

QQuerying with Ontological Terminologies and their Annotations (QUOTA)

Yi Sun

Ontology Research Group
Miami University

OCWIC'2007
February 17,
Deer Creek Resort



Outline

- Introduction & Motivation
 - Ontology Overview
 - Basic Concepts
 - Why Query Annotations?
- Related Research
 - GeneInfoViz
 - GO Categorizer
 - eVOC
- QUOTA



Ontology Overview

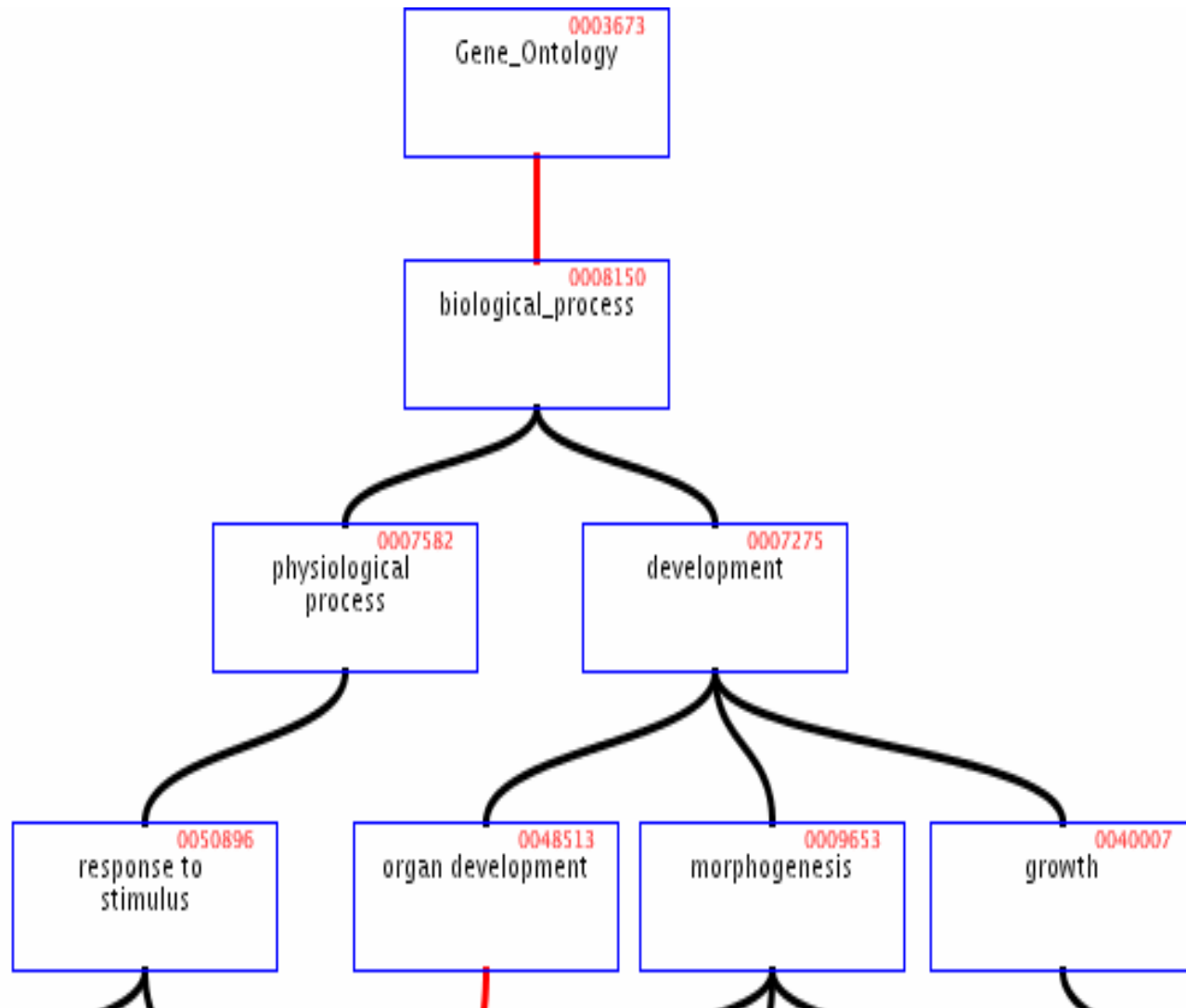
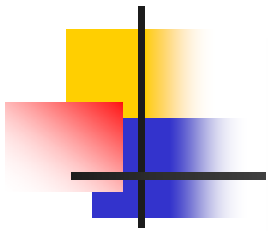
- What is Ontology?
 - Conceptualization
- What's in Ontology?
 - Concepts
 - Definitions of concepts
 - Relationships among concepts
- Why Ontology?



