# IMPROVING GENE NAME DICTIONARIES

*Marci Brandenburg, University of Michigan Health Sciences Libraries, Ann Arbor, Michigan*

**M Library**
**UNIVERSITY OF MICHIGAN**

**NCIBI**
**NATIONAL CENTER FOR INTEGRATIVE BIOMEDICAL INFORMATICS**

## BACKGROUND:

The National Center for Integrative Biomedical Informatics (NCIBI) is one of seven National Centers for Biomedical Computing (NCBC) in the NIH Roadmap. The Center is developing the framework of conceptual models, computational infrastructure and an integrated knowledge repository, which modern scientists need to make effective use of the wealth of data flowing from molecular biology and translational research. NCIBI used a process for detecting and selecting gene information from sentences within scientific articles. These sentences were then tagged as referring to a specific gene. However, this gene name curation process failed to tag about 30% of appropriate sentences. NCIBI asked the Health Sciences Libraries for assistance in improving recall.

## OBJECTIVE:

To demonstrate that librarian interaction increases gene name curation efficiency by improving the gene name dictionary currently maintained by the medical center's informatics unit.
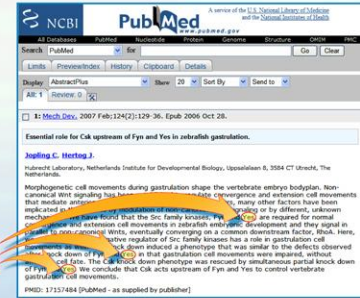
## METHODS:

This project began with a manual review of sentences associated with the genes commonly not tagged and case-by-case identification of synonyms to improve recall. A review of the literature was conducted to determine current practices, such as the creation of rules, mining existing dictionaries, and using natural language processing. Using professional health sciences library experience and subject expertise in biology, I identified potential methods for improving the gene name curation process.



**Yes = common word in the English language**

**Yes = Proto-oncogene tyrosine-protein kinase**

## RESULTS:

Developing a comprehensive gene dictionary is challenging, as each gene name has many synonyms. In addition, it can be difficult to distinguish between gene names and protein names during the curation process. Another challenge is distinguishing between gene name symbols and common words in text. For example, "yes" is a gene symbol for "proto-oncogene tyrosine-protein kinase YES." Therefore, even with an extremely comprehensive gene name dictionary, there are still problems.

There have been a number of methods used to create gene name dictionaries, but none have resolved all of the above problems. A variety of heuristics have been used to eliminate gene or protein name variations and ambiguities. Examples of such heuristics are: using both numeric and Roman numerals for gene names containing numbers (Koike and Takagi), removing Greek letters from names (Zhou, et al.), and converting all letters to lower-case (Mika and Rost). A combination of methods seems to have yielded the best results. However, a comprehensive dictionary for gene name curation is yet to be developed.

| Table 1. Examples of gene names, showing the gene name symbol before and after applying a specific rule. | | | |
|---|---|---|---|
| **Gene Name (Full)** | **Gene Name Symbol (Original)** | **Rule** | **Gene Name Symbol After Applying Rule** |
| Nitric oxide synthase 2, inducible | Nos2 | Use both numeric and Roman numerals | Nos2, NosII |
| Interleukin 1β | Il1β | Remove Greek letters from name | Il1, Il1B, Il1 Beta |
| Tumor necrosis factor | TNF | Convert all letters to lower-case | tnf |

## REFERENCES:

Koike A. and Takagi T. 2004. Gene/protein/family name recognition in biomedical literature. BioLINK 2004: Linking Biological Literature, Ontologies, and Databases. 9-16.

Mika S. and Rost B. 2004. Protein names precisely peeled off free text. Bioinformatics 20(Suppl. 1): i241-i247.

Zhou GD, et al. 2005. Recognition of protein/gene names from text using an ensemble of classifiers. BMC Bioinformatics 6(Suppl. 1): S7-S13.

## FUTURE WORK:

To automate the process of improving gene name dictionaries by creating additional heuristics and using natural language processing.

## ACKNOWLEDGEMENTS: