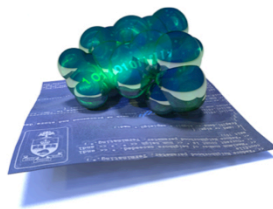A
BIOINFORMATICS
COURSE

# G O
## (GENE ONTOLOGY)

BORIS STEIPE

DEPARTMENT OF BIOCHEMISTRY − DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO

Experimentally *observable properties* of a protein like sequence or structure are straightforward to abstract, store, retrieve and interpret.

Aspects of a *protein's behaviour* like conservation, localization, interactions, regulation of expression etc. require more context – but are also observable and don't pose problems of a principal nature.

However "Function" is a *concept* that aims to integrate a a large variety of observables and not-observables of a gene, and its role in its molecular and cellular context. There is no "natural" set of categories and values that present itself to reason about function.
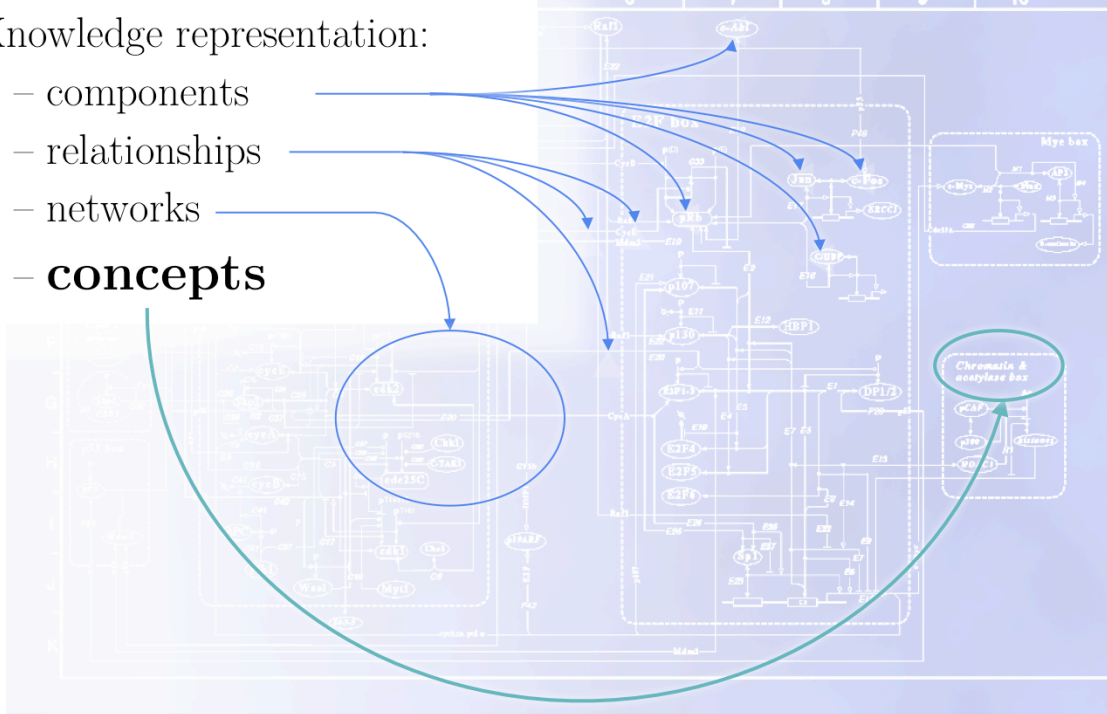
In this sense "Function" is not directly observable.

In general, this is not sufficiently appreciated.

That aspect of biology that realizes the ultimate goal of our endeavours is the most elusive: we can't even properly talk about function. The first task of function analysis therefore is knowledge engineering: to define a "language" in which we may categorize, collect and compare functions.

Knowledge representation:
 – components
 – relationships
 – networks
 – **concepts**

If we believe everything of interest about biology can be expressed as entities*, their attributes, and the relationships between them, we are missing the importance of "concepts" that represent our ideas about *why* systems are composed the way they are.

Considering biomolecular sequences, and the structure of activities of biomolecules, it is obvious that it is possible to order and arrange those to pathways and networks of activities (albeit that is not trivial in practice).

