts.

That is, $\text{Desc}(\mathbf{x}, \mathbf{y})$ means that \mathbf{x} is a descendant of \mathbf{y} , or, in other words, there exists a path from \mathbf{x} to \mathbf{y} . Two concepts \mathbf{x} and \mathbf{y} are called *independent* if $\mathbf{x} \neq \mathbf{y}$, $\neg \text{Desc}(\mathbf{x}, \mathbf{y})$, and $\neg \text{Desc}(\mathbf{y}, \mathbf{x})$. In the following, the scope of the discussion is a multi-rooted intersection area I . The IS-A links pointing out of I are utilized only at the stage of setting the schema's SUBCLASS_OF relationships.

Definition 10: (Articulation concept): An articulation concept (of intersection area I) is either:

1. (Base case) An intersection concept; or
2. (Recurrence) A concept \mathbf{w} for which there are two independent articulation concepts \mathbf{x} and \mathbf{y} such that: (a) $\text{Desc}(\mathbf{w}, \mathbf{x})$, $\text{Desc}(\mathbf{w}, \mathbf{y})$, and (b) no paths from \mathbf{w} to \mathbf{x} and no paths from \mathbf{w} to \mathbf{y} contain other articulation concepts. \square

The role of an articulation concept in a multi-rooted intersection area is to root and name partial areas, just like the naming concepts were used for areas in the overall CV.

Definition 11: (Direct articulation descendant [DARD]): Let \mathbf{v} and \mathbf{w} be articulation concepts (in I) such that $\text{Desc}(\mathbf{w}, \mathbf{v})$. The concept \mathbf{w} is called a direct articulation descendant (DARD) of \mathbf{v} if there exists a path from \mathbf{w} to \mathbf{v} that does not contain another articulation concept. \square

With the definitions of articulation concept and DARD now in place, we can define the partial areas into which the multi-rooted intersection area is partitioned. Each intersection area will be divided into several partial areas (or *p-areas*, for short).

Definition 12: (P-area): A p-area (within an intersection area I) is a set of concepts containing an articulation concept \mathbf{v} and all of \mathbf{v} 's descendants (within I) excluding its DARDs and their respective descendants. \square

Again, it is important to note that a p-area will contain a single articulation concept which will be the p-area's one and only root. Any descendants that are also articulation concepts will define new p-areas. As with the areas of the CV overall, the root concept is used as the name of the p-area.

Let us now demonstrate the above formalism in the partitioning of two example multi-rooted intersection areas from a CV. The first one is X Area shown in Figure 11, where M Area and N Area are also shown. Note that X Area has three roots. Its p-areas appear in Figure 12.

If there is no overlap among the descendants of the roots (intersection concepts) of a multi-rooted intersection area, then the only articulation concepts are the roots themselves. This is, in fact, the case with the intersection area of Figure 11. In such a situation, every concept within the area is neatly grouped together with the unique root that is its ancestor. As a result, these groups form the p-areas of the original multi-rooted intersection area. Every p-area is singly rooted.

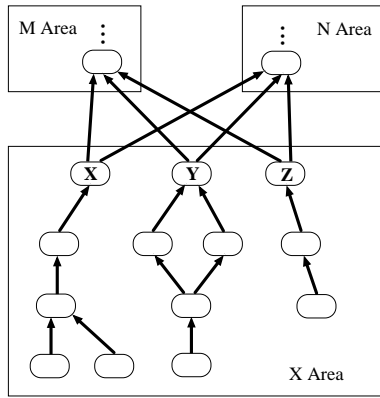


Figure 11: Intersection area with three roots

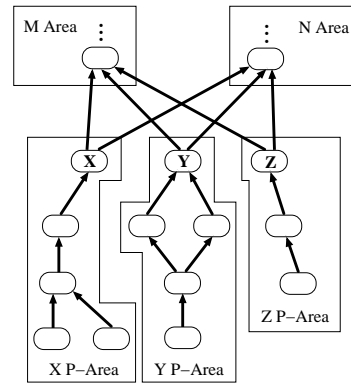


Figure 12: Intersection area's p-areas

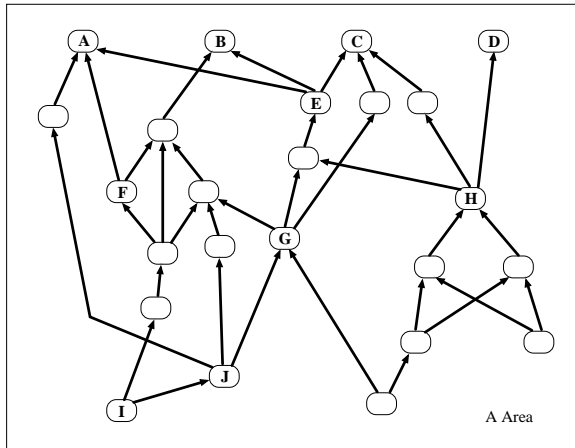


Figure 13: Multi-rooted intersection area with descendant overlap

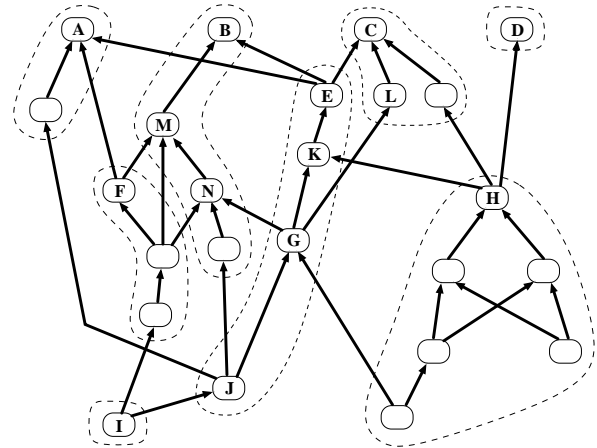


Figure 14: P-areas for intersection area in Figure 13 (with additional concepts labeled)

