# How to use MeSH-related Packages

Koki Tsuyuzaki[1,4], Gota Morota[2], Takeru Nakazato[3] and Itoshi Nikaido[4].

January 15, 2016

[1]Department of Medical and Life Science, Tokyo University of Science.
[2]Department of Animal Science, University of Nebraska-Lincoln
[3]Database Center for Life Science, Research Organization of Information and Systems.
[4]Bioinformatics Research Unit, RIKEN Advanced Center for Computing and Communication.

k.t.the-answer@hotmail.co.jp, dritoshi@gmail.com

## Contents

# 1 Introduction

This document provides the way to use MeSH-related packages; *MeSH.db*, *MeSH.AOR.db*, *MeSH.PCR.db*, *MeSH.XXX.eg.db*-type packages, *MeSHDbi*, and *meshr* packages. MeSH (Medical Subject Headings) is the NLM (U. S. National Library of Medicine) controlled vocabulary used to manually index articles for MEDLINE/PubMed [1] and is a collection of a comprehensive life science vocabulary. MeSH contains more than 25,000 clinical and
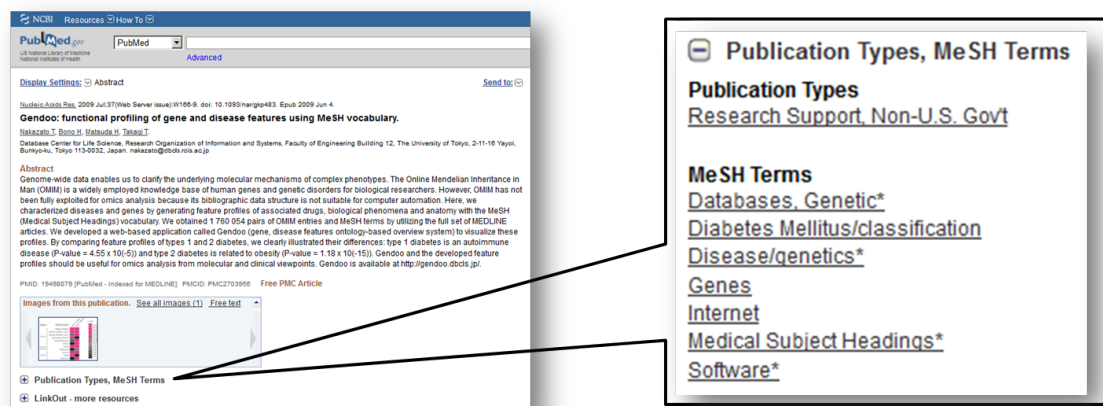


Figure 1: MeSH Term

biological terms. The amount of MeSH term is about twice as large as that of GO (Gene Ontology)[2] and its categories are also wider. MeSH in 2014 proposed its 19 categories and *MeSH.db* provides 16 of them, which are actually assigned to some MeSH terms. Each category is expressed as single capital alphabet as abbreviation defined by NLM. Therefore MeSH is an expected to be much detailed and exhaustive gene annotation tool. Some software or databases using MeSH are now proposed [3, 4, 5, 6].

This vignette introduces R/Bioconductor packages for handling MeSH in R. Original MeSH data is accessible by NLM FTP site (`http://www.nlm.nih.gov/mesh/filelist.html`). The data are downloadable as plain-text format (ASCII MeSH; d2015.bin / q2015.bin). These files were pre-processed by our data-processing pipeline (figure 2) and corresponding information is summarized as a table in SQLite3 file and packed into *MeSH.db*, *MeSH.AOR.db*, and *MeSH.PCR.db*.

| Abbreviation | Category |
|:---:|:---|
| A | Anatomy |
| B | Organisms |
| C | Diseases |
| D | Chemicals and Drugs |
| E | Analytical, Diagnostic and Therapeutic Techniques and Equipment |
| F | Psychiatry and Psychology |
| G | Phenomena and Processes |
| H | Disciplines and Occupations |
| I | Anthropology, Education, Sociology and Social Phenomena |
| J | Technology and Food and Beverages |
| K | Humanities |
| L | Information Science |
| M | Persons |
| N | Health Care |
| V | Publication Type |
| Z | Geographical Locations |