However there is no automatic way to structure and categorize sets of such relationships into a higher level "understanding", into insight *why* the observed entities and relationships give rise to their behaviour, and, *what the "purpose" of the biomolecular systems is*, that {purpose / meaning / objective / fitness function} that has been selected by evolution.

* an "entity" can be a gene, protein or metabolite, a compartment, an activity …

3

## Gene Ontology

Community organized consortium.
Maintains a controlled vocabulary of terms (attributes).
Establishes relationships between attributes in a DAG (Directed Acyclic Graph).
Annotates genes.
Provides data and tools.

Component ontologies:
*Cellular Component,*
*Molecular Function,*
*Biological Process.*

Relationships:
is_a,
has_a,
part_of...

GO is the gold standard for function annotation and used by practically all major molecular databases.

http://www.geneontology.org

In order to speak about function, we need a common language. A common language means that we use the same terms to describe the same facts. To establish a common language about function for molecular biology is the goal of the Gene Ontology (GO) consortium.

The term "ontology" comes from the domain of *knowledge engineering.* An ontology collects terms that describe facts, their definitions, and their relationships.

GO has grown over many years and it is a well funded, large and mature project.

The graph of GO definitions contains three distinct component ontologies. Their root nodes are Cellular Componen, Molecular Function, and Biological Process.

# Gene Ontology Annotations

Browse annotations
with QuickGO ...

http://www.ebi.ac.uk/QuickGO/

... or download ...



Gene Ontology Annotation (GOA) is part of GO's biocuration effort, aiming to associate proteins in UniProtKB with their repsective GO terms. After searching for genes of interest in the QuickGO browser, you can select GO terms that you would like to explore, find all proteins annotated to these nodes, and filter for particular taxons.

What type of evidence is the annotation based on?

Experimental (wet lab)

Non-Experimental (computational)

Author statement from publication

Curator Statements

Is annotation based on genetic mutations or allelic variation?
no / yes

Is a single gene being mutated or compared to other alleles of the same gene?
no / yes → **IMP**

Is annotation based on a genetic interaction with another gene?
yes → **IGI**

Is annotation based on a direct 1 to 1 physical interaction with another gene product?
no / yes → **IPI**

Is annotation based on a direct assay for the function, process, or component of the gene product?
no / yes → **IDA**

Is annotation based on the expression pattern of the gene product?
yes → **IEP**

Will each annotation be individually reviewed & confirmed by a human annotator?
yes / no → **IEA**

Is the computation based purely on the sequence of the gene product?
no / yes → **ISS ISA ISM ISO**

Does the computation include consideration of the genomic context of the gene?
no / yes → **IGC**

Is the computation an integrated analysis, typically including experimental data sets, and often including multiple data types?
yes → **RCA**

Is annotation based on an author statement that cites a published reference as the source of the information?
no / yes → **TAS**

Is annotation based on an author statement that does not cite a published reference as the source of the information?
yes → **NAS**

Is there a GO annotation in another aspect that allows a biocurator to make an inference based on that GO term for an aspect without evidence?
no / yes → **IC**

Are there No Data (ND) to support a GO annotation in a given GO aspect? (see note on use of ND)
yes → **ND**

■ For curator reviewed annotations
□ For annotations NOT reviewed by a curator

**Note on use of ND evidence code:**
Unlike the other evidence codes, the No Data (ND) code does not indicate evidence or a method from a specific reference. Rather, it indicates that the annotator looked at the available information and determined that nothing is known about the gene for a given aspect of GO (molecular function, biological process, or cellular component). The annotater will always look at all available literature for the gene. Depending on the resources and annotation philosophy of the annotating group, the annotater may also look at sequence comparison data to determine if any predictions may be made based on the sequence.

GO as based on definitions, but annotations in GOA reqiure evidence. GOA evidence codes make the evidence explicit for every annotation. They are a crucial part of GOA that you **must** be familier with to work with the data and subset it according to your application's needs.