The partitioning becomes more complex when the descendants of the intersection concepts overlap and create additional articulation concepts. That case is demonstrated by the intersection area A Area shown in Figure 13. This area has four roots: **A**, **B**, **C**, and **D**. By Definition 10, these are articulation concepts. The concepts **E**, **F**, **H**, and **I** are also articulation concepts. The eight p-areas for this intersection area are demarcated in Figure 14 by dashed bubbles.

It is interesting to note that the concept **G** is not an articulation concept even though it has paths without articulation concepts to the independent articulation concepts **B** and **C**. However, the articulation concept **E** lies on a path from **G** to **B** (and also on a path from **G** to **C**). Thus, **G** is not an articulation concept and, in fact, belongs to the p-area rooted at **E**. This follows from item (2b) of Definition 10. For the same reason, the concept **J** is not an articulation concept. It, too, belongs to **E**'s p-area.

In the following, we prove that the p-areas partition the multi-rooted intersection area.

Lemma 6: Let a non-articulation concept \mathbf{v} have multiple articulation ancestor concepts $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$ ($n > 1$). For any i and j such that $i \neq j$, if \mathbf{z}_j is a descendant of \mathbf{z}_i , then \mathbf{v} does not belong to the p-area rooted at \mathbf{z}_i .

Proof: If concept \mathbf{z}_j is a descendant of \mathbf{z}_i , then there must exist a DARD \mathbf{z}_k (possibly \mathbf{z}_j itself) of \mathbf{z}_i on some path between \mathbf{z}_i and \mathbf{z}_j . Concept \mathbf{z}_k is an ancestor of \mathbf{v} since \mathbf{z}_j is an ancestor of \mathbf{v} . By Definition 12, \mathbf{v} will be excluded from the p-area rooted at \mathbf{z}_i as a descendant of the DARD \mathbf{z}_k . ■

Lemma 7: Every concept belongs to at least one p-area.

Proof: Case 1: If \mathbf{v} is an articulation concept, then it is the root of its own p-area. Case 2: Assume to the contrary that a non-articulation concept \mathbf{v} does not belong to a p-area. Since \mathbf{v} is in an intersection area, it must have one or more articulation ancestors $\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_n$ ($n \geq 1$) which are roots of p-areas P_1, P_2, \dots, P_n , respectively. By assumption, \mathbf{v} does not belong to any of the P_i 's. By Definition 12, \mathbf{v} is excluded from P_i ($1 \leq i \leq n$) because there exists a DARD of \mathbf{r}_i (call it $\mathbf{r}_j, i \neq j$) such that $\text{Desc}(\mathbf{v}, \mathbf{r}_j)$. This implies $\text{Desc}(\mathbf{r}_j, \mathbf{r}_i)$. Similarly, we see that there exists an articulation concept \mathbf{r}_k to exclude \mathbf{v} from the p-area P_j . Repeating this n times, we can form a sequence with $n + 1$ articulation concepts starting at \mathbf{r}_i . Each concept in this sequence is a descendant of its predecessor. However, by assumption, \mathbf{v} has only n articulation ancestors. Therefore, some concept must appear more than once in the sequence. This implies there is a cycle in the IS-A hierarchy of the CV—a contradiction. ■

Lemma 8: P-areas are disjoint.

Proof: Case 1: If \mathbf{v} is an articulation concept, then by Definition 12, it is the root of a p-area of its own. Case 2: Assume to the contrary that a non-articulation concept \mathbf{v} belongs to two p-areas X and Y , rooted at \mathbf{r}_X and \mathbf{r}_Y , respectively. Since \mathbf{v} belongs to the p-area X , then by Definition 12 there is no path from \mathbf{v} to \mathbf{r}_X which contains other articulation concepts. Similarly, since \mathbf{v} belongs to the p-area Y , there is no path from \mathbf{v} to \mathbf{r}_Y which contains other articulation concepts. Note that the roots \mathbf{r}_X and \mathbf{r}_Y must be independent concepts. Otherwise, if \mathbf{r}_X is a descendant of \mathbf{r}_Y , then by Lemma 6, \mathbf{v} does not belong to p-area Y rooted at \mathbf{r}_Y . Similarly, \mathbf{r}_Y cannot be a descendant of \mathbf{r}_X . But by Definition 10, if concept \mathbf{v} has two independent articulation concepts ancestors \mathbf{r}_X and \mathbf{r}_Y such that no path from \mathbf{v} to \mathbf{r}_X or \mathbf{r}_Y contains another articulation concept, then \mathbf{v} itself is an articulation concept—a contradiction. ■

Lemmas 7 and 8 together give us:

