

Lexically-suggested hyponymic relations among medical terms and their representation in the UMLS

Olivier Bodenreider, Anita Burgun, Thomas C. Rindflesch

National Library of Medicine, 8600 Rockville Pike, Bethesda, Maryland, 20894 – U.S.A
{olivier|burgun|tcr}@nlm.nih.gov

Abstract

Objective: Among the various methods for identifying thesaurus relations from text corpora, methods based on head modifier relation are interesting in the context of medical terminologies, especially for those terms which differ from one another by only one modifier. Adjectival modifiers play a particular role because they usually introduce a hyponymic relation. This study focuses on comparing lexically-suggested hyponymic relations among medical terms to inter-concept relationships represented in the Unified Medical Language System (UMLS) Metathesaurus.

Methods: Adjectival modifiers were identified from 63,000 medical terms from the UMLS Metathesaurus, and transformed terms were generated by removing them from the original terms. Candidate hyponymic relations were then tested against inter-concept relationships recorded in the UMLS Metathesaurus.

Results: In 50% of the cases, suggested hyponymic relations were present in the UMLS Metathesaurus. In 25% of the cases, the original term and the transformed terms were “siblings” in the UMLS. In the remaining 25%, no relationship was recorded in the UMLS between these two terms.

The lack of relationships observed in the UMLS Metathesaurus is discussed. Additional methods for automatically assessing the suggested hyponymic relations are proposed. Further research directions are briefly presented.

1. Introduction

Hierarchy is one of the major principles used to structure terminologies. In practice, many terminologies use different kinds of relations to create “hierarchies”, reflecting their organizational principles for a given purpose. Strictly, hierarchy is based on a relation of dominance that comprises the taxonomic relation (‘is a’) and the meronymic relation (‘part of’). Although both hierarchical relations support inferencing, the taxonomic relation is often considered primary due to its conceptual prevalence, and it is commonly represented in terminologies. In many cases, in a given terminology, only some of the taxonomic relationships among terms are represented. In order to augment this information, we propose the use of lexical techniques based on the textual structure of terms but independent of the

organizational structure of the terminology. Specifically, we explore methods for enhancing hyponymic relations, the equivalent of taxonomic relations.

The Unified Medical Language System[®] (UMLS[®]) Metathesaurus[®] is an extensive terminology in the biomedical domain, intended to help health professionals and researchers use biomedical information from disparate sources (Lindberg, Humphreys, & McCray 1993). While the structure of each source vocabulary is preserved, terms that are equivalent in meaning are clustered into a unique concept. Furthermore, inter-concept relationships, either inherited from the source vocabularies or specifically generated, give the UMLS Metathesaurus additional semantic structure. In addition, semantic information such as a semantic group is provided for each UMLS concept.

Although 1,041,938 pairs of hierarchically related concepts are recorded in the Metathesaurus, several studies have shown that numerous relationships are not represented (Bodenreider, Burgun et al. 1998; Cimino 1998; Hole & Srinivasan 2000), limiting, for example, the effectiveness of navigation in the UMLS and the performance of applications based on these relationships (Bodenreider, Nelson, Hole, & Chang 1998). Moreover, the nature of hierarchical relationships in the Metathesaurus is not always made explicit by its constituent vocabularies. Additional relationships acquired independently from the structure of the Metathesaurus would thus provide a means for validating existing relationships, for making precise unspecified hierarchical relations, and for adding relationships not currently represented.

Various methods based on linguistic phenomena have been proposed for automatically acquiring hyponymic relations from texts, in the general context of building ontology from text corpora or for automatic thesaurus construction. Hearst identifies a set of lexico-syntactic patterns that indicate a hyponymy relation (Hearst 1992). For example, in “*such X as Y...*” and “*...X, {and other|or other|including|especially} Y...*”, Y is a hyponym of X. Several authors exploit the semantics of the head modifier relation for detecting term similarities from large corpora (see, for example, Ruge 1997). This method is based on a dependency analysis of the text phrases, in which the head and its modifiers are identified. Terms having many heads and modifiers in common with other terms are usually semantically similar. Relations found among terms in such a set include synonymy, hyponymy-hypernymy and meronymy-holonymy.