**Experimental Evidence Codes**
**DA**: Direct Assay
**IPI**: Physical Interaction
**IMP**: Mutant Phenotype
**IGI**: Genetic Interaction
**IEP**: Expression Pattern

**Computational Analysis Evidence Codes**
**ISS**: Sequence or Structural Similarity
**ISO**: Sequence Orthology
**ISA**: Sequence Alignment
**ISM**: Sequence Model
**IGC**: Genomic Context
**RCA**: Reviewed Computational Analysis

**Author Statement Evidence Codes**
**TAS**: Traceable Author Statement
**NAS**: Non-traceable Author Statement

**Curator Statement Evidence Codes**
**IC**: Inferred by Curator
**ND**: No biological Data available

http://www.geneontology.org/page/evidence-code-decision-tree

GO Component Ontologies – GOA Example

Simple search by gene name.
Example:
Res2 – (The *S. pombe* Mbp1 orthologue)

Cellular Component Ontology

GO is divided into three complementary views of function: *Cellular Components* ...

"These terms describe a component of a cell that is part of a larger object, such as an anatomical structure (e.g. rough endoplasmic reticulum or nucleus) or a gene product group (e.g. ribosome, proteasome or a protein dimer)."

http://www.geneontology.org/page/ontology-documentation

GO Component Ontologies – GOA Example

Simple search by gene name.
Example:
Res2 – (The *S. pombe* Mbp1 orthologue)

Molecular Function Ontology

... *Molecular Function* ...

"Molecular function terms describes activities that occur at the molecular level, such as "catalytic activity" or "binding activity". GO molecular function terms represent activities rather than the entities (molecules or complexes) that perform the actions, and do not specify where, when, or in what context the action takes place. Molecular functions generally correspond to activities that can be performed by individual gene products, but some activities are performed by assembled complexes of gene products. Examples of broad functional terms are "catalytic activity" and "transporter activity"; examples of narrower functional terms are "adenylate cyclase activity" or "Toll receptor binding".
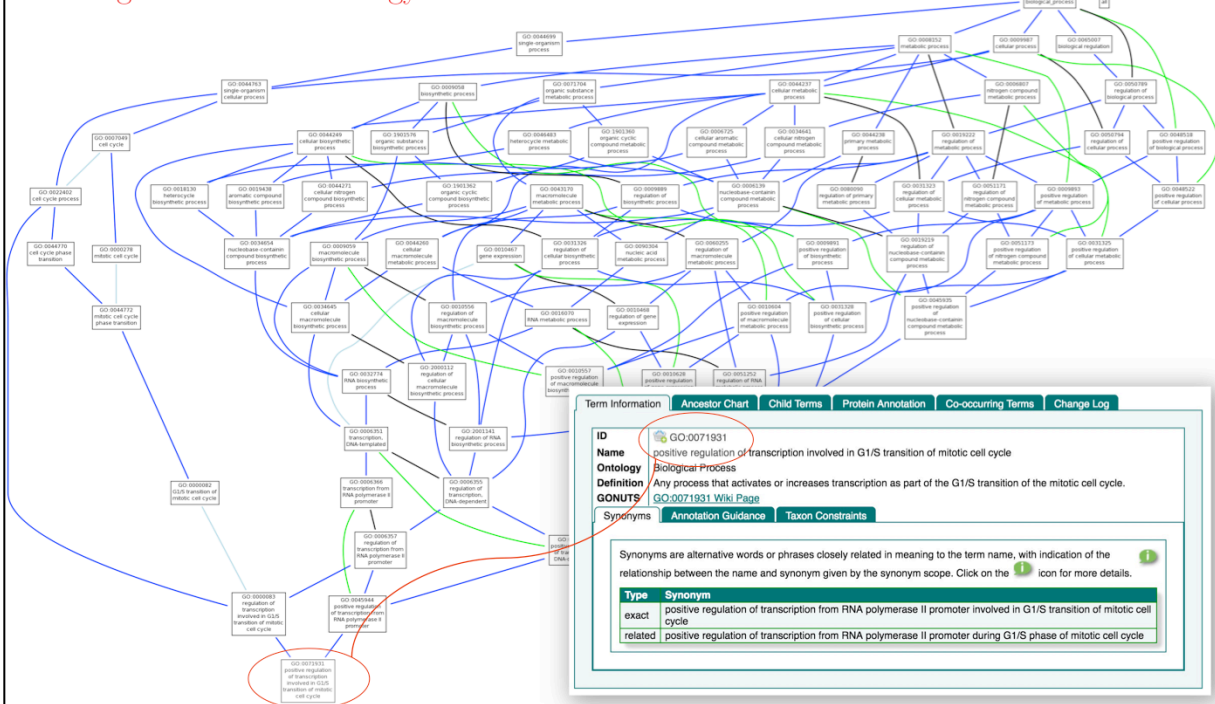
It is easy to confuse a gene product name with its molecular function; for that reason GO molecular functions are often appended with the word "activity"."

