

Technical advance

## A tool for sharing annotated research data: the "Category 0" UMLS (Unified Medical Language System) vocabularies

Jules J Berman\*

Address: Cancer Diagnosis Program, National Cancer Institute, NIH, Rockville, MD, USA

Email: Jules J Berman\* - bermanj@mail.nih.gov

\* Corresponding author

Published: 16 June 2003

Received: 21 January 2003

*BMC Medical Informatics and Decision Making* 2003, **3**:6

Accepted: 16 June 2003

This article is available from: <http://www.biomedcentral.com/1472-6947/3/6>

© 2003 Berman; licensee BioMed Central Ltd. This is an Open Access article: verbatim copying and redistribution of this article are permitted in all media for any purpose, provided this notice is preserved along with the article's original URL.

### Abstract

**Background:** Large biomedical data sets have become increasingly important resources for medical researchers. Modern biomedical data sets are annotated with standard terms to describe the data and to support data linking between databases. The largest curated listing of biomedical terms is the the National Library of Medicine's Unified Medical Language System (UMLS). The UMLS contains more than 2 million biomedical terms collected from nearly 100 medical vocabularies. Many of the vocabularies contained in the UMLS carry restrictions on their use, making it impossible to share or distribute UMLS-annotated research data. However, a subset of the UMLS vocabularies, designated Category 0 by UMLS, can be used to annotate and share data sets without violating the UMLS License Agreement.

**Methods:** The UMLS Category 0 vocabularies can be extracted from the parent UMLS metathesaurus using a Perl script supplied with this article. There are 43 Category 0 vocabularies that can be used freely for research purposes without violating the UMLS License Agreement. Among the Category 0 vocabularies are: MESH (Medical Subject Headings), NCBI (National Center for Bioinformatics) Taxonomy and ICD-9-CM (International Classification of Diseases-9-Clinical Modifiers).

**Results:** The extraction file containing all Category 0 terms and concepts is 72,581,138 bytes in length and contains 1,029,161 terms. The UMLS Metathesaurus MRCON file (January, 2003) is 151,048,493 bytes in length and contains 2,146,899 terms. Therefore the Category 0 vocabularies, in aggregate, are about half the size of the UMLS metathesaurus.

A large publicly available listing of 567,921 different medical phrases were automatically coded using the full UMLS metathesaurus and the Category 0 vocabularies. There were 545,321 phrases with one or more matches against UMLS terms while 468,785 phrases had one or more matches against the Category 0 terms. This indicates that when the two vocabularies are evaluated by their fitness to find at least one term for a medical phrase, the Category 0 vocabularies performed 86% as well as the complete UMLS metathesaurus.

**Conclusion:** The Category 0 vocabularies of UMLS constitute a large nomenclature that can be used by biomedical researchers to annotate biomedical data. These annotated data sets can be distributed for research purposes without violating the UMLS License Agreement. These vocabularies may be of particular importance for sharing heterogeneous data from diverse

biomedical data sets. The software tools to extract the Category 0 vocabularies are freely available Perl scripts entered into the public domain and distributed with this article.

---

## Background

Scientific progress requires the free distribution of research findings. The National Institutes of Health (NIH) has recently stated, in a public notice, the importance of sharing research [1]. "The NIH expects and supports the timely release and sharing of final research data from NIH-supported studies for use by other researchers." [1]. The proposed data sharing policy has arrived at a time when scientists are challenged by confidentiality issues, intellectual property considerations and a variety of federal regulations that constrain the free distribution of research data. In reply to the NIH draft statement, the American Association of Medical Colleges responded, "AAMC believes that effective policies to promote data sharing will require creative, discipline-specific solutions to these complicated problems" [2].