Though vocabularies may be extensive, terms rarely contain more than a few words, making the methods based on discourse structures inefficient for identifying hyponymic relations. For example, the five lexico-syntactic patterns indicated by Hearst are found no more than a total of 2534 times among the 1.5 million terms of the UMLS Metathesaurus. Our methodology is related to the work of Ruge; however, instead of using the head modifier relation for identifying semantically similar terms (having many heads and modifiers in common), we propose to take advantage specifically of the property of adjectival modifiers to introduce a hyponymic relation. For example, since the terms “acute appendicitis” and “appendicitis” only differ by the adjectival modifier “acute” modifying the head “appendicitis”, “acute appendicitis” is identified as a candidate hyponym of “appendicitis”. Therefore, (1) terms that contain adjectival modifiers are potential hyponyms, and (2) removing adjectival modifiers from a term T_1 will create a term T_2 in hypernymic relation to T_1 .

The objective of this study is to compare hyponymic relations among medical terms suggested by lexical techniques to inter-concept relationships represented in the UMLS Metathesaurus. The methodology we employ is to isolate a set of Metathesaurus terms

Lexically-suggested hyponymic relations in the UMLS

containing modification. We then call on natural language processing techniques to associate such terms with variants containing no modification. We assume that, overwhelmingly, modified terms are hyponymic, and that, ideally, the Metathesaurus should stipulate a taxonomic relation between terms containing modification and the appropriate unmodified variants. In order to determine to what extent hyponymy is represented in the Metathesaurus, we calculate how many modified terms are in fact related hyponymously with the appropriate unmodified variant. Finally, we provide an initial analysis of some of the missed hyponymy in the UMLS.

2. Material and Methods

The method may be summarized as follows. Starting with a list of terms, a syntactic analysis of the terms allows us to identify adjectival modifiers. Transformed terms are created by removing any combinations of adjectival modifiers from the original terms. The UMLS Metathesaurus is then queried to determine whether a relationship exists between the original term and the UMLS concept to which the transformed term is mapped.

2.1. Material

The UMLS Metathesaurus contains over 1.5 million terms drawn from more than fifty medical vocabularies, and organized in some 730,000 concepts (UMLS 2000). In order to address the large size of the Metathesaurus, we limited our study to terms from SNOMED International, one of the source vocabularies in the UMLS (Côté 1998). We further selected from SNOMED terms from two major components of clinical medicine: diseases and procedures. We also removed from this set all terms containing a comma (10% of our original set). Commas usually signal a permuted form (e.g., “Glucose measurement, urine”) or, more generally, a complex term (e.g., “Patient transfer, in-hospital, unit-to-unit”) whose structure is usually not suitable for natural language processing tools. Our final list contains 63,234 terms (39,075 disease terms and 24,159 procedure terms), corresponding to 42,663 concepts in the Metathesaurus.

2.2. Establishing a list of adjectival modifiers

The study of adjectival modification in the SNOMED terms under consideration was based on an underspecified syntactic analysis (Rindflesch, Rajan, & Hunter 2000) that draws on a stochastic tagger (Cutting, Kupiec, Pedersen, & Sibun 1992) as well as the SPECIALIST Lexicon, a large syntactic lexicon of both general and medical English that is distributed with the UMLS. Although not perfect, this combination of resources effectively addresses the phenomenon of part-of-speech ambiguity in English, and, for example, correctly identifies “open” as an adjective (rather than a verb) in the term “open wound”.

The resulting syntactic structure identifies the head and modifier for the noun phrase analyzed. Each modifier is also labeled as being either adjectival, adverbial, or nominal. Although all types of modification in the simple English noun phrase were labeled, only adjectives and adverbs were selected for further analysis in this study.

