

SemNExT: A Framework for Semantically Integrating and Exploring Numeric Analyses

Evan W. Patton,¹ **Elisabeth Brown**,² **Matthew Poegel**,²
Hannah De los Santos,² **Chris Fasano**,³ **Kristin P. Bennett**,²
² and **Deborah L. McGuinness**¹

¹ Dept. of Computer Science, RPI

² Dept. of Mathematical Sciences, RPI

³ Neural Stem Cell Institute



Overview

- Motivation & Domain
- Example
- Architecture
 - Datasets
 - Statistics
 - Linking
 - Provenance
- Conclusions



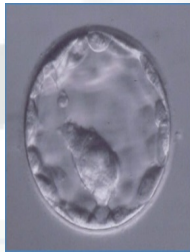
Motivation

- Neural Stem Cell Institute collects significant amounts of data on the state of brain development in the form of *gene reads*
- Questions:
 - How to *cluster genes* together based on *activity over time*?
 - How to associate those gene clusters with *diseases of interest/gene ontology annotations*?
 - Are there other relationships that may not be obvious from the underlying data alone?

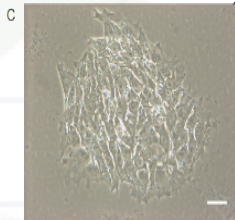


Understanding origins of disease in human cerebral cortex development

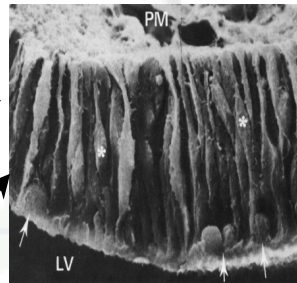
Fertilization



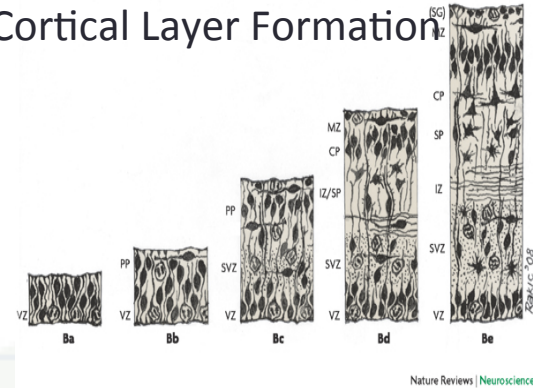
Human Embryonic Stem Cells



Neural Differentiation



Cortical Layer Formation



- Create molecular signature of normal cortical development *in a dish*
- Understand stages of development and when layers form
- Analyze mutated genes associated with disease to understand developmental origins
- Understand which pathways are associated with stages and diseases.

Joint work on CORTECON Data with Dr. Chris Fasano and Sally Temple at Neural Stem Cell Institute



Domain

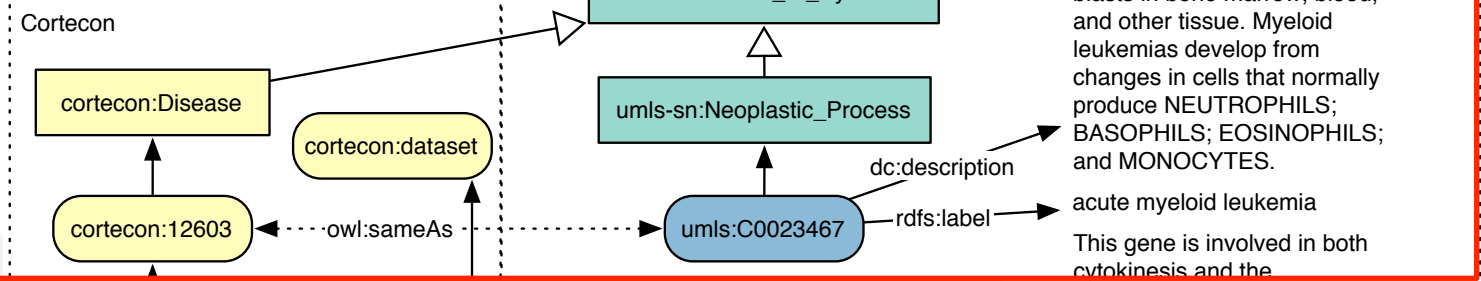
- Brain development data, primarily focused on ***RNA-seq counts*** from RNA sequencing
- ***Genes*** encode ***proteins*** and ***enzymes***, which perform various biological functions
- Proteins often act as part of ***pathways*** and ***interact*** with one another, may cause certain ***diseases***, and be affected by ***compounds***, such as ***drugs***

*For more info, see J. van de Leemput et al. "CORTECON: a temporal transcriptome analysis of in vitro human cerebral cortex development from human embryonic stem cells." *Neuron*: 83.1 (2014): 51-68.*