Scientists need innovative methods to overcome the many barriers to data sharing. One of the most important efforts in bioinformatics is data annotation. Data annotation involves appending descriptive information to experimental data with the intention of creating an encapsulated data object that can be intelligently linked or integrated with other data. In genomic databases, annotation of a gene sequence might involve including the physical location in the genome from which the sequence was derived. This might involve using a standard way of describing a location inside a chromosome (e.g. 1p21.3). Annotation might also involve adding information about the translated proteins or the disease that might be associated with a mutation found in the sequence. These annotations are a critical part of the sequence data, because the annotations assist in the discovery of the biological relevance of the sequence. When experimental data are distributed, it is important to distribute the annotation along with the raw data. Unfortunately, many medical vocabularies used for annotation are encumbered under proprietary license agreements. This often means that scientists cannot freely share their annotated biomedical data sets.

The purpose of this article is to provide researchers with a free and legal strategy for sharing research data that has been annotated with a standardized biomedical terminology [3]. Standardized vocabularies are used to annotate text and data with medical terms identified by unique identifier codes that can be parsed by computers and linked to equivalent terms contained in other data sets. For instance, renal cell carcinomas are known by a variety of different terms, including: hypernephroma, clear cell carcinoma of kidney, renal cell adenocarcinoma, rcc

(abbreviated form), Grawitz tumor (eponymous term), etc. Medical nomenclatures typically assign a unique number to "renal cell carcinoma" that is shared by all the synonymous terms. In the case of UMLS, this identifier is C0007134. Some medical nomenclatures include relationships linking a concept to more general (parent) terms and more specific (descendant terms). A data set that annotates each inclusion of "renal cell carcinoma" or "rcc" or "renal cell adenocarcinoma" with a unique concept identifier ensures that different representations of the medical concept can be retrieved and integrated.

## The UMLS Metathesaurus

The UMLS is the largest curated medical nomenclature in existence. It is composed of more than 90 different biomedical vocabularies and contains 2,146,899 million medical terms mapping to 875,255 medical concepts. The UMLS files are available at no cost from the National Library of Medicine's Web site. Anyone wishing to obtain the UMLS metathesaurus must obtain and sign the UMLS License Agreement and register as a UMLS user. The UMLS metathesaurus, the UMLS License Agreement, and detailed instructional documents are available from the following URLs:

License Agreement:

<http://www.nlm.nih.gov/research/umls/license.html>

Current version of Metathesaurus:

<http://www.nlm.nih.gov/research/umls/>

The UMLS License Agreement contains specific language describing usage restrictions. Listed below, using the enumeration and exact text found in the UMLS, are pertinent paragraphs from the 2003 UMLS License Agreement.

1. The NLM hereby grants a nonexclusive, non-transferable right to LICENSEE to use the UMLS products and incorporate them in any computer applications or systems designed to improve access to biomedical information of any type subject to the restrictions in other provisions of this Agreement. The list of licensees authorized to use the UMLS products is public information.
3. LICENSEE is prohibited from distributing the UMLS products or subsets of these products, including individual vocabulary sources within the Metathesaurus, except as an integral part of computer applications developed by

LICENSEE for a purpose other than redistribution of data contained in the UMLS products.

10. LICENSEE shall acknowledge NLM as its source of the UMLS data, citing the year of the UMLS data, in a suitable and customary manner but may not in any way indicate or imply that NLM or any of the organizations whose vocabulary data are included in the UMLS has endorsed LICENSEE or its products.

11. Some of the Material in the UMLS Metathesaurus is from copyrighted sources. If LICENSEE uses any data from the UMLS Metathesaurus:

a) the LICENSEE is required to display in full, prior to providing user access to any Metathesaurus data, the following wording in order that its users be made aware of these copyright constraints:

"Some material in the UMLS Metathesaurus is from copyrighted sources of the respective copyright claimants. Users of the UMLS Metathesaurus are solely responsible for compliance with any copyright restrictions and are referred to the copyright notices appearing in the original sources, all of which are hereby incorporated by reference."

and to display a list of all of the vocabularies obtained from the UMLS Metathesaurus that are used in the LICENSEE's application.

b) the LICENSEE is prohibited from altering data obtained from the UMLS Metathesaurus, but may include data from other sources in applications that also contain UMLS data. The LICENSEE may not imply in any way that data from other sources is part of the UMLS Metathesaurus or of any of its vocabulary sources.

c) the LICENSEE is required to include in its applications identifiers from the UMLS Metathesaurus such that the original source vocabularies for any data obtained from the UMLS Metathesaurus can be determined by reference to a complete version of the UMLS Metathesaurus.