Modifiers were identified in 64% of the terms. The number of modifiers per term ranged from one to ten. 89% of the terms with any modification at all were found to have one or two modifiers. 5,416 unique adjectives (62,393 total occurrences) and 69 unique adverbs (509 total occurrences) were extracted from the set of terms. The rationale for extracting adverbs in addition to the adjectives is that in modifying adjectives, adverbs contribute semantically to modification in the phrase.

2.3. Transforming terms

When modifiers were identified in a term O , a set of transformed terms (T_1, T_2, \dots, T_n) was created by removing from term O any combinations of modifiers found in it, whether the syntactic structure of the transformed term is correct or not. The number of transformed terms is $2^m - 1$, m being the number of modifiers. For example, the term “autoimmune hemolytic anemia” contains the two modifiers “autoimmune” and “hemolytic”, so that the following three transformed terms are generated “autoimmune anemia”, “hemolytic anemia”, and “anemia”. 104,199 terms were generated by applying this transformation to our original set.

2.4. Mapping transformed terms to the UMLS

The transformed terms were mapped to the UMLS by first attempting an exact match between the input term and Metathesaurus concepts. If an exact match failed, normalization was then attempted. This process makes the input and target terms potentially compatible by eliminating such inessential differences as inflection, case and hyphen variation, as well as word order variation. One fourth of the transformed terms (26,346) were successfully mapped to a UMLS concept.

2.5. Excluding non-hyponymic relations

In about 10% of the cases, concepts for the original term C_o and for the transformed term C_t did not belong to the same semantic group. For example, “cleft palate”, a disorder, is not semantically compatible with “palate”, a body part. Assuming, arguably, that two terms in hyponymic relation must belong to the same semantic group, we excluded such pairs of concepts from further processing. A similar principle was used for selecting one concept in case of multiple mappings. For example, “renal calculus” is correctly associated with “calculus”, the stone, and not with calculus in mathematics.

2.6. Checking the relations against the UMLS Metathesaurus

For each pair of concepts (C_o, C_t) corresponding respectively to the original term O and to one of the transformed terms T generated from O , the Metathesaurus was queried for relationships:

- C_o is a synonym of C_t (concept identifiers are the same),
- C_o is an ancestor of C_t (‘child of’ or ‘narrower than’ relationships, possibly on more than one generation),
- C_o is a sibling of C_t (the two concepts share a common first-generation ancestor),
- C_o is otherwise related to C_t (‘other’ relationship).

In the UMLS, synonymy is a relation among terms, and synonymous terms are clustered into a concept. For this reason, in this study, a pair of terms is first checked for synonymy. If the two terms are synonymous, no other relationship between terms is represented in the UMLS. The other three kinds of relationships, however, are inter-concept relationships, and, for a given pair of concepts, more than one kind of relationships may be represented in the Metathesaurus. For example, concepts represented as hierarchically related in one vocabulary may be siblings in another vocabulary. Since this study focuses on hyponymy, relationships are searched in the order mentioned above, stopping at the first relationship encountered. If two concepts are both siblings and hierarchically related, only the latter is recorded here. The two concepts are declared unrelated if no relationship is found.

Lexically-suggested hyponymic relations in the UMLS

Modifier	Occur. mappings	Relationship in the Metathesaurus				
		Synonym	Ancestor	Sibling	Other	Unrelated
secondary	925	0%	21%	62%	0%	16%
congenital	903	16%	30%	9%	3%	41%
chronic	795	2%	45%	41%	0%	13%
acute	785	2%	43%	37%	1%	17%
metastatic	752	0%	15%	70%	0%	15%
malignant	531	1%	54%	3%	0%	42%
open	350	1%	35%	34%	0%	30%
closed	298	1%	42%	35%	0%	22%
benign	278	5%	55%	12%	0%	29%
upper	231	1%	38%	38%	5%	17%
acquired	214	8%	46%	20%	3%	24%
primary	205	11%	47%	25%	1%	16%
familial	196	11%	43%	23%	4%	19%
pulmonary	187	3%	57%	11%	4%	26%
partial	181	4%	51%	23%	1%	20%
idiopathic	178	5%	54%	24%	3%	14%
abdominal	167	1%	24%	16%	0%	59%
renal	167	3%	48%	16%	2%	31%
retinal	153	5%	33%	3%	2%	57%
neonatal	146	1%	47%	25%	0%	27%
recurrent	130	2%	31%	54%	0%	13%
bilateral	128	0%	52%	38%	0%	10%
cerebral	126	6%	49%	13%	4%	27%
hereditary	126	6%	56%	18%	0%	20%
cervical	122	0%	35%	42%	1%	22%
peripheral	116	9%	45%	34%	1%	11%
infectious	113	9%	35%	5%	4%	48%
spinal	110	14%	45%	7%	0%	34%
thoracic	110	0%	32%	37%	0%	31%
complete	109	12%	60%	14%	1%	14%
TOTAL	28,851	4%	43%	24%	1%	27%

Table 1 – Distribution of the semantic relations introduced by the 30 most frequent adjectival modifiers, as represented in the Metathesaurus

3. Results

The distribution of the relationships of the original concept (C_o) to the transformed concept (C_t) is given in Table 1 for the most frequently occurring modifiers. For example, 21% of all terms containing the adjective “secondary” are associated as hyponyms with a similar term not containing “secondary”. Under the assumption that (almost) all modified terms in the Metathesaurus should be overtly linked to an unmodified hypernym, Table 1 indicates that a large number of such terms are not linked to the appropriate hypernym. About 60% of the

terms containing “complete” as a modifier are linked to the appropriate unmodified hypernym; however, only 15% of the terms containing “metastatic” are associated with the unmodified hypernym. The last line of the table contains the numbers for the 3607 modifiers and indicates that, overall, more than half of the possible hyponymy links in the Metathesaurus are missing.

Other columns in Table 1 indicate that when hyponymy is not represented, some other relationship often appears. A small number of modified (C_t) and unmodified (C_o) terms are treated as synonyms in the Metathesaurus, while almost a quarter of the total share only a common first-generation ancestor (“siblings”). Finally, no relationship at all is found in the Metathesaurus between the corresponding concepts, C_o and C_t , in roughly another quarter of the cases. In the following section we discuss the etiology of the unmarked hyponymy in the UMLS.

4. Discussion

The major finding in this study is that less than half of the hyponymic relations suggested by lexical techniques are actually represented as hierarchical relationships in the UMLS Metathesaurus. We present an analysis of the causes for “missing” relationships in the UMLS. We then present some common features or patterns observed among missing relationships, that could be used for the automatic validation of lexically-suggested hyponymic relations in the context of the UMLS. Finally, we present some future directions for this work.

4.1 *Hyponymic relations not represented in the UMLS Metathesaurus*

The issue of missing relations in the UMLS Metathesaurus has been often addressed (see, for example, Bodenreider, Burgun et al. 1998; Cimino 1998; Hole & Srinivasan 2000). Here, a manual review of some 15,000 (C_o , C_t) pairs would be necessary to fully evaluate the validity of the hyponymic relations suggested by the presence of adjectival modifiers and not represented in the UMLS Metathesaurus. However, by withdrawing from processing the (C_o , C_t) pairs where the two concepts belong to different semantic groups, this method already provides a mechanism that prevents some false positive hyponymic relations from being identified. A manual review of a limited number of missing relationships in the UMLS Metathesaurus suggests five major causes, often intertwined: a lack of organization within one source vocabulary, a lack of links across sources, underspecified terms, missing synonyms, and the existence of micro-relations.

4.1.1 *Lack of structuration within a source*

By design, some terminologies allow a limited depth for organizing terms. Traditional medical classifications, for example, have a single-tree structure and use the position of the tree for identifying terms, usually with a limited number of digits for the code. As a consequence, terms of differing granularity can appear at the same level of the classification. Figure 1 shows an example of this phenomenon: “acute infantile eczema” is a hyponym of the three terms “acute eczema”, “infantile eczema” and “eczema”. Only the relationship to “disease of the skin and subcutaneous tissues”, provided by SNOMED, is represented in the UMLS for “acute infantile eczema”.