## 1.1 About MeSH

*MeSH.db* provides the corresponding table which contains MeSH ID, MeSH term, MeSH category, synonym, qualifier ID, and qualifier term. Qualifier term means more rough annotation (subheadings) than MeSH. MeSH has hierarchical structure like GO. Such structure is provided as *MeSH.AOR.db* (AOR: ancestor-offspring Relationships) and *MeSH.PCR.db* (PCR: parent-child Relationships as corresponding table.
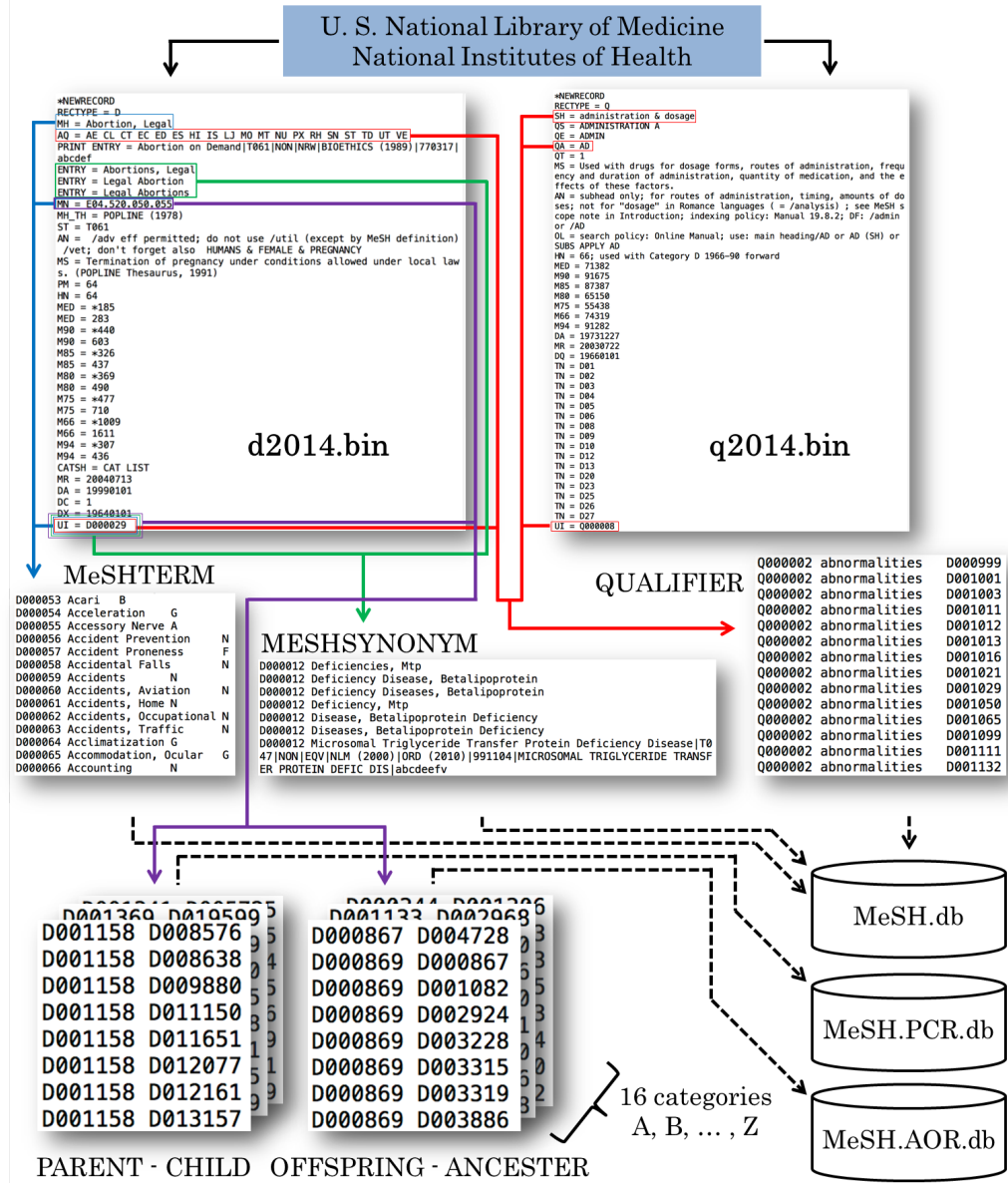


Figure 2: Data pre-process for MeSH.db

## 1.2 The correspondence between MeSH ID and NCBI Entrez Gene ID

MeSH.XXX.eg.db (XXX is an abbreviation of species name such as Hsa: Homo sapiens) packages provide the correspondence between Entrez Gene IDs and NLM MeSH IDs. Such correspondence in wide variety of organisms are summarized as each MeSH.XXX.eg.db by three way of methods, Gendoo[4], gene2pubmed, and RBBH (reciprocal BLAST best Hit).