The majority of vocabularies included in the UMLS contain restrictions on their use. An example of these restrictions are transcribed here from the UMLS License Agreement.

#### Category 3 Restrictions:

LICENSEE's right to use material from the source vocabulary is restricted to internal use at the LICENSEE's site(s) for research, product development, and statistical analysis only. Internal use includes use by employees, faculty, and

students of a single institution at multiple sites. Notwithstanding the foregoing, use by students is limited to doing research under the direct supervision of faculty. Internal research, product development, and statistical analysis use expressly excludes: use of material from these copyrighted sources in routine patient data creation; incorporation of material from these copyrighted sources in any publicly accessible computer-based information system or public electronic bulletin board including the Internet; publishing or translating or creating derivative works from material from these copyrighted sources; selling, leasing, licensing, or otherwise making available material from these copyrighted works to any unauthorized party; and copying for any purpose except for back up or archival purposes.

LICENSEE may be required to display special copyright notices before displaying data from the vocabulary source. Applicable notices are included in the list of UMLS Metathesaurus Vocabulary sources, that is part of this Agreement.

Category 3 restrictions prohibit using UMLS-annotated data sets for any of the following purposes: 1) distribution to colleagues, 2) posting on a publicly available site (such as an internet web site), 3) submission to shared data set repositories, or submitted as supplemental data in support of research articles, 4) submission to scientific journals.

The inclusion of Category 3 vocabularies in the UMLS may seem like an egregious oversight. However, the original purpose of the UMLS was to provide a way of mapping between the different vocabularies in the UMLS metathesaurus. Using UMLS, an institution that annotated one database with a proprietary nomenclature and another database with a different proprietary nomenclature, could map equivalent terms in the two databases. In the 1980s, when the UMLS was first released, it was common for large institutions to have site licenses for several proprietary nomenclatures.

In the 1980s, it was not obvious that individual scientists in 2003 would have the technologic facility to create immense data sets annotated with concepts parsed from nomenclatures containing millions of terms. It was not obvious that scientists would seek to share and merge these data sets into immense collections of publicly available data.

The purpose of this article is to provide a simple Perl script specifically designed to extract the Category 0 UMLS vocabularies, and to describe the value of the Category 0 vocabularies as a data sharing tool.