Lexically-suggested hyponymic relations in the UMLS

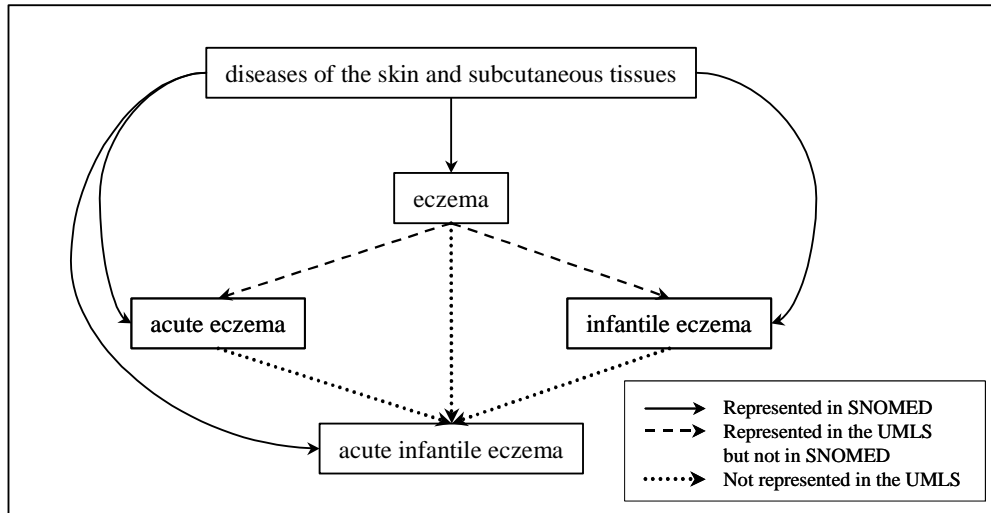


Figure 1 – Hierarchical relationships for the term “acute infantile eczema” in the UMLS.

4.1.2 Lack of links across vocabularies

The UMLS not only merges different vocabularies into a unified structure, but also attempts to link terms across sources both by clustering synonymous terms from various vocabularies into a unique concept, and by creating additional inter-concept relationships.

In the example above (Figure 1), the partially organized list of terms from SNOMED acquires an additional structure through relationships contributed by other source vocabularies or by the UMLS editors, so that “acute eczema” and “infantile eczema” are recorded as hyponyms for “eczema” in the Metathesaurus. In some cases, however, when a specialized term appears only in one vocabulary (e.g. “acute infantile eczema”), it may fail to be linked to some hypernym.

Moreover, some of the source vocabularies in the UMLS provide terms but do not contribute relationships at all, even among their terms. Terms that are synonymous with existing terms are easily integrated. Some 70,000 UMLS concepts, however, remain without any hierarchical relationships.

4.1.3 Underspecified terms

The UMLS Metathesaurus provides several examples of confusion between the generic concept represented by a term T and the most frequent meaning of T. This phenomenon is extremely frequent in the biomedical domain, where numerous modifiers are implicit in medical terms. For example, “hip dislocation” and “acquired hip dislocation” are synonyms in the Metathesaurus while, in fact, hip dislocation may be either congenital or acquired by traumatism, even if the typical, most frequent form for hip dislocation is traumatic. As a result, “congenital hip dislocation” becomes a hyponym of “hip dislocation”, while “acquired hip dislocation” is a synonym of “hip dislocation”. In addition, “congenital hip dislocation” also becomes a hyponym of “acquired hip dislocation”, which is incorrect.

Except for “acute” and “chronic”, differences in the frequency of opposite adjectives confirm the importance of this phenomenon (e.g., “congenital”: 903, “acquired”: 214).

4.1.4 Missing synonymy

The methodology we employ is based on terms, and can suggest a hierarchical relationship between the concepts C_o and C_t only if at least one term O of C_o can be related to at least one term T of C_t . For example, “chronic uremia” and “chronic renal failure” are synonyms, but “hypertensive renal failure” and “chronic hypertensive uremia” have no synonyms. For this reason, this method is able to identify a hyponymic relation between “chronic hypertensive uremia” and “chronic uremia”, but fails to identify the relation between “chronic hypertensive uremia” and “hypertensive renal failure”.

