# Flexible Enrichment with Cortana – Software Demo

**Marvin Meeng**                                                    MEENG@LIACS.NL

LIACS, Leiden University, Niels Bohrweg 1, 2333 CA, Leiden, The Netherlands

**Arno Knobbe**                                                     KNOBBE@LIACS.NL

LIACS, Leiden University, Niels Bohrweg 1, 2333 CA, Leiden, The Netherlands

## Abstract

The software demonstration will introduce an open-source package, called *Cortana*, that simplifies and unifies the procedure of enriching a list, or ranking, of biological entities with background knowledge. A host of software applications already exists that can do this enrichment [2, 5, 9]. However, most are focused on enriching only a single kind of biological entity, eg. genes in the case of gene set enrichment [6], or they are concerned with only a single source of background knowledge, eg. biological processes from GO. As a result, these tools are by design limited in their ability to integrate background knowledge from a variety of domain-crossing sources. The software tool introduced here, *Cortana*, will allow integration of knowledge extracted from both existing sources, like the online knowledge bases [3, 4, 8] that are used by other enrichment tools, as well as custom made ones, created by, or available to, the end user.

## 1. Cortana - Main

The rationale behind using *Cortana* as the basis for the proposed enrichment procedure is that it is a generic data mining tool. As such, it benefits from recent developments in the data mining field. These include, for example, statistically sound validation methods and a range of well-understood, and well-tested quality measures. Finally, it can deal with a variety of data types, including nominal, numeric, ordinal and binary. Fur-

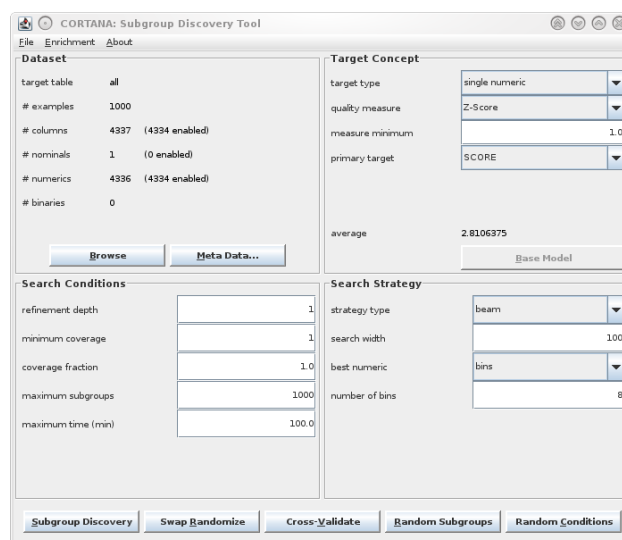*Figure 1.* Main window with parameters set for mining run.

thermore, as it is not purpose-build, it can be used for a wide variety of data mining tasks, biological enrichment being among them, and this allows for many different primary and background sources to be used.

## 2. Cortana - Bioinformatics facilities

To smoothly integrate the enrichment functionality into *Cortana* a 'bioinformatics module' was added, allowing easy deployment of the tool for enrichment tasks. The typical workflow is described in the next section. The end user can use both existing and self defined, or custom, background sources. The tool is made available with domain files from two knowledge bases. The first consist of the three *GO*-domains: *biological process*, *cellular component* and *molecular function*, the well known binary gene/GO-term associations provided by http://www.geneontology.org [3]. The other is a set of so-called *association-matrices* derived from the literature by means of text mining. To

create these, a set of common concepts is determined through text mining of a large corpus of PubMed articles [1, 7], and then association scores between these concepts are calculated and stored in a matrix, or actually multiple matrices, each one specific to a certain domain. Alternatively, a user can use custom background sources. However, in this case, the end user will also need to supply a mapping file, allowing the entities in the input ranking to be matched to the entities in a background source.