**Table 1: UMLS Category 0 Vocabularies, with the number of terms**


---

|   |
|---|
| AIR – AI/RHEUM – 685 terms  |
| AOD – Alcohol and Other Drug Thesaurus – 20685 terms                            |
| CCS – Clinical Classifications Software – 1608 terms                            |
| COSTAR – Computer Stored Ambulatory Records – 3461 terms                        |
| CSP – Computer Retrieval of Inform. Scientific Projects (CRISP) – 20434 terms   |
| DXP – DxPlain – 10113 terms   |
| HCPCS – Health Care Financing Admin. Common Proc. Coding System – 5091 terms    |
| HL7 – Health Level 7 Vocabulary – 635 terms                                     |
| ICD9CM – International Classification of Disease: 9th Revision – 20069 terms    |
| ICPC2P – International Classification of Primary Care – 13383 terms             |
| ICPCBAQ – ICPC Basque translations – 695 terms                                  |
| ICPCDAN – ICPC Danish translations – 723 terms                                  |
| ICPCDUT – ICPC Dutch translations – 723 terms                                   |
| ICPCFIN – ICPC Finnish translations – 722 terms                                 |
| ICPCFRE – ICPC French translations – 723 terms                                  |
| ICPCGER – ICPC German translations – 723 terms                                  |
| ICPCHEB – ICPC Hebrew translations – 485 terms                                  |
| ICPCHUN – ICPC Hungarian translations – 718 terms                               |
| ICPCITA – ICPC Italian translations – 723 terms                                 |
| ICPCNOR – ICPC Norwegian translations – 722 terms                               |
| ICPCPOR – ICPC Portuguese – 723 terms   |
| ICPCSPA – ICPC Spanish – 723 terms  |
| ICPCSWE – ICPC Swedish – 723 terms  |
| LCH – Library of Congress – 6652 terms  |
| LNC – LOINC Logical Observations Identifiers, Names and Codes – 79522 terms     |
| MCM – Gloss. of Meth. Terms for Clin. Epid. Stud. of Human Disorders – 43 terms |
| MSH – Medical Subject Headings – 516793 terms                                   |
| MTH – UMLS Metathesaurus – 32072 terms  |
| MTHHH – Metathesaurus Hierarchical HCPCS terms – 322                            |
| MTHICD9 – NLM-generated entry terms for ICD-9 I8507                             |
| MTHMST – Metath. Vers. of Min. Standard Terminol Dig. Endoscopy – 1945 terms    |
| MTHMSTFRE – MTHMST French – 1833 terms  |
| MTHMSTITA – MTHMST Italian – 1799 terms   |
| NCBI – National Center for Bioinformatics Taxonomy – 136466 terms               |
| NCI – National Cancer Institute Thesaurua – 2276 terms                          |
| PDQ – Physician Data Query Online System – 19624 terms                          |
| QMR – Quick Medical Reference – 943 terms                                       |
| RAM – Randolph A. Miller Clinically Related Concepts – 258 terms                |
| RXNORM – RxNorm of National Library of Medicine – 33906 terms                   |
| SPN – Standard Product Nomenclature of U.S. Food and Drug Admin – 5064 term     |
| SRC – UMLS Metathesaurus Source Terminologies – 695 terms                       |
| UWDA – University of Washington Digital Anatomist – 79463 terms                 |
| VANDF – U.S. Department of Veterans Affairs – 15795 terms                       |

---

## Methods

The January 1, 2003 release was used to extract the Category 0 vocabularies. Category 0 vocabularies contained in the UMLS are listed in Table 1.

## Software

All of the software scripts used are written in Perl, an open source, freely available cross-platform programming language with interpreters available for virtually every type of computer operating system. Perl can be downloaded from <http://www.activestate.com> or <http://www.cpan.org>.

Detailed information on Perl is available at <http://www.cpan.org> or <http://www.perldoc.com>

The Perl scripts distributed with this article are entered into the public domain by the author and included as a supplemental file with this article.

JHARCOLL, a corpus of medical text was used to obtain an indication of the utility of the Category 0 vocabularies to capture text concepts. JHARCOLL is a public domain collection of 568,035 different phrases extracted from

surgical pathology reports [4]. An example of 19 consecutive lines from the JHARCOLL file is:

- paranuclear staining
- paranuclear vacuoles
- paraortic lymph node
- paraortic lymph nodes
- paraortic mass
- paraortic node
- paraortic nodes
- paraortic region
- paraortic region lymph nodes
- paraosseous hemangioma
- paraovarian adhesion
- paraovarian adhesions
- paraovarian cyst
- paraovarian cyst aspiration
- paraovarian cyst fluid
- paraovarian cyst sac
- paraovarian cyst wall
- paraovarian cystadenofibroma
- paraovarian cysts

The JHARCOLL file is intended to be a corpus of actual phrases found in pathology text that can be used to design and test software that codes, translates, or otherwise manipulates medical text. JHARCOLL is available at:

<http://65.222.228.150/jjb/jharcoll.tar.gz>

The author's Parse module, in Perl, was used to match concepts included in JHARCOLL against concepts from the Category 0 vocabularies as well as from the entire UMLS metathesaurus. The Parse module is described in detail elsewhere [4].