4.1.5 Synonymy versus Hyponymy

Certain hyponymic relations not represented explicitly in the Metathesaurus are lacking due to an interaction between synonymy and hyponymy. In certain cases, the difference between these two phenomena is not absolute. In clear instances of synonymy, the following situation obtains (Cruse 1986): X is a synonym of Y if any proposition P containing X has equivalent truth-conditions to another proposition P' , which is identical except that X is replaced by Y .

A broader conception of the notion of synonymy can be developed that is based on the notion that synonymy is scalar rather than absolute. On this basis synonymous terms are defined by the conjunction of two properties: (1) they manifest a high degree of semantic overlap and (2) they have a low degree of implicit contrastiveness. Since they differ in respect to some semantic traits, a pair of synonymous terms can be incompatible, compatible, or hyponymic/hypernymic. Cruse therefore appeals to the notion of plesionymy, which refers to synonyms that less than absolute.

Plesionyms yield sentences with different truth-conditions, and if the terms are in a hyponymous relation, there may be unilateral entailment. For example “posttransfusion viral hepatitis” and “posttransfusion hepatitis” are considered synonymous in the Metathesaurus, although our processing indicates that they are in a hyponymic relationship. Such a state of affairs suggests plesionymy. In a plesionymous relationship there is one term that it is possible to assert while simultaneously the other term is denied: “it is a posttransfusion hepatitis” but “it is not a posttransfusion viral hepatitis” whereas “it is not a posttransfusion viral hepatitis” implies “it is a posttransfusion hepatitis”. The two concepts evince “capital traits” in common but “posttransfusion viral hepatitis” is a hyponym of the other term, and the relation is called micro-hyponymy.

Within a set of Metathesaurus synonymous terms, several kinds of micro-relations are often represented. Moreover, some items have to be considered synonyms for information retrieval while they must be clearly distinguished for clinical purposes. As a result, in 4% of the (C_o , C_t) pairs, the relation represented in the UMLS is synonymy rather than hyponymy.

4.2 Assessing hyponymic relations not represented in the UMLS

In most cases of missing hyponymic relations, the configurations discussed in figures 2, 3, and 4 indicate that for a given concept, all hierarchical relationships but one (i.e., the dotted line) are represented in the UMLS. This may provide additional clues, or patterns, for automatically assessing the lexically-suggested hyponymic relations. Three common situations are presented.

Lexically-suggested hyponymic relations in the UMLS

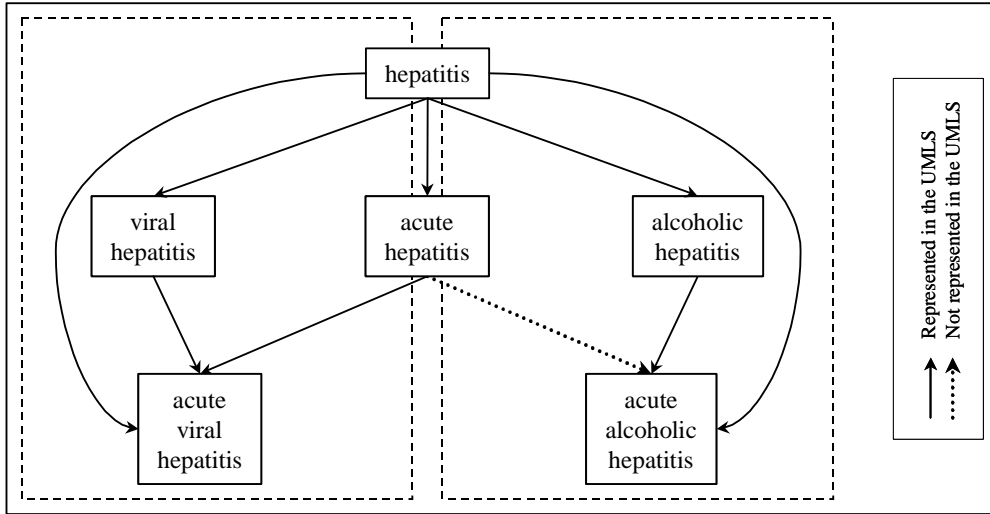


Figure 2 – Missing link: triangular pattern.