Theorem 4: The p-areas of a multi-rooted intersection area partition the area. ■

In the following, we present the algorithm that groups concepts in a multi-rooted intersection area I into p-areas. \mathcal{A}_v is a set of concepts within one p-area rooted at v . $\mathcal{A}_{\text{PAREA}}$ is a set of \mathcal{A}_v 's. At the end of the partitioning process, $\mathcal{A}_{\text{PAREA}}$ will be returned. Every concept v has an associated counter for “unprocessed parents” called “p-counter[v].” This algorithm uses one auxiliary function “Num_parents_of” which takes a concept v and a multi-rooted intersection area I as input, and returns the number of v 's parents in I . Every concept v in p-area X rooted at x has an associated variable denoted as $\mathbf{art}[v]$ with value x . Indentation indicates a block. The algorithm does top-down processing to find p-areas in I . A concept can be processed only if all its parents have been processed. Therefore, initially all roots of a multi-rooted intersection area are ready to be processed.

```

set_of_p-areas FUNCTION P-AREA_Partition(area  $I$ )
BEGIN
  // initialization
   $Q := \text{newqueue}()$ ;           //  $Q$  contains concepts ready to be processed.
   $\mathcal{A}_{\text{PAREA}} := \text{newset}()$ ;   //  $\mathcal{A}_{\text{PAREA}}$  will hold all p-areas.
  FOR EACH concept  $v$  in  $I$  DO // Initialize p-counter and  $\mathbf{art}[v]$ .
    p-counter[ $v$ ] := Num_parents_of( $v, I$ );
     $\mathbf{art}[v] := ''$ ;
  FOR EACH root  $v$  of  $I$  DO // Put roots of  $I$  into the queue since
    enqueue( $Q, v$ );       // they are ready to be processed.

  WHILE ( NOT emptyqueue( $Q$ ) ) DO
     $v := \text{dequeue}(Q)$ ;
     $\mathcal{S}_{\text{ART}} := \text{newset}()$ ; //  $\mathcal{S}_{\text{ART}}$  will hold the  $\mathbf{arts}$  of  $v$ 's parents.
    FOR EACH parent  $c$  of  $v$  in  $I$  DO
      insert( $\mathcal{S}_{\text{ART}}, \mathbf{art}[c]$ );
    IF ( |  $\mathcal{S}_{\text{ART}}$  | = 1 ) THEN
      // All parents of  $v$  are in one p-area.
      // Concept  $v$  will be in the same p-area as its parents.
      Let  $u$  be the (only) element in  $\mathcal{S}_{\text{ART}}$ ; // Concept  $v$ 's parents are in the p-area
                                              // rooted at  $u$ .
      insert( $\mathcal{A}_u, v$ ); // Insert  $v$  into set  $\mathcal{A}_u$ .
       $\mathbf{art}[v] := u$ ; // Define  $\mathbf{art}[v]$  to be the articulation concept  $u$ .
    ELSE IF ( |  $\mathcal{S}_{\text{ART}}$  | > 1 AND
       $\exists u$  in  $\mathcal{S}_{\text{ART}}$  which is a descendant of all other nodes in  $\mathcal{S}_{\text{ART}}$  ) THEN
      // see Note 1 below
      insert( $\mathcal{A}_u, v$ );
       $\mathbf{art}[v] := u$ ;
    ELSE
      // Concept  $v$  is an articulation concept.
      // In case |  $\mathcal{S}_{\text{ART}}$  | = 0,  $v$  is an intersection concept (root of  $I$ ).
       $\mathcal{A}_v := \text{newset}()$ ; // Create a new set  $\mathcal{A}_v$  to hold concepts in this p-area.
      insert( $\mathcal{A}_v, v$ ); // Insert articulation concept  $v$  into set  $\mathcal{A}_v$ .
      insert( $\mathcal{A}_{\text{PAREA}}, \mathcal{A}_v$ ); // Insert  $\mathcal{A}_v$  into  $\mathcal{A}_{\text{PAREA}}$ .
       $\mathbf{art}[v] := v$ ; // Define  $\mathbf{art}[v]$  to be itself since  $v$  is an articulation concept.
  
```

```

// After we process concept  $\mathbf{v}$ , we need to decrease the p-counter of  $\mathbf{v}$ 's children by one.
// After the decrease, if any p-counter is equal to zero, we put the associated concept
// into the queue since it is ready to be processed.
FOR EACH child  $\mathbf{k}$  of  $\mathbf{v}$  IN I DO
    p-counter[ $\mathbf{k}$ ]--;
    IF ( p-counter[ $\mathbf{k}$ ] = 0 ) THEN
        enqueue( $\mathcal{Q}$ ,  $\mathbf{k}$ );
RETURN  $\mathcal{A}_{\text{PAREA}}$ ;
END  $\square$ 

```

Note 1: By Theorem 4, \mathbf{v} must belong to one p-area. In this case, \mathbf{v} is not an articulation concept. Concept \mathbf{v} will belong to a p-area rooted at \mathbf{u} which is a descendant of all other elements in \mathcal{S}_{ART} . We can always find such an element \mathbf{u} in \mathcal{S}_{ART} because of the following three conditions: \mathcal{S}_{ART} has more than one element, there are no two independent concepts, and a CV is acyclic. \square