The Parse module is available to the public at:

<http://65.222.228.150/jjb/parse.tar.gz>

The January 1, 2003 UMLS Metathesaurus was downloaded. The two UMLS files needed to extract the Category 0 vocabularies are:

MRCON 151,048,493 1-01-03 7:01 am

MRSO 105,590,732 1-01-03 7:01 am

The UMLS distribution contains Metamorphosys, a Java application designed to extract user-specified vocabularies from the UMLS metathesaurus.

<http://umlsinfo.nlm.nih.gov/MetamorphoSys3.html>

Readers who prefer to use the UMLS extraction tool can input the list of Category 0 vocabularies provided herein, bypassing the GOODCUI.PL Perl script [see Additional file 1]. Those who prefer to use the GOODCUI.PL script should execute the program on a computer with at least 250 MBYTE RAM memory.

**Algorithms**

The GOODCUI.PL script uses two files from the UMLS metathesaurus: MRCON and MRSO. A few records excerpted from MRCON shows terms that map to the UMLS CUI (Concept Unique Identifier) code C0338106, Colon Adenocarcinoma. In this case, as in the case for most CUIs, multiple terms map to the same Concept.

C0338106|ENG|P|L0278121|PF|S0585866|Adenocarcinoma of colon|3|

C0338106|ENG|P|L0278121|VO|S0361300|COLON ADENOCARCINOMA|0|

C0338106|ENG|P|L0278121|VO|S0518037|adenocarcinoma of the colon|0|

C0338106|ENG|P|L0278121|VO|S0766399|colon, adenocarcinoma of the|0|

C0338106|ENG|S|L0493848|PF|S0766389|colon cancer, adenocarcinoma|0|

An excerpt from MRSO lists records that match to the same UMLS Code C0338106, Colon Adenocarcinoma

C0338106|L0278121|S0361300|DXP|SY|NOCODE|0|

C0338106|L0278121|S0518037|PDQ|PT|208/02601|0|

C0338106|L0278121|S0585866|MDR|LT|10001167|3|

C0338106|L0278121|S0585866|RCD|PT|X78gO|3|

C0338106|L0278121|S0766399|PDQ|SY|208/02601|0|

C0338106|L0493848|S0766389|PDQ|SY|208/02601|0|

C0338106|L0493848|S1648401|CCPSS|PT|0005940|3|

Column 4 of the MRSO file contains the abbreviated name of the source vocabulary for each term in the metathesaurus. The Perl script, MRSOCNT.PL, extracts the names of the source vocabularies from the MRSO file and attaches the computed number of terms (defined as unique character strings) contributed by each vocabulary. The list of vocabulary abbreviations is then compared against the published listing of source vocabularies by category restriction, available at the UMLS web site. The composed list of Category 0 vocabularies is used by the GOODCUI.PL script to identify extraction records.

In MRSO, the fourth column of MRSO records lists the vocabulary source while other columns list UMLS codes (C, L and S numbers), specifying a unique record. Starting with a list of all the Category 0 vocabularies and matching against column 4 in each MRSO record, it is easy to create a shortened version of MRSO that consists only of records extracted from the Category 0 vocabularies. Once obtained, each category 0 MRSO record can be matched against the MRCON file, containing the actual vocabulary terms for each set of codes. The GOODCUI.PL script employs this strategy to create a final file (GOODNEW2.TXT) consisting of codes and terms for all UMLS Category 0 records [see Additional file 1].