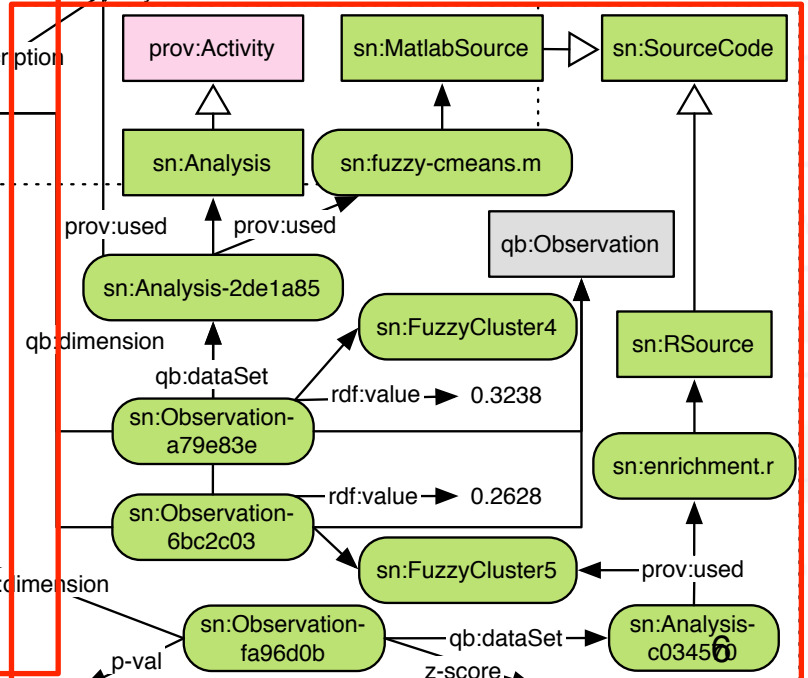
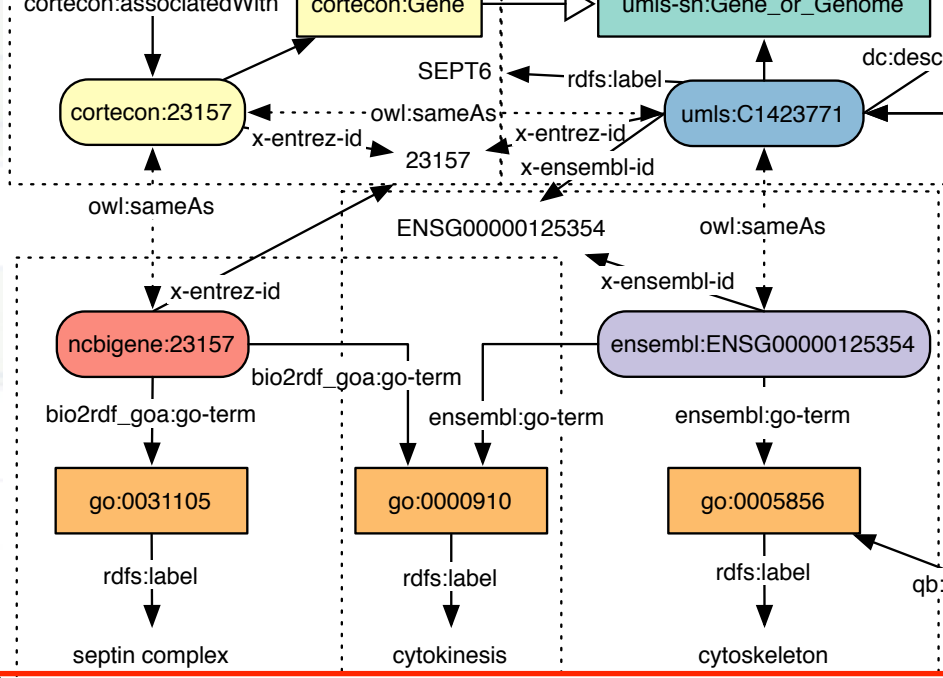