Let us illustrate the partitioning process of an intersection area (Figure 14). Assume the roots of this area were inserted into the queue in the order: \mathbf{A} , \mathbf{B} , \mathbf{C} , and \mathbf{D} . When we process \mathbf{M} (dequeue it from \mathcal{Q}), we will generate a set \mathcal{S}_{ART} with one element $\mathbf{art}[\mathbf{B}]$ (equal to \mathbf{B}). Since there is only one element in \mathcal{S}_{ART} , \mathbf{M} will be inserted into $\mathcal{A}_{\mathbf{B}}$ and $\mathbf{art}[\mathbf{M}]$ will be assigned \mathbf{B} . Next, we decrement the p-counters of \mathbf{M} 's children. \mathbf{F} and \mathbf{N} are inserted into the queue since their p-counters are zero. When we process \mathbf{F} , the \mathcal{S}_{ART} of \mathbf{F} will have two elements, \mathbf{A} and \mathbf{B} , which are $\mathbf{art}[\mathbf{A}]$ and $\mathbf{art}[\mathbf{M}]$, respectively. Since \mathbf{A} and \mathbf{B} are independent, \mathbf{F} is an articulation concept and will be inserted into a new set $\mathcal{A}_{\mathbf{F}}$; $\mathbf{art}[\mathbf{F}]$ will be assigned \mathbf{F} . For \mathbf{G} , the algorithm will generate \mathcal{S}_{ART} with three elements: $\mathbf{art}[\mathbf{N}] = \mathbf{B}$, $\mathbf{art}[\mathbf{K}] = \mathbf{E}$, and $\mathbf{art}[\mathbf{L}] = \mathbf{C}$. Since \mathbf{E} is a descendant of \mathbf{B} and \mathbf{C} , \mathbf{G} is not an articulation concept. Instead, \mathbf{G} will be inserted into $\mathcal{A}_{\mathbf{E}}$, and $\mathbf{art}[\mathbf{G}]$ will be assigned \mathbf{E} .

6.2 Singly-Rooted Schema

After a multi-rooted intersection area is partitioned into its respective p-areas, a separate class (called a *p-area class*) is defined for each of these p-areas. The p-area classes are intended as concept representations that promote better dissemination of the semantics of the underlying CV. To achieve this purpose, each class is singly rooted. The root captures the semantics of the class because all other concepts in the class are its descendents and are thus its conceptual specializations. Therefore, the root serves as a suitable name.

The subclass relationships of a p-area class are defined with respect to the parentage of the (unique) root as for an area class. The schemas for the intersection areas of Figure 11 and Figure 13 are shown, respectively, in Figure 15 and Figure 16. In Figure 15, we see three p-area classes (along with two area

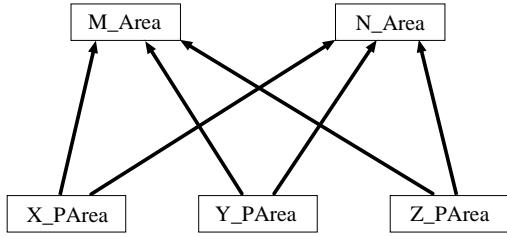


Figure 15: Schema for intersection area in Figure 11

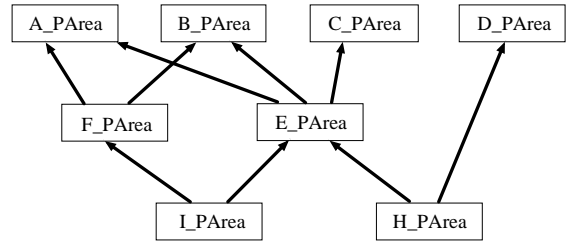


Figure 16: Schema for intersection area in Figure 13

classes) and six subclass relationships. Figure 16 has eight p-area classes and nine subclass relationships. Short-cuts are omitted, as before. Also, in the context of this enhanced schema, we extend the notion of schema-level browsing path (Definition 9) to include p-area classes as well as area classes.

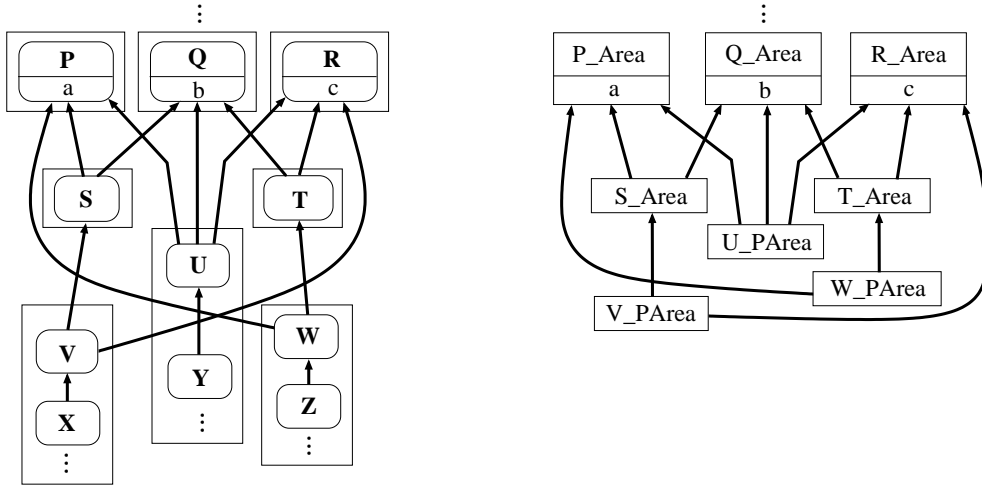


Figure 17: (a) P-areas for Figure 8 (a); (b) Classes for those p-areas

The enhanced partitioning of the multi-rooted intersection area classes solves the problem of establishing informative subclass relationships in the OOVR schema. That is, these subclass relationships more properly reflect the IS-A relationships which cross p-areas of the underlying CV. In Figure 17 (a), we show the p-areas for the CV from Figure 8 (a). Its schema appears in Figure 17 (b). For any concept-level browsing path in Figure 17 (a), there exists a corresponding schema-level browsing path in Figure 17 (b). For instance, for the concept-level path (P, S, V, X), the corresponding schema-level path is (P_Area, S_Area, V_PArea).

The resulting schema is called the singly-rooted schema. By having singly-rooted p-areas, we achieve a schema where each class is named after its unique root and has an extension of semantically uniform

concepts. The name properly captures the contents of the class, and this, in turn, promotes a more accurate abstraction. Additionally, schema browsing is facilitated in all cases, solving the problems of Section 5. The growth of the size of the schema due to the inclusion of the p-area classes is not a concern, as they offer a more refined abstraction of the CV.

Theorem 5: There are no cycles in the singly-rooted schema.

Proof: Suppose to the contrary that the schema contains a cycle of SUBCLASS_OF relationships. Classes represent p-areas or areas. Both kinds are singly rooted. Assume the cycle has $m + 1$ classes Z_0, Z_1, \dots, Z_m such that Z_i is a SUBCLASS_OF Z_{i+1} ($i < m$) and Z_m is a SUBCLASS_OF Z_0 (see Figure 4). Let the roots of Z_0, Z_1, \dots, Z_m be $\mathbf{r}_{z_0}, \mathbf{r}_{z_1}, \dots, \mathbf{r}_{z_m}$, respectively. For each class Z_i ($0 \leq i < m$), there is an IS-A connection from the root \mathbf{r}_{z_i} of Z_i to a concept \mathbf{w}_{i+1} in Z_{i+1} . Also, \mathbf{r}_{z_m} IS-A \mathbf{w}_0 in Z_0 . Since $\text{Desc}(\mathbf{w}_{i+1}, \mathbf{r}_{z_{i+1}})$, then $\text{Desc}(\mathbf{r}_{z_i}, \mathbf{r}_{z_{i+1}})$. Similarly, we see that $\text{Desc}(\mathbf{r}_{z_i}, \mathbf{r}_{z_{i+1}}), \text{Desc}(\mathbf{r}_{z_{i+1}}, \mathbf{r}_{z_{i+2}}), \dots, \text{Desc}(\mathbf{r}_{z_{i-1}}, \mathbf{r}_{z_i})$. But this implies that a cycle of concepts exists in the CV—a contradiction. ■

7 Applying the Singly-Rooted Methodology to a Large CV

In this section, we will apply our modeling approach to an existing OOVr with multi-rooted intersection classes [20]. This OOVr was originally obtained from the MED [6], which contains about 48,000 concepts and 61,000 IS-A links (1996 version). The MED OOVr’s original schema consisted of 90 area classes: 53 property-introducing classes and 37 intersection classes. Of the 37 intersection classes, 14 are multi-rooted. Each of these is listed in Table 1 along with its number of roots and number of constituent concepts. The average size of these intersection classes is 1,975 concepts. The number of concepts belonging to such an intersection class can be quite large: Class 2 contains 29 roots and 19,364 concepts (Table 1).

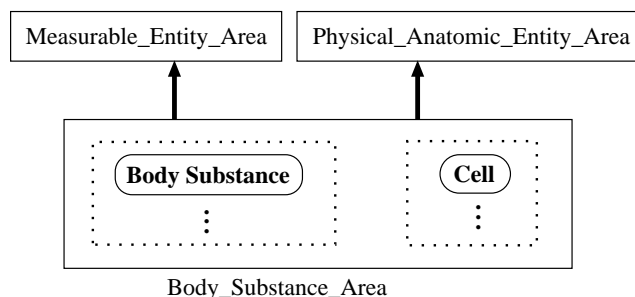


Figure 18: Multi-rooted intersection class with disjoint descendants of roots in original MED OOVr schema

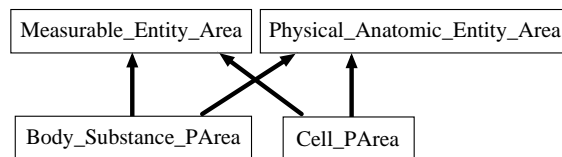


Figure 19: Refined schema from Figure 18

Intersection class	Number of roots	Number of concepts
Intersection class 1	2	106
Intersection class 2	29	19,364
Intersection class 3	965	6967
Intersection class 4	31	324
Intersection class 5	11	182
Intersection class 6	9	170
Intersection class 7	15	40
Intersection class 8	42	43
Intersection class 9	2	175
Intersection class 10	3	3
Intersection class 11	2	52
Intersection class 12	3	96
Intersection class 13	124	124
Intersection class 14	16	16
Total: 14	Total: 1254	Total: 27,662