Gendoo is the web-application based on text-mining of PubMed. Co-occurrence relations in PubMed document are exhaustively retrieved and much relevant correspondence are filtered by some information science techniques.

gene2pubmed is the correspondence between Entrez Gene IDs and NLM PubMed IDs. These relationship is manually assigned by NCBI curator teams. We also summarized the relationship between MeSH Terms and PubMed IDs from licensed-PubMed, then merged as Gene IDs - MeSH IDs correspondence.

For some minor species including non-model organisms, which have no sufficient databases for annotation, we defined 15 well-annotated organisms and 100 minor-organisms, then conducted RBBH between all possible combinations using BLASTP search.

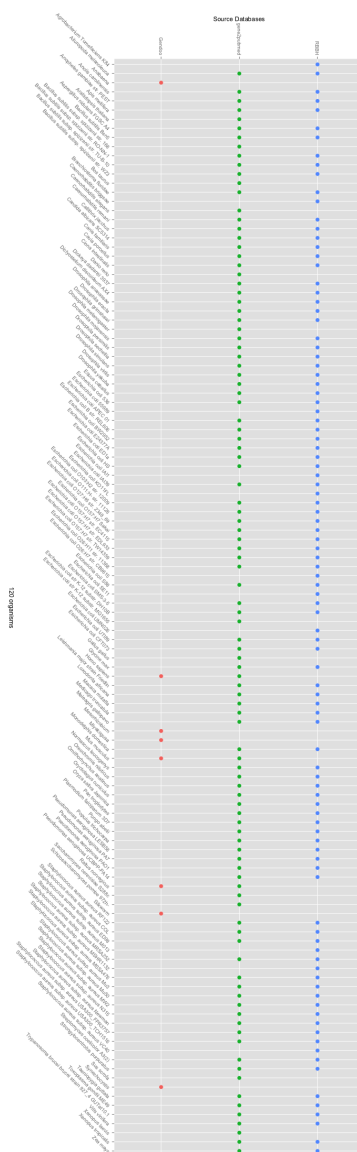| Method | Way of corresponding Entrez Gene IDs and MeSH IDs |
|---|---|
| Gendoo | Text-mining |
| gene2pubmed | Manual curation by NCBI teams |
| RBBH | sequence homology with BLASTP search (E-value $< 10^{-50}$) |

Figure 3: 120 organisms for MeSH.XXX.eg.db and those source databases

## 1.3 Database interface package for MeSH-related packages

We also implemented a database interface (DBI) package named *MeSHDbi*. This package is important because of two reasons. First reason is a unification of DBI functions for MeSH-related packages. *MeSH.db*, *MeSH.AOR.db*, *MeSH.PCR.db*, and $MeSH.XXX.eg.db$ packages inherit the MeSHDb-class defined by *MeSHDbi* and behavior of these packages is uniformly designed. Second reason is supporting construction of user's original MeSH.XXX.eg.db package. Due to the rapid development of DNA sequence technology, wide variety of genome sequences are more and more determined and the correspondence of Gene IDs and MeSH IDs may be designed by many databases [3, 4, 5, 6]. Therefore, we prepared the function to create MeSH.XXX.eg.db package for a situation in which users can retrieved the relationship between Gene IDs and MeSH IDs by some means.

# References

[1] S. J. Nelson and et al. The MeSH translation maintenance system: structure, interface design, and implementation. *Stud. Health Technol. Inform.*, 107: 67-69, 2004.

[2] M. Ashburner and et al. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.*, 25(1): 25-29, 2000.

[3] T. Nakazato and et al. BioCompass: a novel functional inference tool that utilizes MeSH hierarchy to analyze groups of genes. *In Silico Biol.*, 8(1): 53-61, 2007.

[4] T. Nakazato and et al. Nucleic Acids Res. *Gendoo: functional profiling of gene and disease features using MeSH vocabulary.*, 37: W166-W169, 2009.

[5] D. J. Saurin and et al. GeneMeSH: a web-based microarray analysis tool for relating differentially expressed genes to MeSH terms. *BMC Bioinformatics*, 11: 166, 2010.

[6] M. A. Sartor and et al. Metab2MeSH: annotating compounds with medical subject headings. *Bioinformatics*, 28(10): 1408-1410, 2012.