Example

Disease



Gene



Analysis

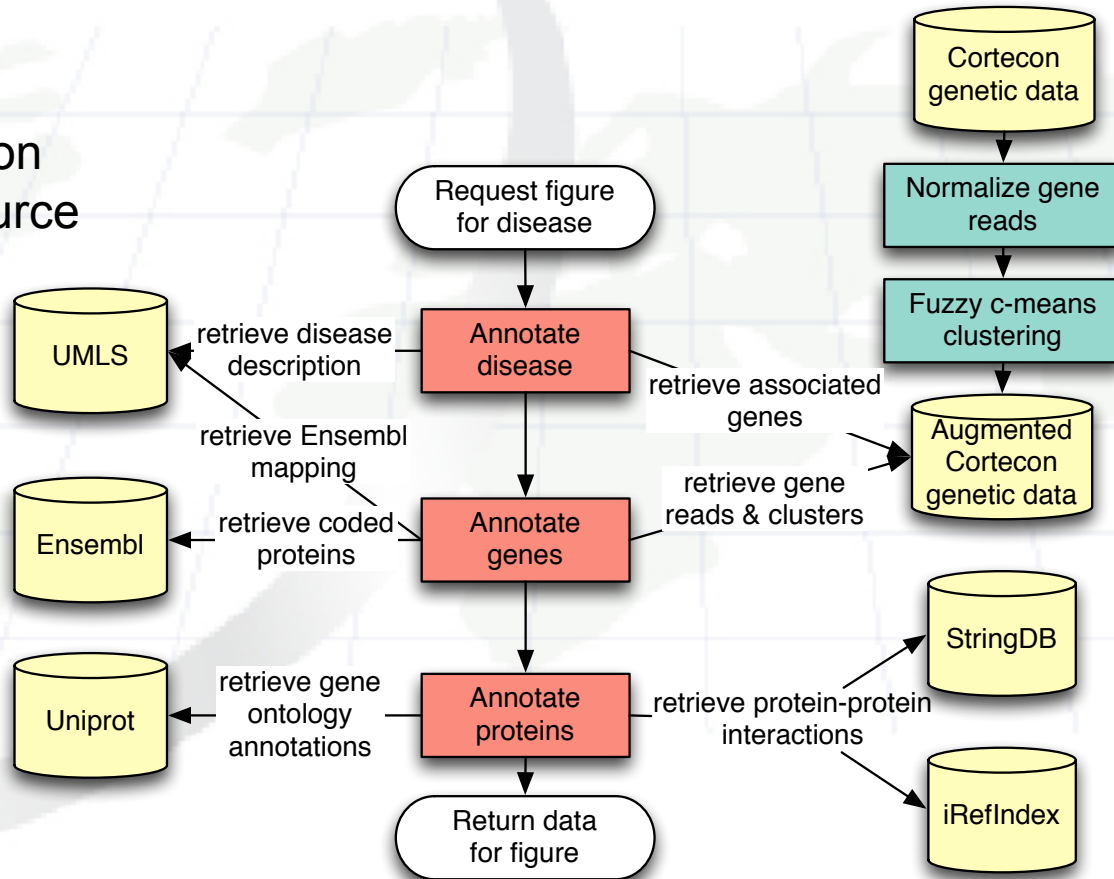




Architecture

SemNeXT annotates **entities** from **data sources**, performs **analyses**, and visualizes the results as **Chord and Heat Map (ChEM)** diagrams.

- Analysis
- Annotation
- Data Source





Datasets

Relational Databases

1. NSCI Cortecon
 - Genes, diseases, RNAseq read data
2. KEGG
 - Genes, pathways, proteins
3. StringDB
 - Protein interactions
4. Ensembl
 - Genes, proteins, gene ontology annotations
5. Unified Medical Language System
 - Genes, diseases

RDF Quad Stores

6. Bio2RDF
7. ReDrugS (RPI)
 - IRefIndex
 - Proteins, protein interactions
 - Drugbank
 - Proteins, drugs
 - Online Mendelian Inheritance in Man
 - Diseases, genes
8. Uniprot
 - Genes, proteins, protein interactions, gene ontology annotations

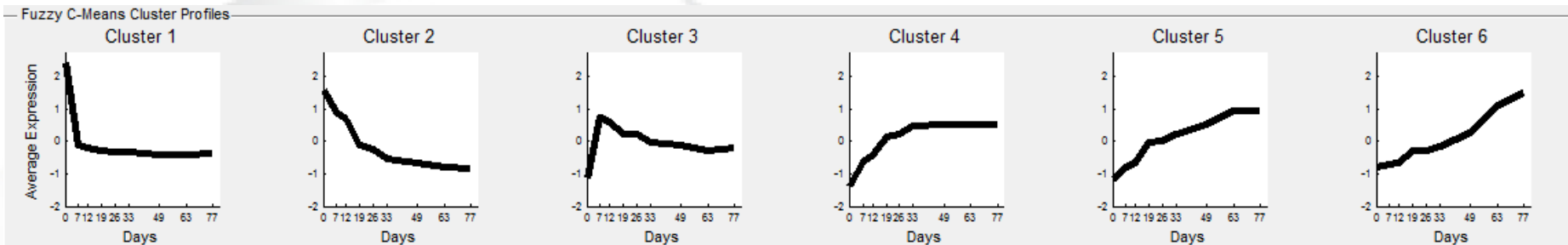


Statistics

- We used ***fuzzy c-means clustering*** and ***singular value decomposition*** to cluster genes and order them based on activation
 - The fuzzy c-means clustering technique identified one more gene cluster than the CORTECON analysis
 - What is significant about this cluster?
- ***Enrichment analysis*** was used to determine when a disease/pathway is enriched or depleted for a particular set of genes
- Statistical analyses are made available at dereferenceable, content-negotiable URIs