Gene Ontology (GO)

- Developed by the Gene Ontology Consortium (www.geneontology.org)
- GO Components
 - Cell Component (CC)
 - Molecular Function (MF)
 - Biological Process (BP)

GO Example

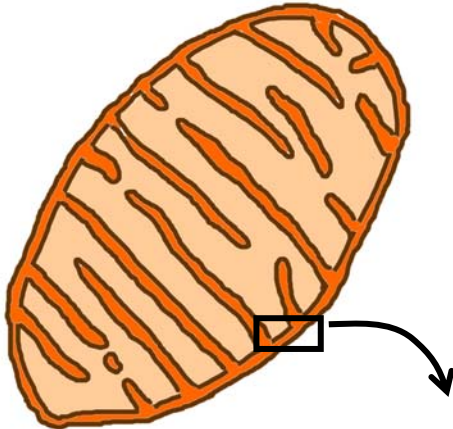




Annotation

- What is Annotation?
 - Association between product items and Ontology terms
- Multiple Annotations
 - e.g. in GO:
 - Gene → located in one or more cell components
 - Gene → perform one or more molecular functions
 - Gene → be active in one or more biological processes

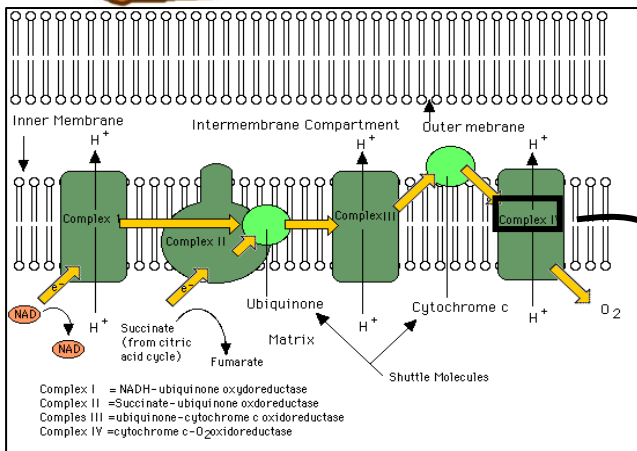
Example: Mitochondrial P450



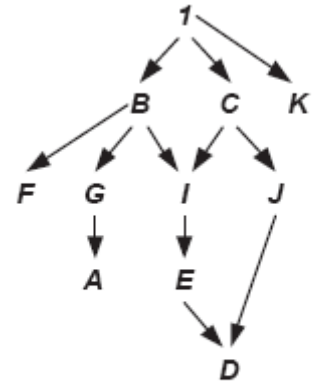
Cellular component:
mitochondrial inner membrane
GO:0005743

Biological process:
Electron transport
GO:0006118

Molecular function:
monooxygenase activity
GO:0004497



Basic Concepts used in Research



- Interval rank

- Improvement on Level

- $\text{Interval_rank}(n) = (\text{distRoot}(n) + \text{distBottom}(n)) / 2$

- Information Content

- Metric proposed by Seco, Veale and Hayes (2004)