GOODLST.PL expands and normalizes the GOODNEW2.TXT file, producing a file that consists of term in all lowercase letters, with the "s" truncated from words that end in "s", and with noun and adjectival forms of terms added where they are absent. For example:

Colonic Adenocarcinoma -> colonic adenocarcinoma

colonic adenocarcinomas -> colonic adenocarcinoma

adenocarcinoma of colon -> colonic adenocarcinoma

These transformations, applied to the entire vocabulary, may facilitate implementations of algorithms designed to match free text terms against terms found in the vocabulary.

Perl scripts using the Parse module (GOODCNT.PL and MRCONCNT.PL) examine each of the half-million phrases in JHARCOLL, looking for identical matches from the Category 0 vocabularies or the entire UMLS metathe-

saurus, respectively. These Perl script require the author's freely available Parse module, described in detail elsewhere [4] and distributed as supplemental files with this article.

## Results

There were 43 Category 0 vocabularies in the January 1, 2003 release of the UMLS metathesaurus. Several of the largest and most useful Category 0 vocabularies are:

MESH (Medical Subject Headings) – 516,793 terms

National Center for Bioinformatics Taxonomy – 136,466 terms

LOINC (Logical Observations Identifiers, Names and Codes) – 79,522 terms

MTHICD9 – NLM-generated entry terms for ICD-9 18,507

The extraction script (GOODCUI.PL) created a file containing all Category 0 terms and concepts, GOODNEW2.TXT, in under two minutes using a 1.6 GHz personal computer.

GOODNEW2.TXT is 72,581,138 bytes in length and contains 1,029,161 terms. The UMLS January, 2003 MRCON file is 151,048,493 bytes in length and contains 2,146,899 terms. Therefore the Category 0 vocabularies, in aggregate, are about half the size of the UMLS metathesaurus.

The JHARCOLL file, containing 567,921 different phrases mapped to 1,069,184 UMLS terms and to 711,279 Category 0 vocabulary terms. This indicates that the Category 0 vocabularies matched to about 67% of the terms successfully matched by the complete UMLS.

The number of phrase matches for both the Category 0 vocabularies and for the complete UMLS exceeded the number of different phrases in JHARCOLL. Therefore, on average, JHARCOLL phrases had more than one match against either nomenclature. It was further determined that 545,321 JHARCOLL phrases had at least one match against UMLS terms while 468,785 JHARCOLL phrases matched against Category 0 terms. This indicates that when the two vocabularies are evaluated by their fitness to find at least one term for a phrase, the Category 0 vocabularies are about 86% as useful as the complete UMLS metathesaurus.

## Discussion

Term three of UMLS License Agreement [see Background] prohibits users "from distributing the UMLS products or subsets of these products." This is a reasonable condition, enforcing the National Library of Medicine's role as the

sole curator and distributor of UMLS. However, Term 3 stipulates an exception when the subsets are "an integral part of computer applications developed by LICENSEE for a purpose other than redistribution of data contained in the UMLS products." When a researcher includes a UMLS term as a data set annotation, and distributes the data set to a colleague, her purpose is to disseminate her research. In this case, annotation terms are integrated into the data and do not appear as complete source vocabularies or as subsets of the UMLS that would be suitable as a medical vocabulary. Therefore, distributing data sets annotated with Category 0 terms would not violate the UMLS License Agreement.

For the most part, the Category 0 terms have been contributed directly by U.S. Federal Agencies or by organizations that receive U.S. Federal funds for the purpose of creating publicly available vocabularies. The two largest contributors to Category 0 terms in UMLS are the National Library of Medicine's MESH (Medical Subject Headings) and the NCBI's Taxonomy.

Both MESH and the NCBI Taxonomy vocabularies can be obtained individually from the National Library of Medicine website <http://www.nlm.nih.gov>. Mesh is used by the National Library of Medicine to index all biomedical abstracts included in MedLine, and has been used to index medical terms found throughout the internet [5]. Moore et al have shown that MESH is a useful vocabulary for capturing clinical concepts in surgical pathology reports [6]. MESH alone is a sufficient medical vocabulary for many indexing purposes. The National Center for Bioinformatics Taxonomy contains over 130 thousand terms used to assign standard names for organisms, biological properties, and molecules. Combined, MESH and NCBI Taxonomy would serve to capture most of the concepts included in any biomedical text or data set.