Table 1: Multi-rooted intersection classes in the original MED OOV schema

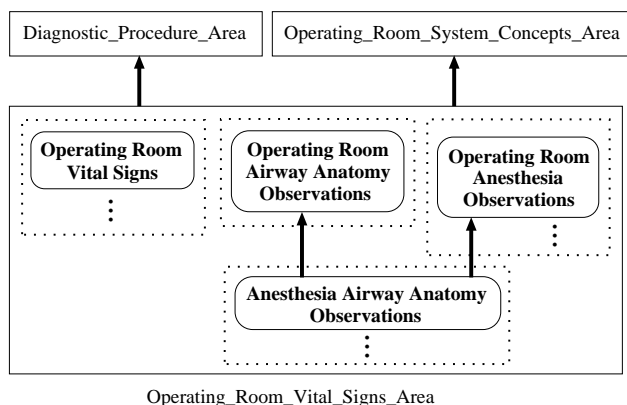


Figure 20: Multi-rooted intersection class with non-disjoint descendants in original MED OOV schema

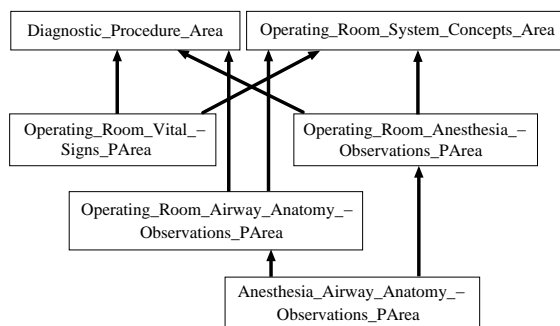


Figure 21: Singly-rooted schema from Figure 20

In Table 1, there are nine classes in which the descendants of the various roots form disjoint sets. For example, Class 1, *Body_Substance_Area*, has two roots, **Body Substance** and **Cell**, and 106 total concepts (Figure 18). Due to the disjointness, the only articulation concepts in the area are the roots. Thus, applying our methodology, we get two p-area classes *Body_Substance_PAria* and *Cell_PAria*, which together replace *Body_Substance_Area*. Both classes have the superclasses *Measurable_Entity_Area* and *Physical_Anatomic_Entity_Area* (Figure 19). The other eight classes are 6, 7, 8, 9, 10, 11, 13, and 14 from the table.

The remaining classes in Table 1 have roots whose descendants overlap. One of them is Class 12, *Operating_Room_Vital_Signs_Area*, which has three roots, **Operating Room Vital Signs**, **Operating**

Intersection class in the original schema	Number of roots	Number of corresponding p-area classes in the singly-rooted schema	Average number of concepts per p-area class
Intersection class 1	2	2	53
Intersection class 2	29	168	115
Intersection class 3	965	1038	7
Intersection class 4	31	48	7
Intersection class 5	11	16	11
Intersection class 6	9	9	19
Intersection class 7	15	15	3
Intersection class 8	42	42	1
Intersection class 9	2	2	87
Intersection class 10	3	3	1
Intersection class 11	2	2	26
Intersection class 12	3	4	24
Intersection class 13	124	124	1
Intersection class 14	16	16	1
Total: 14	Total: 1254	Total: 1489	Average: 19

Table 2: P-area classes in the singly-rooted MED OOV schema

Room Anesthesia Observations, and **Operating Room Airway Anatomy Observations**, and a total of 96 concepts (Figure 20). This class has four articulation concepts, three of them being the roots. The other is **Anesthesia Airway Anatomy Observations**. The four p-area classes which supplant the original intersection class are shown in Figure 21.

Table 2 summarizes the results of applying the revised OOV approach to the MED. Previously, we had fourteen multi-rooted intersection classes with 1,254 roots in total. The new schema contains 1,489 p-area classes. Their average size is nineteen concepts. This should be compared to the average size of 1,975 concepts of the intersection classes that were replaced. This more detailed abstraction level provides a set of smaller and more manageable semantic units and facilitates better navigation of the CV.

8 Conclusions

We have presented a technique which allows a semantic network-based controlled vocabulary (CV) to be converted into an equivalent OODB representation called an OOV. We described the theoretical aspects of our approach and gave an algorithmic specification for its implementation. At its foundation are the notions of area, articulation concept, and p-area, which are used to partition a CV and induce an OODB schema comprising structurally and semantically uniform units. The OODB schema consists of two kinds of classes, area classes and p-area classes. In our development of the methodology, we have solved the difficult problem

of how to formally assign concepts of a multi-rooted area class to a set of singly-rooted p-area classes. Each class in the OOCR schema contains concepts that have the exact same set of properties, making them all structurally uniform. Additionally, every class is singly-rooted and therefore exhibits semantic uniformity.

A major advantage of the OOCR representation is the abstract layer provided by its schema. The singly-rooted schema obtained guarantees that each class comprises a logical unit of concepts. The unique root is used as the name of a class to capture the overarching nature of the class's concepts. Utilizing this abstraction, a user can more readily browse the CV network and comprehend its content. We have presented the results of applying our OOCR approach to an existing CV called the MED. In the future, our methodology may help refine the organization of the metathesaurus [30] of the UMLS by partitioning the sets of concepts associated with semantic types into smaller logical units [15].