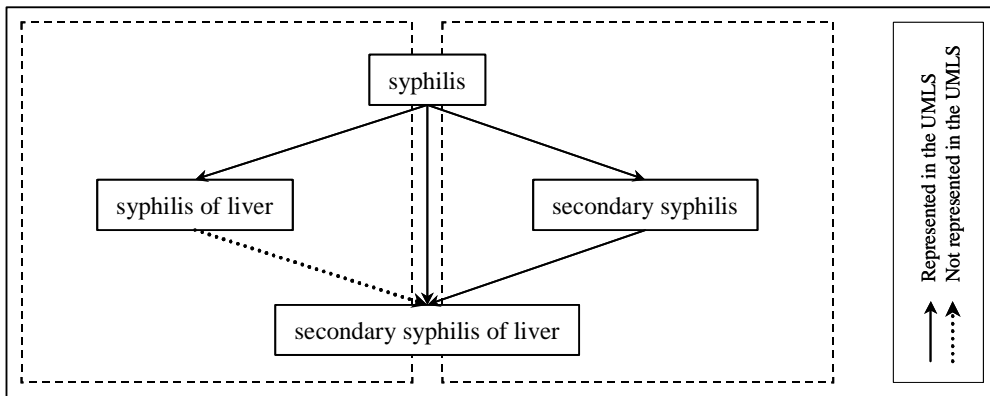


Figure 3 – Missing link: diamond-shaped pattern.

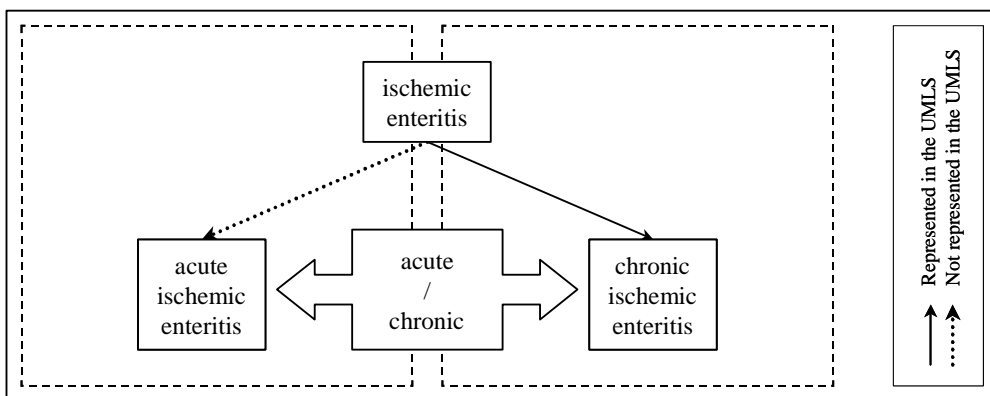


Figure 4 – Missing link: triangular pattern with opposites.

In the first example, the hyponymic relation between “acute alcoholic hepatitis” and “acute hepatitis”, represented in figure 2, is not only suggested through the “acute” modifier, but confirmed by an equivalent mirror-image (triangular pattern), differing only by one modifier (“viral” instead of “alcoholic”). In other words, two orthogonal kinds of hyponyms are derived from “hepatitis” through adjectival modification: “acute” introduces a notion of evolution

over time, while “alcoholic” and “viral” refer to the etiology. When combined, an etiology and an evolution mode for a disease are expected to create terms in hyponymic relation with both of their constituent terms.