## 3. Cortana - Workflow

This section describes the typical workflow when enriching a (gene-)ranking with background knowledge. One starts by opening a file containing a ranking of biological entities. Typically this will be the result of an analysis of data obtained in a microarray experiment, yielding a list of *differentially expressed genes*, although other



Figure 2. Selecting a domain to add to ALL gene ranking.

types of biological entities (eg. proteins) can be dealt with analogously. Figure 3 shows part of the ranking used in this example, with details such as the ENTREZ-identifier of the gene in question, its score resulting from the primary data analysis, its resulting rank, and the gene-symbol. Figure 1 shows *Cortana*'s main window, in which some information about the loaded data is displayed, such as the number of genes in the ranking, and the total number of descriptive values per gene that is available for enrichment. After loading the ranking, one or more files containing background knowledge from various domains can be included in the analysis. This is achieved by pressing 'Add CUI Domain', which will show a list of available CUI domains, as seen in Figure 2. Note that adding multiple background sources is possible. If one uses a custom background source, one will be asked to also indicate which file to use, to map the entities in the ranking to the corresponding information in the background source.
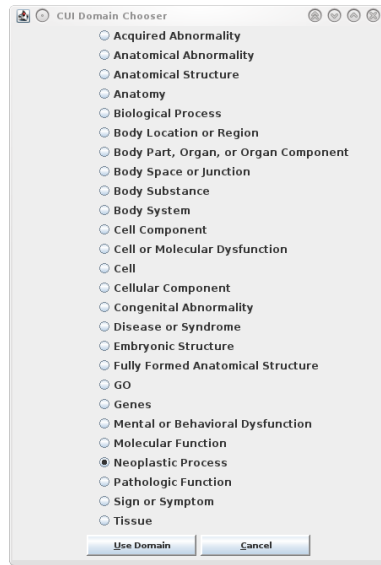


Figure 3. Browse combined table of ALL gene ranking and domain (*neoplastic process* in this case).



Figure 4. Result of enriching ALL gene ranking with Neoplastic Process.

*Cortana* will now combine the original data with the selected domain(s), the result of which can be checked by selecting 'Browse' (see Figure 3) or 'MetaData'. Then, one sets the parameters to be used for the mining run, and starts the actual enrichment by pressing 'Subgroup Discovery', again see Figure 1. After the mining finishes, a result window like Figure 4 will be shown, listing those subgroups, or concepts, that were found to be interesting, according to the selected parameters. In this case, setting the search parameter *refinement depth* to *2*, see Figure 1, allowed two conditions to be combined to form subgroups.

## 4. Benefits of the presented tool

*Cortana* is a flexible data mining tool for exploratory data analysis. With the addition of a bioinformatics module, its range of generic data mining facilities can now be employed for enrichment in the biological and medical domain. The tool comes with a large collection of background knowledge that can be involved in the enrichment process. This collection includes the customary functional annotations from GO (binary concepts), as well as the association matrices describing the level of association between each entity

in the ranking and concepts from a number of domains (numeric concepts). In total, the 28 background files comprise 26.8 GB of background information available for enrichment. The bioinformatics module of *Cortana* has been employed in a number of medical and biological applications, including enrichment of neuroblastoma gene rankings, leukaemia (ALL and AML) gene rankings, and metabolic syndrome gene and metabolites rankings.

# References

[1] Jelier, Schuemie, Veldhoven, et al., *Anni 2.0: a multipurpose text-mining tool for the life sciences*, Genome Biology 2008, 9(6):R96.

[2] Glynn Dennis Jr. et al., *DAVID: Database for Annotation, Visualization, and Integrated Discovery*, Genome Biology, 2003 4(9):R60, `http://david.abcc.ncifcrf.gov/home.jsp`.

[3] The Gene Ontology, 2011, `http://www.geneontology.org`.

[4] Ensembl, 2011, `http://www.ensembl.org`.

[5] Zeeberg B.R. et al., *GoMiner: A Resource for Biological Interpretation of Genomic and Proteomic Data*, Genome Biology, 2003 4(4):R28, doi:10.1186/gb-2003-4-4-r28, `http://discover.nci.nih.gov/gominer/index.jsp`.

[6] Subramanian, et al., *Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles*, PNAS October 25, 102:43, 2005.

[7] Jelier, Goeman, Hettne, Schuemie, den Dunnen & 't Hoen, *Literature-aided interpretation of gene expression data with the weighted global test*, Briefings in Bioinformatics, doi:10.1093/bib/bbq082, 2010.

[8] KEGG: Kyoto Encyclopedia of Genes and Genomes, 2011, `http://www.genome.jp/kegg/`.

[9] Search for Enriched Gene Sets, 2011, `http://kt.ijs.si/software/SEGS`.

[10] Trajkovski, Lavrač & Tolar, *SEGS: Search for enriched gene sets in microarray data*, Journal of Biomedical Informatics, 41(4), 588–601, 2008, `http://dx.doi.org/10.1016/j.jbi.2007.12.001`.