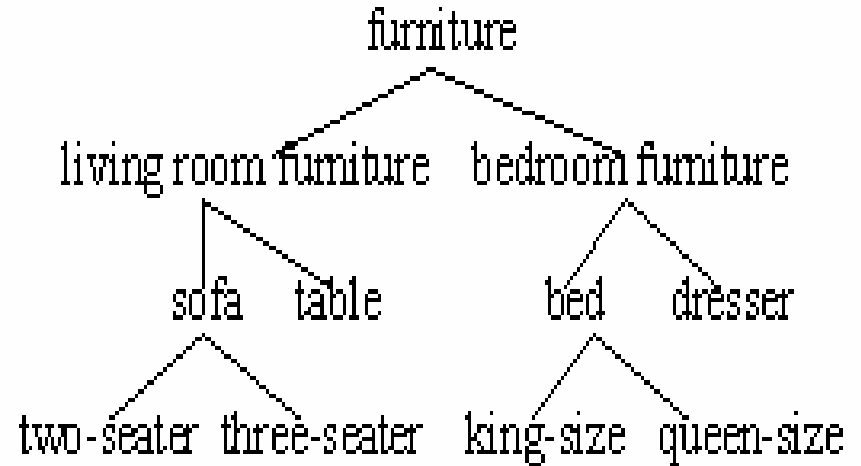
- $\text{IC}(n) = \log [(\text{hypo}(n) + 1) / \text{max}] / \log (1 / \text{max})$

hypo(n) is the number of hyponyms (descendents)

max is the maximum number of concepts in the taxonomy

- Semantic similarity

Semantic Similarity



- Use a value (similarity) to determine closeness of two objects (concepts)
- Distance Based:

$$\text{sim}(c1, c2) = \frac{2 * \text{len}(\text{root}, c3)}{(\text{len}(c1, c3) + \text{len}(c2, c3) + 2 * \text{len}(\text{root}, c3))}$$

where c3 is the lowest common subsumer (LCS)

- Information Content Based:

$$\text{sim}(c1, c2) = \frac{2 * \log \text{IC}(c3)}{[\log \text{IC}(c1) + \log \text{IC}(c2)]}$$



Why Query Annotations?

- Retrieve information and make prediction based on current information
- For example, in GO:
Given a set of genes, what are the terms that best represent these genes, what is the relationship among these genes?



Related Research

- GeneInfoViz
- GO Categorizer
 - eVOC



GeneInfoViz

- Implemented by Mi Zhou and Yan Cui, 2004
- Basic Functionality: Given a set of genes
 - 1: Display the GO terms that are annotated to these genes and highlight them in GO DAG
 - 2: Create relationship network between genes based on overlapping terms

GeneInfoViz: example

GeneInfoViz: Constructing and Visualizing Gene Relation Networks

Select Organism

Human

Zebra Fish

Mouse

Cow

Rat

C. Elegan

Fruit Fly

[Q&A](#)

Enter your list of genes here

[Click for Sample Gene List \(LocusIDs\)](#)

```
4605 3002 6659 4609
1942 6662 1956 6722
2064 4485 5327 1382
4602 2625 7494
```

Search By

Select Evidence Code

- TAS (Traceable author statement)
- IDA (Inferred from direct assay)
- IMP (Inferred from mutant phenotype)
- IGI (Inferred from genetic interaction)
- IPI (Inferred from physical interaction)
- E (Experimental evidence)
- ISS (Inferred from sequence or structural similarity)
- IEP (Inferred from expression pattern)
- NAS (Non-traceable author statement)
- IEA (Inferred from electronic annotation)
- P (Predicted/computed)
- NR (Not recorded)