References

- [1] American Medical Association, Chicago, IL. *Physicians' Current Procedural Terminology: CPT. 4th ed.*, 1998.
- [2] F. Bancelhon, C. Delobel, and P. Kanellakis, editors. *Building an Object-Oriented Database System: The Story of O₂*. Morgan Kaufmann Publishers, Inc., San Mateo, CA, 1992.
- [3] E. Bertino and L. Martino. *Object-Oriented Database Systems, Concepts and Architectures*. Addison-Wesley Publishing Co., Inc, Redwood City, CA, 1993.
- [4] R. J. Brachman. On the epistemological status of semantic networks. In N. V. Findler, editor, *Associative Networks: Representation and Use of Knowledge by Computers*, pages 3–50. Academic Press, Inc., New York, NY, 1979.
- [5] R. G. G. Cattell and D. K. Barry, editors. *The Object Database Standard: ODMG 2.0*. Morgan Kaufmann Publishers, Inc., San Francisco, CA, 1997.
- [6] J. J. Cimino, P. Clayton, G. Hripcsak, and S. Johnson. Knowledge-based approaches to the maintenance of a large controlled medical terminology. *Journal of the American Medical Informatics Association*, 1(1):35–50, 1994.
- [7] J. J. Cimino, G. Hripcsak, S. B. Johnson, and P. D. Clayton. Designing an introspective, multipurpose, controlled medical vocabulary. In *Proc. Thirteenth Annual Symposium on Computer Applications in Medical Care*, pages 513–517, Washington, DC, Nov. 1989.
- [8] R. A. Côté, D. J. Rothwell, J. L. Palotay, R. S. Beckett, and L. Brochu, editors. *SNOMED International: The Systematized Nomenclature of Human and Veterinary Medicine*. College of American Pathologists, Northfield, IL, 1993.
- [9] S. Even. *Graph Algorithms*. Computer Science Press, Potomac, MD, 1979.
- [10] D. H. Fischer. Consistency rules and triggers for multilingual terminology. In *Proc. TKE'93, Terminology and Knowledge Engineering*, pages 333–342, 1993.
- [11] J. Geller, M. Halper, and Y. Perl. Hybrid diagram/form interface: A two-layered Web-based interface to an OODB vocabulary. In preparation.

- [12] P. M. D. Gray, K. G. Kulkarni, and N. W. Paton. *Object-Oriented Databases: A Semantic Data Model Approach*. Prentice Hall, New York, NY, 1992.
- [13] H. Gu, J. J. Cimino, M. Halper, J. Geller, and Y. Perl. Utilizing OODB schema modeling for vocabulary management. In J. J. Cimino, editor, *Proc. 1996 AMIA Annual Fall Symposium*, pages 274–278, Washington, DC, Oct. 1996.
- [14] H. Gu, M. Halper, J. Geller, and Y. Perl. Benefits of an OODB representation for controlled medical terminologies. *Journal of the American Medical Informatics Association*, 6:283–303, 1999.
- [15] H. Gu, Y. Perl, J. Geller, M. Halper, L. Liu, and J. J. Cimino. Representing the UMLS as an OODB: Modeling issues and advantages. *Journal of the American Medical Informatics Association*, 7(1):66–80, Jan/Feb 2000.
- [16] B. L. Humphreys, D. A. B. Lindberg, H. M. Schoolman, and G. O. Barnett. The Unified Medical Language System: An informatics research collaboration. *Journal of the American Medical Informatics Association*, 5(1):1–11, 1998.
- [17] W. Kim and F. H. Lochovsky, editors. *Object-Oriented Concepts, Databases, and Applications*. ACM Press, New York, NY, 1989.
- [18] D. A. B. Lindberg, B. L. Humphreys, and A. T. McCray. The Unified Medical Language System. *Methods of Information in Medicine*, 32:281–291, 1993.
- [19] L. Liu, M. Halper, J. Geller, and Y. Perl. Controlled vocabularies in OODBs: Modeling issues and implementation. *Distributed and Parallel Databases*, 7(1):37–65, Jan. 1999.
- [20] L. Liu, M. Halper, H. Gu, J. Geller, and Y. Perl. Modeling a vocabulary in an object-oriented database. In *CIKM-96, Proc. 5rd Int'l Conference on Information and Knowledge Management*, pages 179–188, Rockville, MD, Nov. 1996.
- [21] M. E. S. Loomis. *Object Databases: The Essentials*. Addison-Wesley Publishing Co., Reading, MA, 1995.
- [22] E. Mays, C. Apte, J. Griesmer, and J. Kastner. Experience with K-Rep: An object-centered knowledge representation language. In *Proc. IEEE AI Application Conference*, San Diego, CA, Mar. 1988.
- [23] National Library of Medicine, Bethesda, MD. *Medical Subject Headings*. Updated annually.
- [24] N. F. Noy and C. D. Hafner. The state of the art in ontology design: A survey and comparative review. *AI Magazine*, 18(3):53–74, Fall 1997.
- [25] D. E. Oliver, E. H. Shortliffe, and InterMed Collaboratory. Collaborative model development for vocabulary and guidelines. In J. J. Cimino, editor, *Proc. 1996 AMIA Annual Fall Symposium*, page 826, Washington, DC, Oct. 1996.
- [26] ONTOS, Inc. Lowell, MA. *ONTOS DB/Explorer 4.0 Reference Manual*, 1996.
- [27] The OOVr Browser.
URL: <http://object.njit.edu:8080/~newoohvr/JBI/INTERMED/index.html>.
- [28] A. L. Rector, S. Bechhofer, C. A. Goble, I. Horrocks, W. A. Nowlan, and W. D. Solomon. The GRAIL concept modelling language for medical terminology. *Artificial Intelligence in Medicine*, 9:139–171, 1997.
- [29] A. L. Rector, W. A. Nowlan, and A. J. Glowinski. Goals for concept representation in the GALEN project. In *SCAMC'93, the 17th annual Symposium on Computer Applications in Medical Care*, pages 414–418, Washington, USA, 1993.
- [30] P. L. Schuyler, W. T. Hole, M. S. Tuttle, and D. D. Sherertz. The UMLS metathesaurus: Representing different views of biomedical concepts. *Bull Med Libr Assoc*, 81(2):217–222, 1993.