Extracted Brain Development Clusters



Placenta

Epithelial

Brain

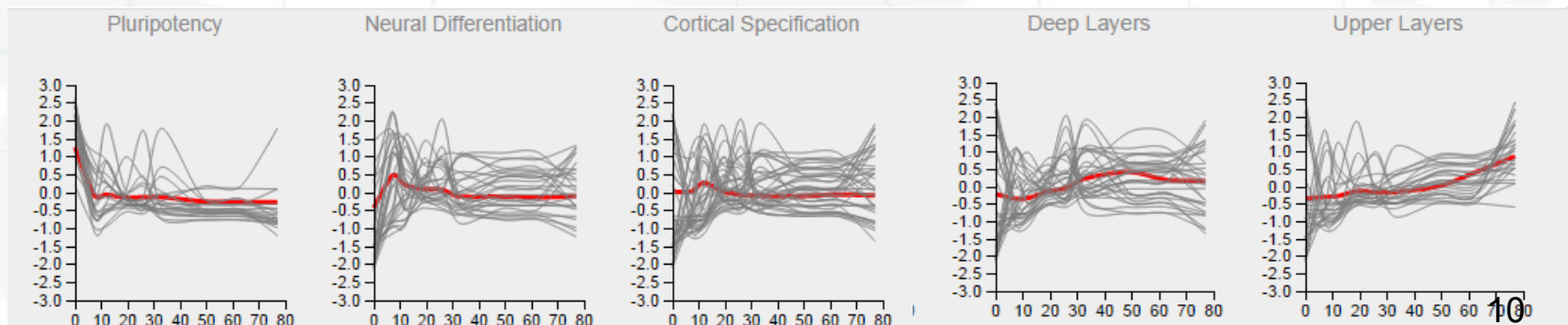
Brain

Brain

Brain

Tissue Enrichment

Stages from Cortecon <http://cortecon.neuralsci.org/>



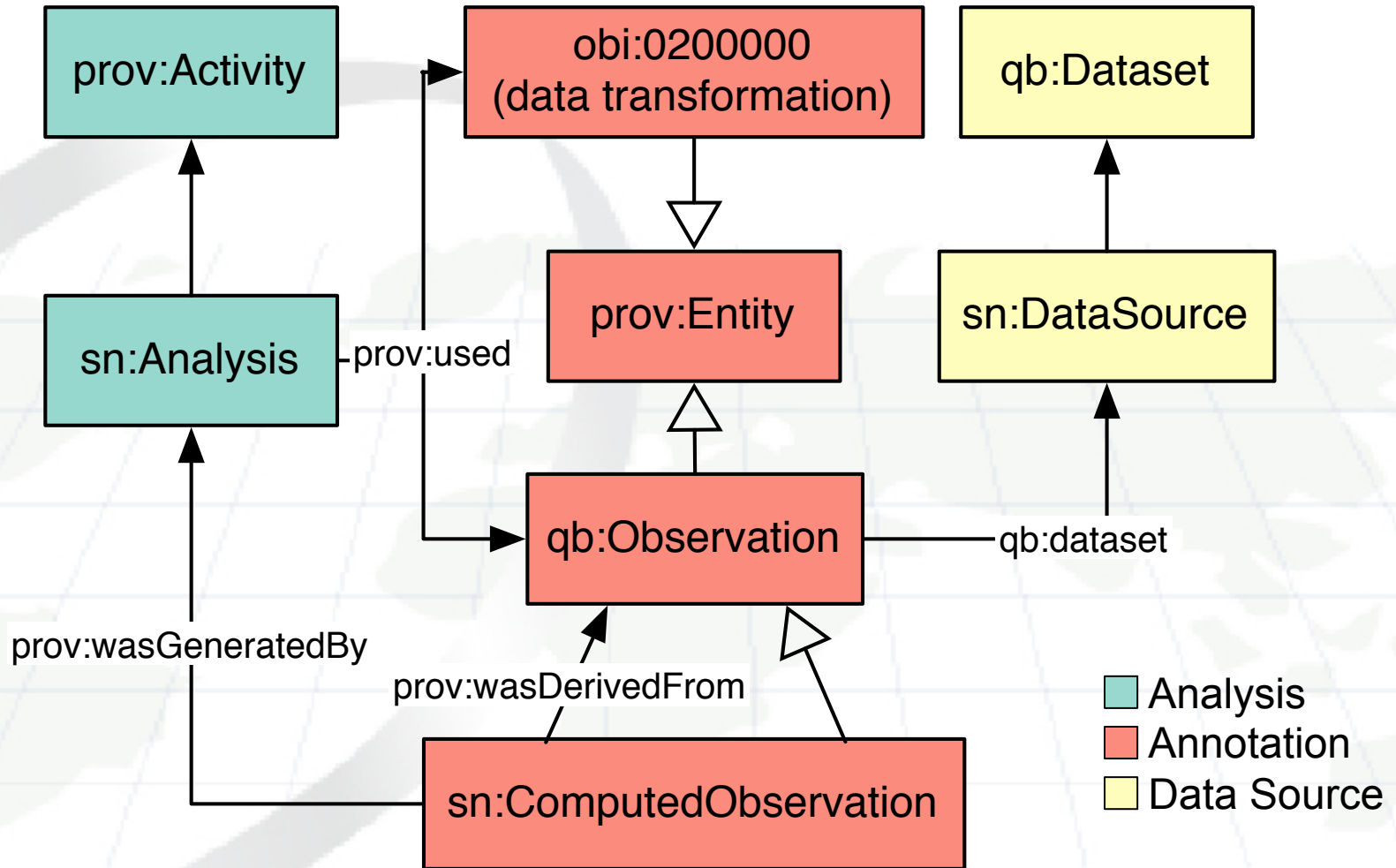


Linking

- Composed of two techniques:
 - Cross-database identifier linking:
 - Typical identifiers include Gene Ontology, NCBI, KEGG, Ensembl identifiers
 - Textual matching in the absence of unique identifiers:
 - “Alzheimer’s disease, familial, 1”, “Alzheimer’s disease, familial, 2”, etc. => “Alzheimer’s disease”
 - Exploit broader and type relations to collapse results into a “summary” entity