GO Component Ontologies – GOA Example

Simple search by gene name. Example:
Res2 – (The *S. pombe* Mbp1 orthologue)

Biological Process Ontology

... and *Biological Process*.

"A biological process term describes a series of events accomplished by one or more organized assemblies of molecular functions. Examples of broad biological process terms are "cellular physiological process" or "signal transduction". Examples of more specific terms are "pyrimidine metabolic process" or "alpha-glucoside transport". The general rule to assist in distinguishing between a biological process and a molecular function is that a process must have more than one distinct steps.
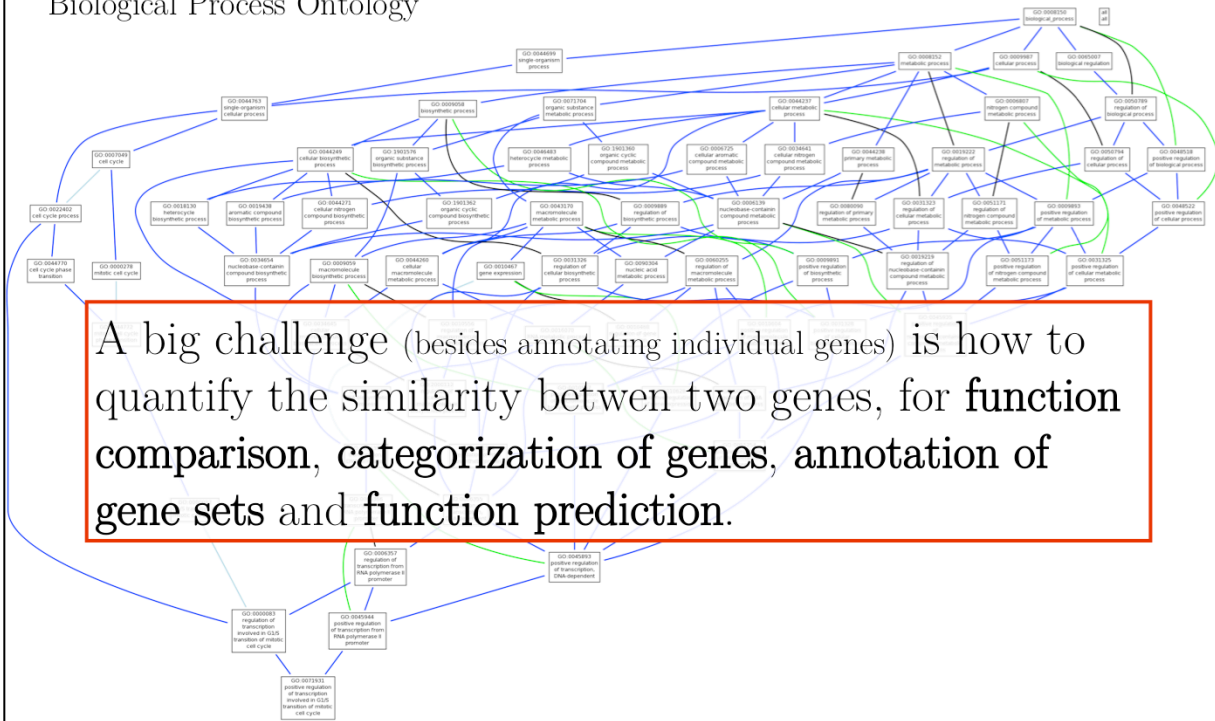
A biological process is not equivalent to a pathway. At present, the GO does not try to represent the dynamics or dependencies that would be required to fully describe a pathway."

Simple search by gene name.mExample:
Res2 – (The *S. pombe* Mbp1 orthologue)

Biological Process Ontology

A big challenge (besides annotating individual genes) is how to quantify the similarity betwen two genes, for **function comparison**, **categorization of genes**, annotation of **gene sets** and **function prediction**.

http://steipe.biochemistry.utoronto.ca/abc

B O R I S . S T E I P E @ U T O R O N T O . C A

DEPARTMENT OF BIOCHEMISTRY & DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO, CANADA