Given that the Category 0 vocabularies are available individually, and at no cost, what is the advantage of using an aggregate nomenclature by combining the Category 0 vocabularies? When data sets are annotated with Category 0 terms, they can be freely shared, and the annotated data can be re-integrated with the FULL set of UMLS knowledge sources, including the Category 1,2 and 3 vocabularies. This is possible because UMLS concept relationships (particularly ancestral and descendant concepts) are always available to UMLS license holders. Because UMLS-coded terms can be related to terms from any UMLS vocabulary, researchers benefit from using the subset of Category 0 UMLS-encoded vocabularies rather than directly employing natively encoded vocabularies (such as MESH or NCBI).

## Conclusions

The Category 0 vocabularies of UMLS constitute a large nomenclature that can be used by biomedical researchers to annotate biomedical data. These annotated data sets can be distributed for research purposes without violating the UMLS License Agreement. These vocabularies may be of particular importance for sharing heterogeneous data from diverse biomedical data sets. The software tools to extract the Category 0 vocabularies are freely available Perl scripts entered into the public domain and distributed with this article.

## Competing Interests

None declared.

## Authors' contributions

This work, consisting of the manuscript and the supplemental Perl scripts, was produced solely by the author as part of his official duties at the U.S. National Institutes of Health.

## Additional material

### Additional File 1

A compressed file containing the Perl scripts used in the article, along with the complete Parse module previously described [4]. The Perl script used to extract the Category 0 terms from UMLS, GOODCUI.PL is included.

Click here for file

[<http://www.biomedcentral.com/content/supplementary/1472-6947-3-6-S1.gz>]

## Acknowledgements

The statements and opinions expressed in the article are those of the author and do not represent the policies of the US. Government or any of its agencies and do not represent legal advice. Users of the UMLS are advised to carefully read the License Agreement and to abide by its wording.

## References

1. **Final NIH Statement on Sharing Research Data** [<http://grants.nih.gov/grants/guide/notice-files/NOT-OD-03-032.html>]. Feb 26, 2003
2. **Comment Letter on NIH Data Sharing Proposal from the American Association of Medical Colleges** [<http://www.aamc.org/advocacy/library/research/corres/2002/051102.htm>]. May 10, 2002
3. Berman JJ: **A Perl script to produce an unencumbered subset of the unified medical language system** Abstract presented at *Advancing Pathology Informatics, Imaging and the Internet* Pittsburgh, PA . Oct 2-4, 2002
4. Berman JJ: **Concept-Match Medical Data Scrubbing: How pathology text can be used in research** *Arch Pathol Lab Med* 2003, **127**:680-686.
5. Malet G, Munoz F, Appleyard R and Hersh W: **A model for enhancing Internet medical document retrieval with "medical core metadata"** *J Am Med Inform Assoc* 1999, **6**:163-172.
6. Moore GW, Miller RE and Hutchins GM: **Indexing by MeSH titles of natural language pathology phrases identified on first encounter using the barrier word method** In: *Computerized Natural Medical Language Processing for Knowledge Representation*

Edited by: Scherrer JR, Cote RA, Mandil SH. Amsterdam: North-Holland; 1989:29-39.

### Pre-publication history

The pre-publication history for this paper can be accessed here:

<http://www.biomedcentral.com/1472-6947/3/6/prepub>

Publish with **BioMed Central** and every scientist can read your work free of charge

*"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."*

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:  
[http://www.biomedcentral.com/info/publishing\\_adv.asp](http://www.biomedcentral.com/info/publishing_adv.asp)