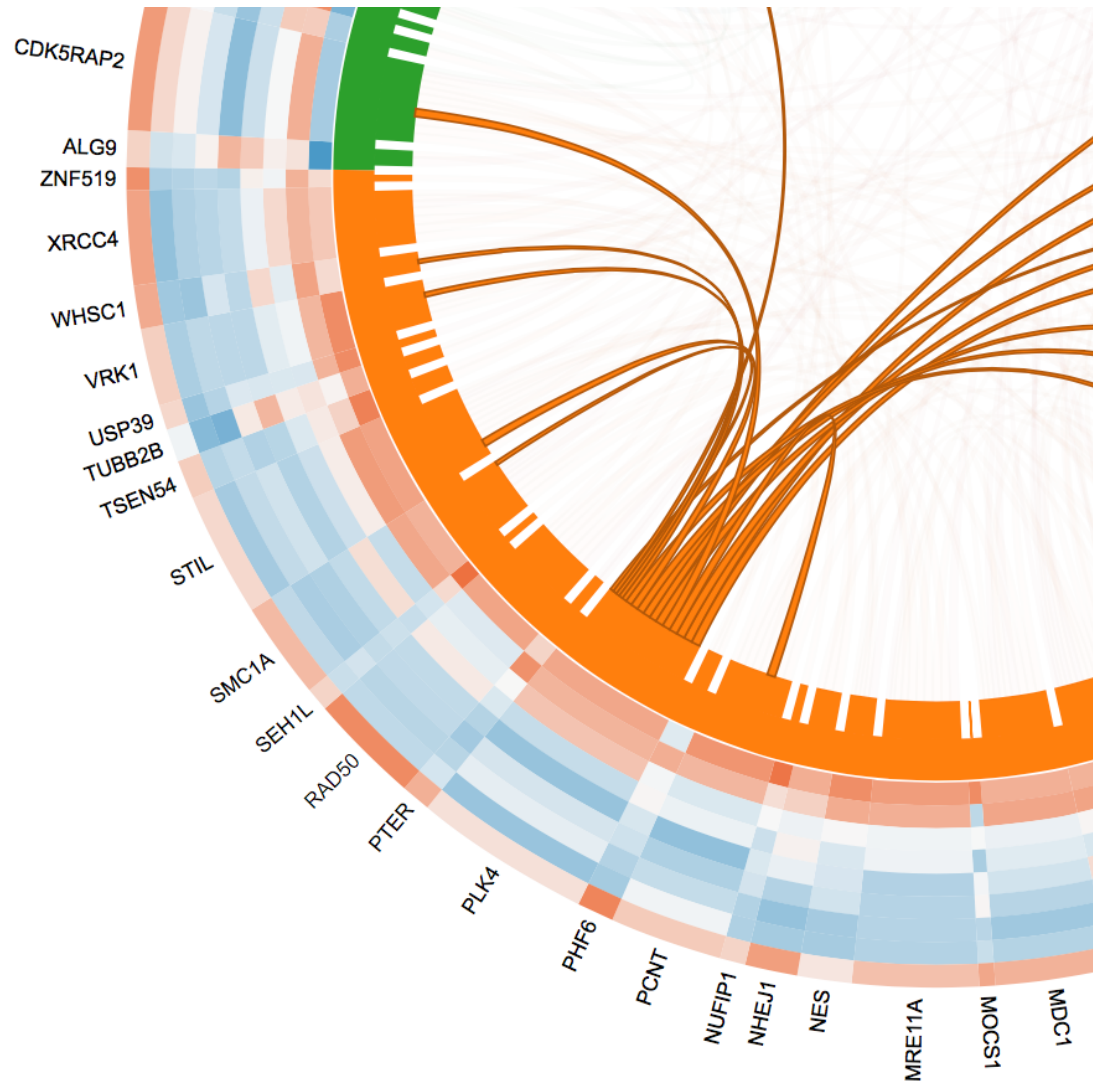


Provenance



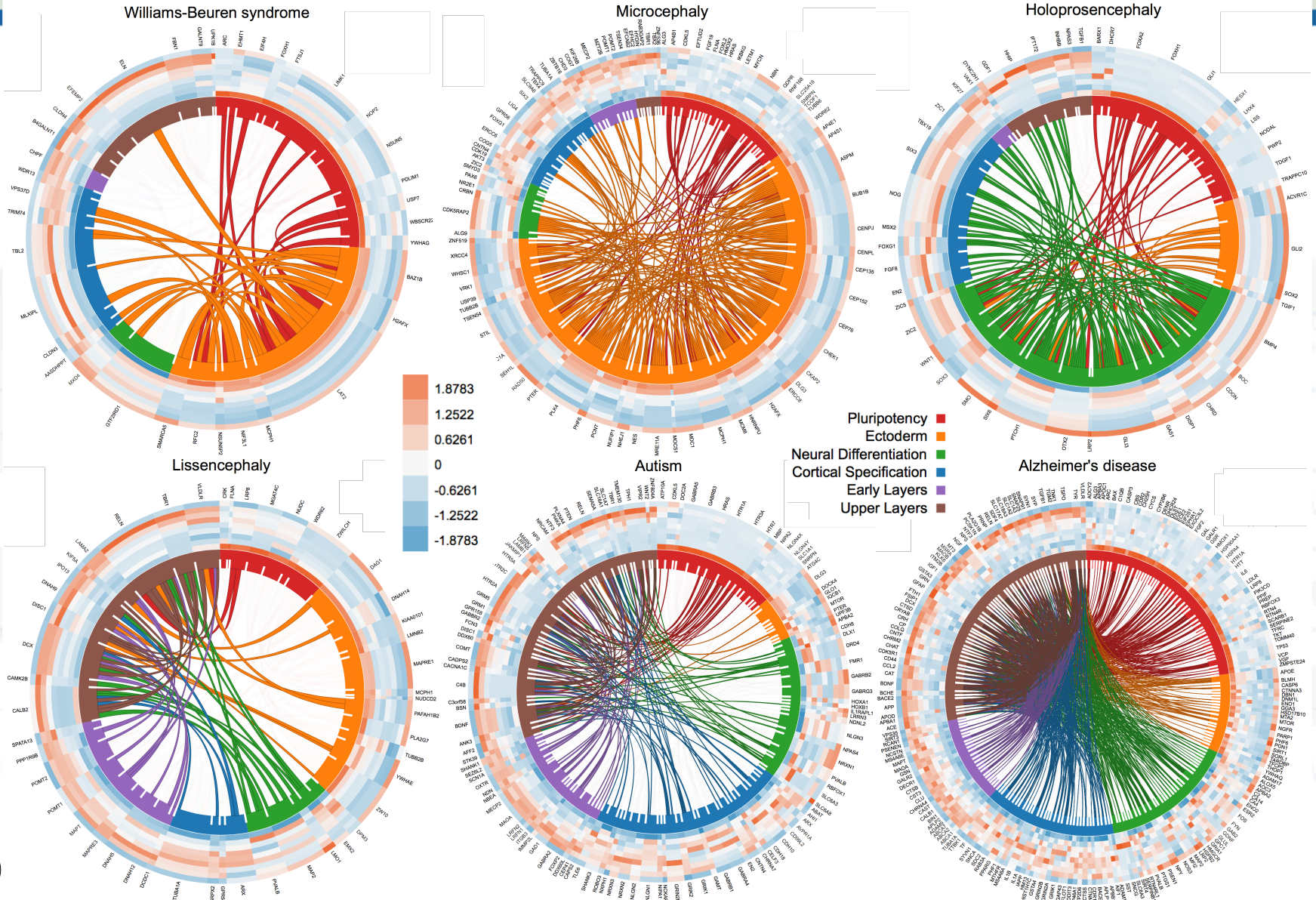


Visualization





Visualization





Results

- Through the application of fuzzy c-means clustering, we discovered a new cluster of genes (Ectoderm) in the early stages of brain development
- Our combination of statistical analysis, linking to structured data sources, and visualization allow us to provide domain scientists with deeper insights into relationships within and between clusters



Conclusions

- We modeled statistical analyses of genomic data using best-in-class ontologies
- Analysis results were linked with additional structured data to provide literature support
- SemNExT combines statistical analysis and linked data in a generalizable way by incorporating new analyses, data sources, and ontologies
- We are using SemNExT with structured knowledge to make sense of a newly identified cluster of genes in neural development



Future Directions

- Applying SemNExT to two related domains of plant microbiome and child development
- Looking for users/collaborators for feedback
 - Contact d1m@cs.rpi.edu or pattoe@rpi.edu



Acknowledgements

- NSF Graduate Research Fellowship:
Mr. Patton
- NSF Grant 1331023: Ms. Brown, Ms.
De los Santos, and Dr. Bennett
- RPI Internal Funding
- Dr. John Erickson for valuable feedback
on the project
- SemStats 2015 workshop chairs and
reviewers



References

- Neural Stem Cell Institute – <http://www.neuralsci.org/>
- NSCI Cortecon – <http://cortecon.neuralsci.org/>
- KEGG – <http://www.genome.jp/kegg/>
- StringDB – <http://string-db.org/>
- Ensembl – <http://www.ensembl.org/>
- UMLS – <https://www.nlm.nih.gov/research/umls/>
- Bio2RDF – <http://bio2rdf.org/>
- ReDrugS – <http://redrugs.tw.rpi.edu/>
- Uniprot – <http://www.uniprot.org/>
- SemNExT Demo – <https://semnext.tw.rpi.edu/chem/>



TWOC

QUESTIONS?