Search

GeneInfoViz: Constructing and Visualizing Gene Relation Networks

Dynamic Visualization of Gene Relation Network Based On

Biological Process

Level 1

Directed Acyclic Graph

Molecular Function

Level 1

Directed Acyclic Graph

Cellular Component

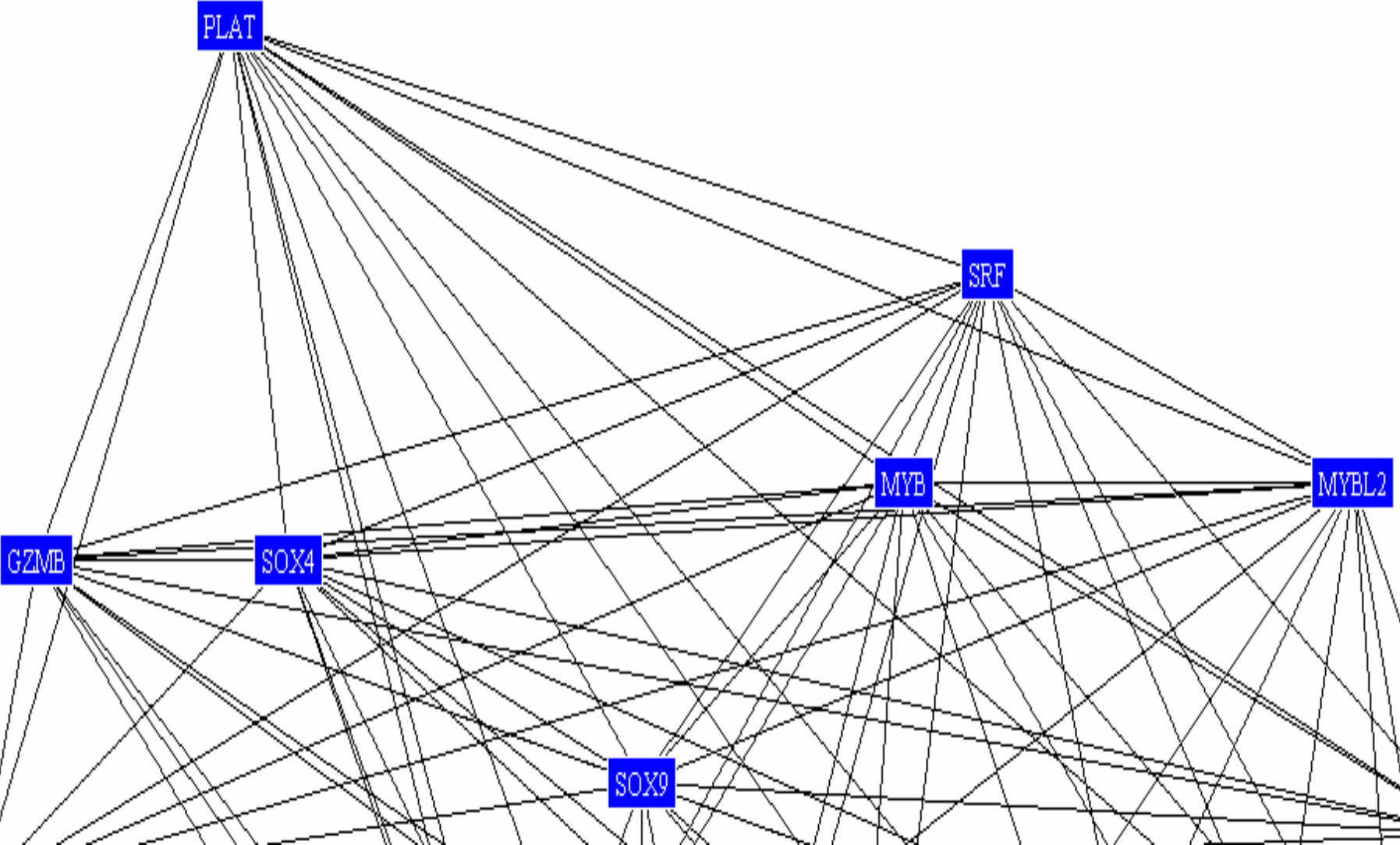
Level 1

Directed Acyclic Graph

Note: Biological Process[P], Molecular Function[F], Cellular Component[C];

