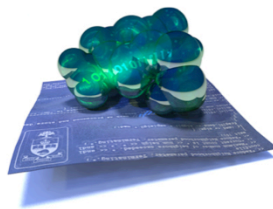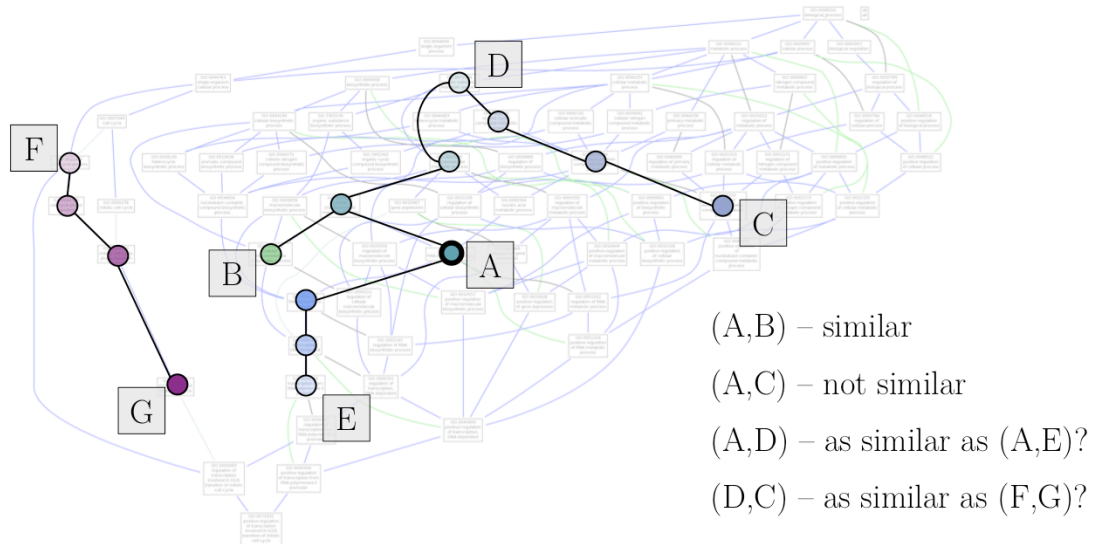# Semantic Similarity

Boris Steipe

DEPARTMENT OF BIOCHEMISTRY – DEPARTMENT OF MOLECULAR GENETICS
University of Toronto

## Calculating semantic similarity:

Once genes have been annotated to a GO term (in GOA),
two genes are considered *similar* if the nodes they are annotated to are close in the
Gene Ontology.



(A,B) – similar

(A,C) – not similar

(A,D) – as similar as (A,E)?

(D,C) – as similar as (F,G)?

It is easy to see that nodes A and B are similar since they share the same parent,
whereas A and C are not – their lowest common ancestor is three steps away.
However matters are not so clear–cut if we are asking about relationships at different
levels of the tree (as between (A,D) and (A,E)), or in parts of the tree that are broad
(as with (D,C)) as compared to parts that are narrow as with (F,G).

A unifying concept is to ask how *specific* each parent is, i.e. which portion of the tree
it has beneath it. This can be expressed as *information* (in the information-
theoretical sense).

Nb.: The GO is not actually a tree but a DAG (Directed Acyclic Graph). This means that  nodes
can have more than one parent but the connections are directed such that there are no cycles in
the graph. A tree *is* a DAG, but not all DAGs are trees. This is necessary, because many
concepts can be valid children of more than one parent concept, but this arrangement also
creates difficulties since there may be multiple pathwaysto higher nodes and to the root.

Calculating semantic similarity:

## Information based measures

Example:
16 nodes
40 annotated genes



$$p_i = \frac{\sum_{C_i} |\text{annotated genes}|}{\sum_{O} |\text{annotated genes}|}$$

$C_i \in \{\text{node of } i \text{ and all child nodes}\}$
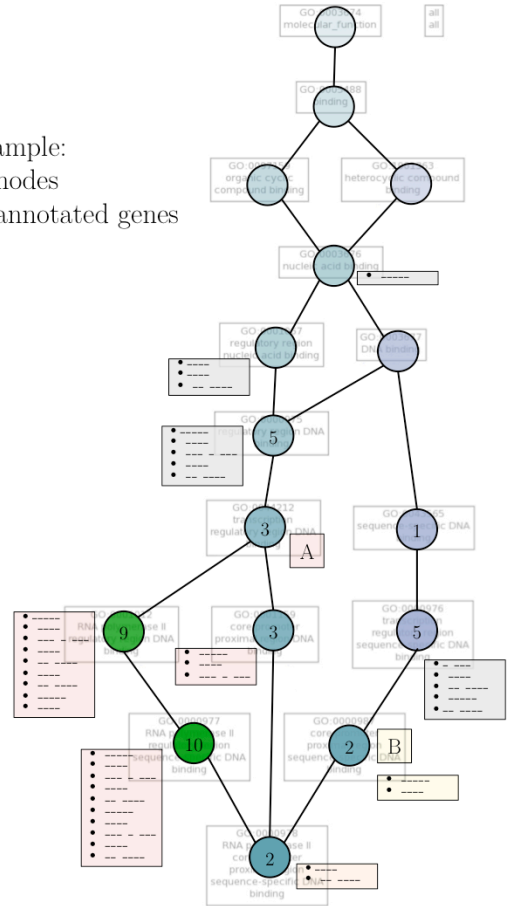$O \in \{\text{all nodes in ontology}\}$

$$p_A = \frac{3 + 9 + 10 + 2}{40} = 0.625$$

$$p_B = \frac{2 + 2}{40} = 0.1$$

$$I = -\log(p)$$

$$I_A = -\log(0.625) = 0.47$$
$$I_B = -\log(0.1) = 2.3$$

The "probability" $p_A$ of a gene A is the number of genes annotated to it and its children, divided by the number of genes in the tree.

"Information" is the negative log of the probability.

3

Calculating semantic similarity:

## Information based measures

Example:
16 nodes
40 annotated genes



One way to define *semantic similarity* is simply the information of the "Most Informative Common Ancestor" (MICA) (Resnik 1999).

$$\text{SIM}(A, B) = I_{\text{MICA}(A,B)} = -\log\left(\frac{36}{40}\right) = 0.11$$

Other approaches take the individual nodes' information into account (e.g. Lin 1998):

$$\text{SIM}(A, B) = \frac{I_{\text{MICA}(A,B)}}{I_A + I_B} = 0.04$$

Resnik, P (1995). Semantic similarity in a taxonomy: An Information-Based measure and its application to problems of ambiguity in natural language. Journal of Artificial Intelligence Research 11, 95–130.

Lin, D. (1998). An Information-Theoretic Definition of Similarity. In Proceedings of the Fifteenth International Conference on Machine Learning (ICML '98), Jude W. Shavlik (Ed.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 296-304

Calculating semantic similarity:

## Information based measures

Example:
16 nodes
40 annotated genes

Still other approaches take the position of the nodes in the graph into account, by down-weighting the contributions of more distant nodes. (Wang 2007).

Wang, J. Z., Du, Z., Payattakool, R., Yu, P. S. & Chen, C.-F. (2007). A new method to measure the semantic similarity of gO terms. Bioinformatics (Oxford, England) 23, 1274–1281.

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

BORIS . STEIPE@UTORONTO.CA

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

o