Another common pattern (diamond shaped) is presented in figure 3. Here again, the symmetrical representation helps assess the lexically-suggested hyponymic relation. A term (“secondary syphilis of liver”) is a hyponym of three terms, themselves in hyponymic relation (“secondary syphilis”, referring to a phase of the disease, and “syphilis of liver”, referring to its location, are hyponyms of “syphilis”). A term derived from the two hyponyms by combination of the two adjectival modifiers is a hyponym of both terms.

Finally, in figure 4, the context is limited to two terms (“acute ischemic enteritis” and “chronic ischemic enteritis”), only one of them having a direct hypernym (“ischemic enteritis”). It is nevertheless possible to take advantage of the opposition between the two modifiers (“acute” and “chronic”) for assessing the lexically-suggested hyponymic relation between “acute ischemic enteritis” and “ischemic enteritis”, not represented in the UMLS.

4.4 Future directions

The method we propose could help complete the set of hierarchical relationships represented in the UMLS. Lexically-suggested hyponymic relations could, for example, become candidate hierarchical relationships to be reviewed by the UMLS editors.

A more complete set of hierarchical relationships would be useful especially for information retrieval purposes where missing relations are known to lower recall performances. Enhanced knowledge of the role played by adjectival modifiers would also help surgically remove modifiers from queries rather than using more aggressive techniques such as approximate matching.

We plan to further study the patterns of missing relationships for automatically assessing the validity of the relations identified by this method.

Acknowledgements

This research was supported in part by an appointment to the National Library of Medicine Research Participation Program administered by the Oak Ridge Institute for Science and Education through an interagency agreement between the U.S. Department of Energy and NLM.

References

- BODENREIDER O., BURGUN A., BOTTI G., FIESCHI M., LE BEUX P., & KOHLER F. (1998). Evaluation of the Unified Medical Language System as a medical knowledge source. *J Am Med Inform Assoc*, 5(1), pp. 76-87.
- BODENREIDER O., NELSON S. J., HOLE W. T., & CHANG H. F. (1998). Beyond synonymy: exploiting the UMLS semantics in mapping vocabularies. *Proc AMIA Symp*, pp. 815-819.
- CIMINO J. J. (1998). Auditing the Unified Medical Language System with semantic methods. *J Am Med Inform Assoc*, 5(1), pp. 41-51.
- CÔTÉ R. A. (1998). *Systematized Nomenclature of Human and Veterinary Medicine: SNOMED International. Version 3.5*. Northfield, (IL) - Schaumburg (IL): College of American Pathologists - American Veterinary Medical Association.

Lexically-suggested hyponymic relations in the UMLS

- CRUSE D. A. (1986). *Lexical semantics*. Cambridge [Cambridgeshire] ; New York: Cambridge University Press.
- CUTTING D. R., KUPIEC J., PEDERSEN J. O., & SIBUN P. (1992). A practical part-of-speech tagger. *Proceedings of the Third Conference on Applied Natural Language Processing.*, pp. 133-140.
- HEARST M. (1992). Automatic acquisition of hyponyms from large text corpora, *Proceedings of COLING'92* (Vol. 2, pp. 539-545). Nantes (France).
- HOLE W. T., & SRINIVASAN S. (2000). Discovering missed synonymy in a large concept-oriented metathesaurus. *Proc AMIA Symp*, pp. 354-358.
- LINDBERG D. A., HUMPHREYS B. L., & MCCRAY A. T. (1993). The Unified Medical Language System. *Methods Inf Med*, 32(4), pp. 281-291.
- RINDFLESCH T. C., RAJAN J. V., & HUNTER L. (2000). Extracting molecular binding relationships from biomedical text, *Proceedings of the 6th Applied Natural Language Processing Conference* (pp. 188-195). San Francisco: Morgan Kaufmann Publishers.
- RUGE G. (1997). Automatic Detection of Thesaurus relations for Information Retrieval Applications. In C. Freksa & M. Jantzen & R. Valk (Eds.), *Foundations of Computer Science: Potential - Theory - Cognition* (pp. 499-506): Springer.
- UMLS. (2000). *UMLS Knowledge Sources* (11th ed.). Bethesda (MD): National Library of Medicine.