UniGeneID	LocusID	Symbol	Gene Name	Gene Ontology Term	Evid
Hs.405662	1382	CRABP2	cellular retinoic acid binding protein 2	lipid binding[F]	IEA
Hs.405662	1382	CRABP2	cellular retinoic acid binding protein 2	retinal binding[F]	IEA
Hs.405662	1382	CRABP2	cellular retinoic acid binding protein 2	transporter activity[F]	IEA
Hs.405662	1382	CRABP2	cellular retinoic acid binding protein 2	epidermis development[P]	TAS
Hs.405662	1382	CRABP2	cellular retinoic acid binding protein 2	regulation of transcription, DNA-dependent[P]	TAS
Hs.405662	1382	CRABP2	cellular retinoic acid binding protein 2	signal transduction[P]	TAS
Hs.405662	1382	CRABP2	cellular retinoic acid binding protein 2	transport[P]	IEA
Hs.516664	1942	EFNA1	ephrin-A1	ephrin receptor binding[F]	TAS
Hs.516664	1942	EFNA1	ephrin-A1	cell-cell signaling[P]	TAS
Hs.516664	1942	EFNA1	ephrin-A1	integral to plasma membrane[C]	TAS
Hs.516664	1942	EFNA1	ephrin-A1	membrane[C]	IEA
Hs.488293	1956	EGFR	epidermal growth factor receptor (erythroblastic 1	ATP binding[F]	IEA
Hs.488293	1956	EGFR	epidermal growth factor receptor (erythroblastic 1	epidermal growth factor receptor activity[F]	NR
Hs.488293	1956	EGFR	epidermal growth factor receptor (erythroblastic 1	receptor activity[F]	IEA

Gene relationship network based on overlap of annotating terms





GO Categorizer

- Implemented by Cliff Joslyn and Susan Mniszewski et al in 2004
- Basic Functionality: Given a set of genes
 - Generate a ranked list of GO terms with values indicating how well they represent these genes
- Significance:
 - With this ranked list, we can easily tell if the terms best representing these genes, for example, cluster at one place or are evenly



eVOC

- Implemented by Janet Kelso et al in 2003
- Basic Functionality: Given a set of terms, return sets of genes that are annotated by each of the terms, and perform set operations based on the logical query.
- For example: Find the genes that are related to both “liver” and “cancer”
 - $\text{Set}(\text{“liver”}) \cap \text{Set}(\text{“cancer”})$



QUOTA

- Motivation & Objective
 - Querying Direction:
- In: Input to QUOTA
- Out: Output from QUOTA

1. IN: Set of Annotated Objects

OUT: Ontological Terms

- Similar to GeneInfoViz Function 1
- Generate the Annotation for each Object
- Provide various set options (AND, NOT, OR)

2. IN: Set of Annotated Objects

OUT: Ranked Ontological Terms

- Originated from GO Categorizer
- Generate ranking score for each term
- Provide more options:
 - Interval rank
 - Information Content
 - Various Semantic Similarity measures

3. IN: Set of Annotated Objects

OUT: Similarity value matrix

- Generate similarity matrix among objects
- Based on Semantic Similarity values between each pair of annotating terms
- Various aggregation methods to integrate all SS values
- For example:
 - $AO1 \rightarrow t1, t3, t5; AO2 \rightarrow t2, t4, t6;$

4. IN: Ontological Terms

OUT: Annotated Objects

- Similar to eVOC
- Given a term, generate a set of objects annotated by this term
- Perform set operations (AND, NOT, OR)



5. IN: Two sets of Annotated Objects

OUT: Semantic Similarity Value

- Input: Two sets of annotated objects
 - e.g. Two diseases, each one is potentially marked by a set of genes through some laboratory experiments
- Output: Semantic similarity between the two sets
 - e.g. A value indicating how similar two diseases are

5. IN: Two sets of Annotated Objects

OUT: Semantic Similarity Value (cont...)

- Based on various Semantic Similarity measures
- Possible implementation:
 - For $A = \{g_1, g_2, g_3\}$ and $B = \{g_5, g_6\}$
 - Calculate $\text{sim}(g_i, g_j)$ for all pairs where g_i in A and g_j in B using various methods in Query 3
 - Use some aggregation function to determine an overall similarity between A and B



Generic Tool

- Generic
 - Works for all Ontologies and Annotation files with valid formats
- A translation program to convert Ontology into our formats has already been developed



Resources & Tools

- Gene Ontology & Annotation:
 - <https://www.geneontology.org>
 - Formats:
 - GO: txt, OWL, XML
 - Annotation: XML, MySQL
- GeneInfoViz: <https://www.genenet.org>
- GO Categorizer: <http://www.c3.lanl.gov/posoc/>
- eVOC: info@egenetic.com

Thank you!