d1m@cs.rpi.edu | pattoe@rpi.edu

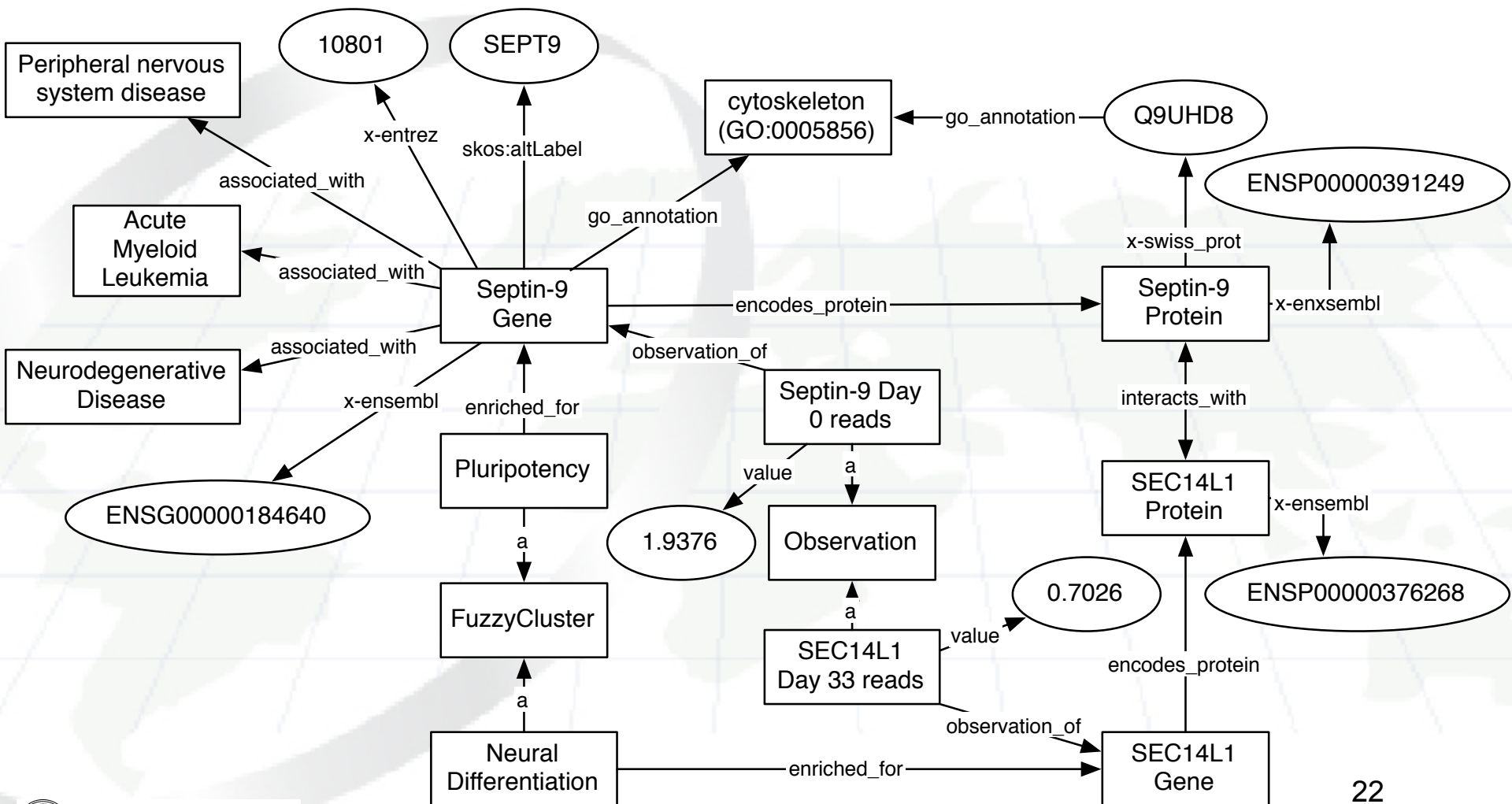


TWIC

Extras

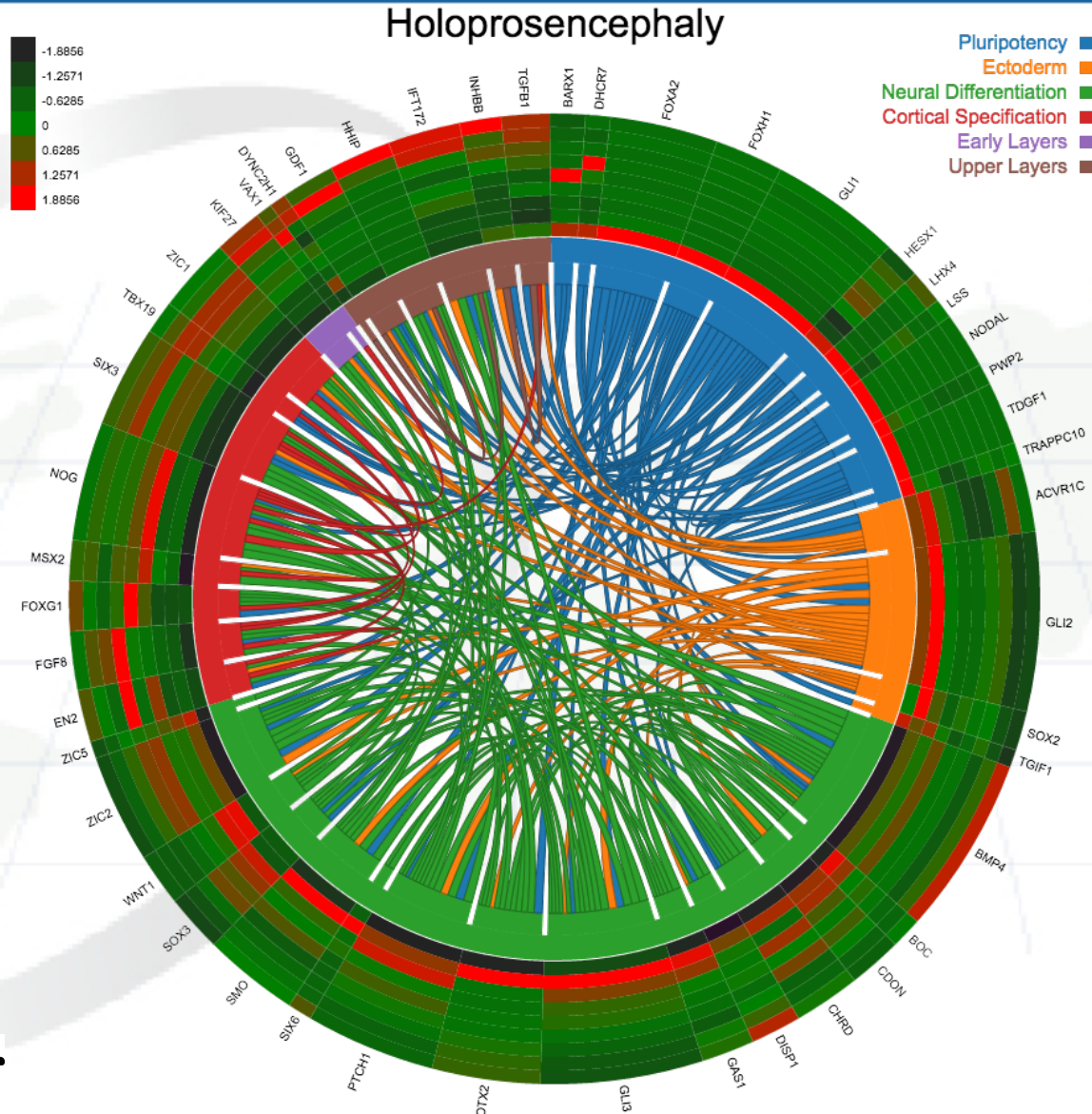


Example



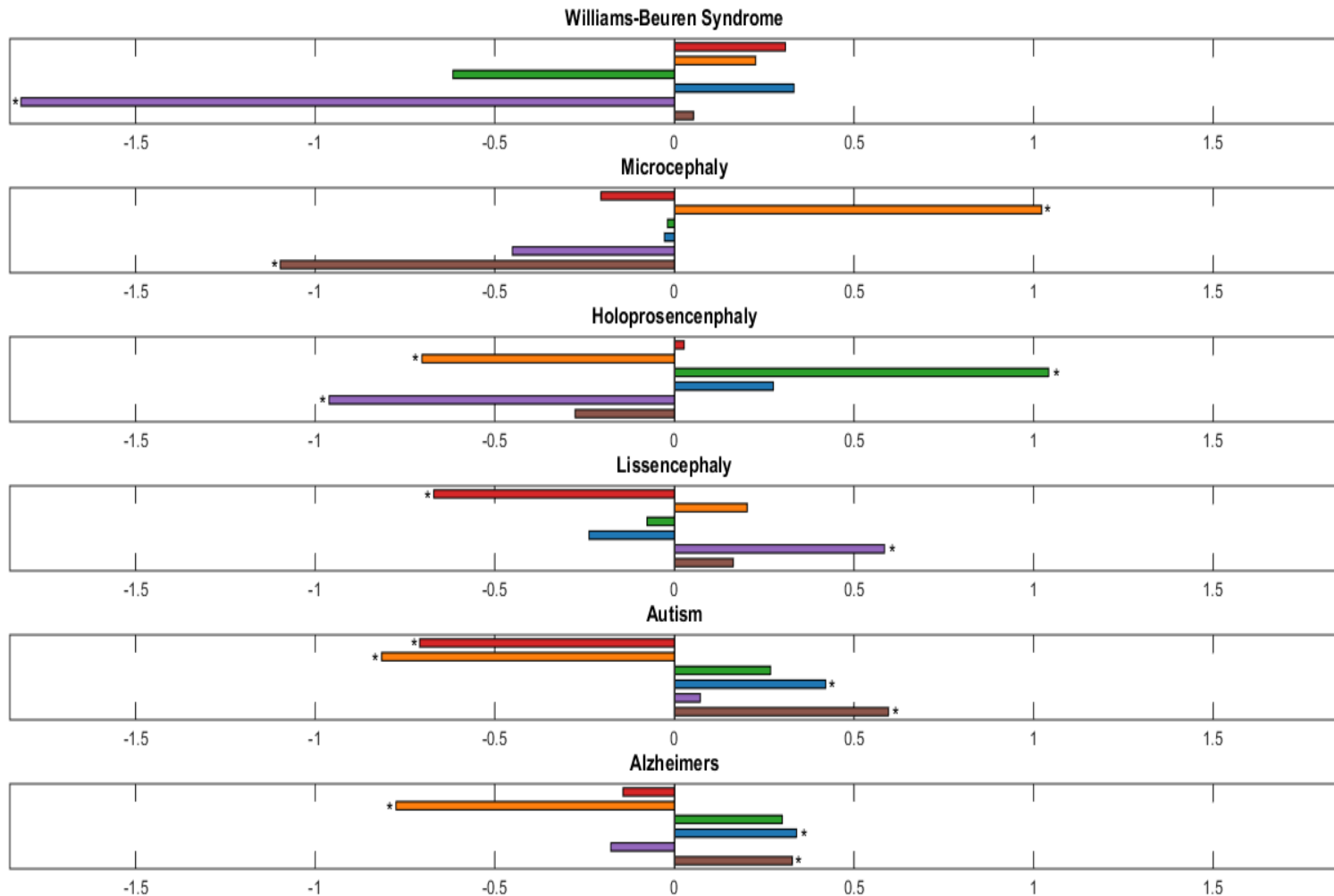
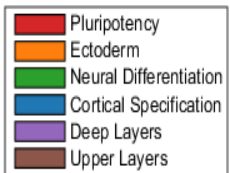


Visualization





Stage Enrichment for Selected Diseases



Less genes than expected

More genes than expected