- [31] V. Soloviev. An overview of three commercial object-oriented database management systems: ONTOS, ObjectStore, and O₂. *SIGMOD Record*, 21(1):93–104, Mar. 1992.
- [32] M. S. Tuttle and S. J. Nelson. The role of the UMLS in ‘storing’ and ‘sharing’ across systems. *Int’l J. Bio-Medical Computing*, 34:207–237, 1994.
- [33] United States National Center for Health Statistics, Washington, DC. *International Classification of Diseases: Ninth Revision, with Clinical Modifications*, 1980.
- [34] S. B. Zdonik and D. Maier. Fundamentals of object-oriented databases. In S. B. Zdonik and D. Maier, editors, *Readings in Object-Oriented Database Systems*, pages 1–32. Morgan Kaufmann Publishers, Inc., San Mateo, CA, 1990.

Li-min Liu is a senior member of the technical staff at Pumpkin Networks, Inc. in California. He received his Ph.D. in computer and information science (CIS) from New Jersey Institute of Technology (NJIT) in 1999; M.S. in CIS from Syracuse University in 1994; and B.S. in CIS from Tamkang University in Taiwan in 1990. Dr. Liu is also a research associate of the OODB & AI Laboratory of the CIS Department at NJIT. His research interests include object-oriented (OO) database systems, OO modeling, knowledge representation, medical informatics, controlled vocabularies, and parallel databases.

Michael Halper received the B.S. degree (with honors) in computer science from New Jersey Institute of Technology (NJIT) in 1985; the M.S. degree in computer science from Fairleigh Dickinson University in 1987; and the Ph.D. degree in computer science from NJIT in 1993. During his graduate studies, he was the recipient of a Garden State Graduate Fellowship from the State of New Jersey. Dr. Halper is an associate professor of computer science at Kean University, and a visiting researcher at the OODB & AI Laboratory of NJIT. His research interests include conceptual and object-oriented data modeling, part-whole modeling, extensible data models, object-oriented database systems, and medical informatics. He has worked on the OOHVR project—funded by the National Institute of Standards and Technology (NIST) Advanced Technology Program—to model controlled vocabularies and terminologies using object-oriented database technology. Dr. Halper has had numerous papers in international journals, conferences, and workshops. He is a member of The Honor Society of Phi Kappa Phi.

James Geller received an Electrical Engineering Diploma from the Technical University Vienna, Austria, in 1979. His M.S. degree (1984) and his Ph.D. degree (1988) in Computer Science were received from the State University of New York at Buffalo. He spent the year before his doctoral defense at the Information Sciences Institute (ISI) of USC in Los Angeles, working with their Intelligent Interfaces group. James Geller is currently professor in the Computer and Information Science Department of the New Jersey Institute of Technology, where he is also Director of the OODB & AI Laboratory and Vice Chair of the M.S. and Ph.D. programs in Biomedical Informatics. Dr. Geller has published numerous journal and conference papers in a number of areas including knowledge representation, parallel artificial intelligence, medical informatics, and object-oriented databases. His current research interests concentrate on object-oriented modeling of medical vocabularies, and on Web mining. James Geller is a past SIGART Treasurer.

Yehoshua Perl received his Ph.D. degree in Computer Science in 1975 from the Weizmann Institute of Science, Israel. He was appointed lecturer and senior lecturer in Bar-Ilan University, Israel in 1975 and 1979, respectively. He spent a sabbatical at the University of Illinois in 1977–78. From 1982 to 1985, he was visiting associate professor at Rutgers University (New Brunswick). Since 1985, he has been in the Computer and Information Science Department at New Jersey Institute of Technology (NJIT) where he was appointed professor in 1987. Dr. Perl spent short research visits at: University of Capetown, University of Paris VI, University of Toronto, University of British Columbia, University of Rome, and GMD-IPSI. He received the Harlan Perlis Research Award of NJIT in 1996.

Dr. Perl is the author of more than 80 papers in international journals and conferences. His publications are in the following areas: object-oriented databases, design and analysis of algorithms, data structures,

design of networks, sorting networks, graph theory, and data compression. From 1995 to 1999, Dr. Perl was mainly involved in the OOHVR (Object Oriented Healthcare Vocabulary Repository) project supported by the National Institute of Standards and Technology, Advanced Technology Program. Highlights of his research include among others: the shifting algorithm technique for tree partitioning, analysis of interpolation search, the design of periodic sorting networks, modeling vocabularies using object-oriented databases, and enhancing the semantics of object-oriented databases.