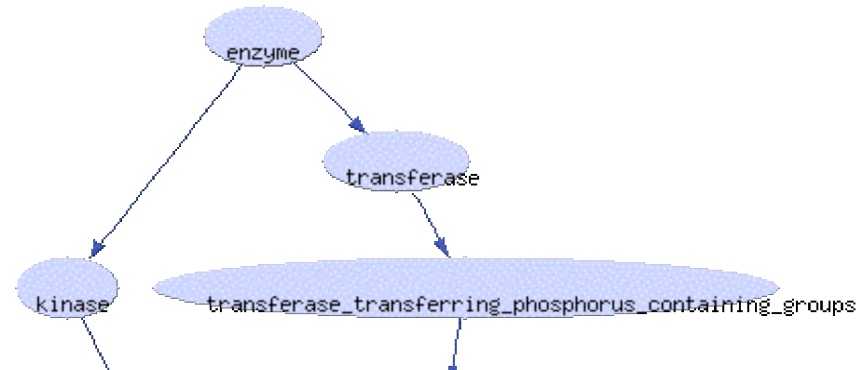# A Knowledge-Based Clustering Algorithm Driven by Gene Ontology

Jill Cheng

Affymetrix, Inc.
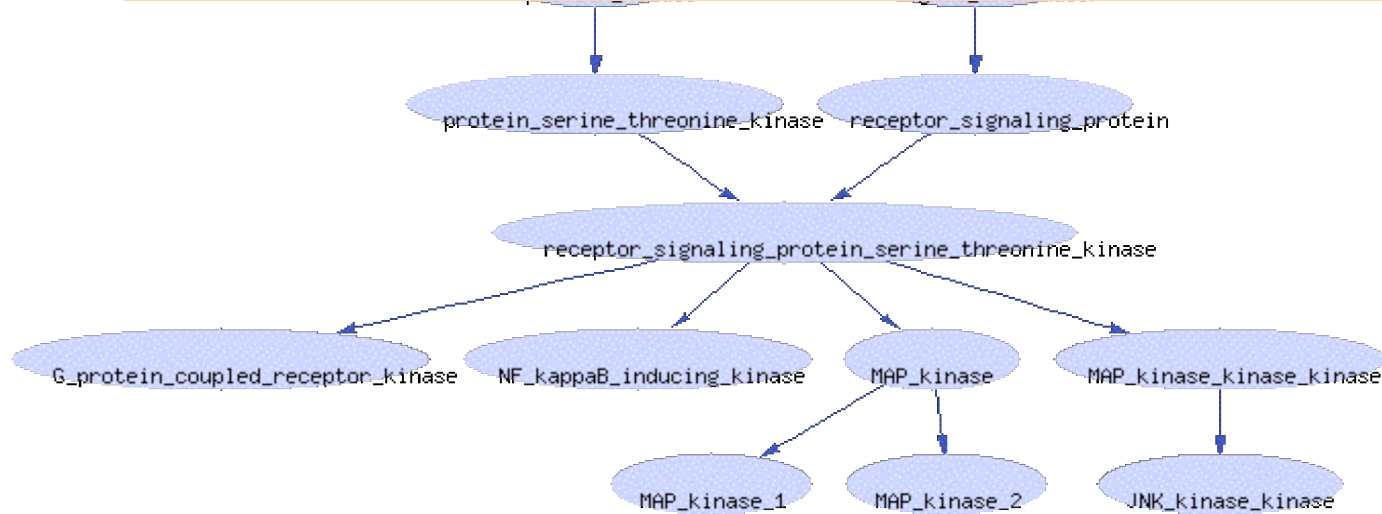
Jan 15, 2004

# The DAG structure of Gene Ontology



enzyme

transferase

kinase

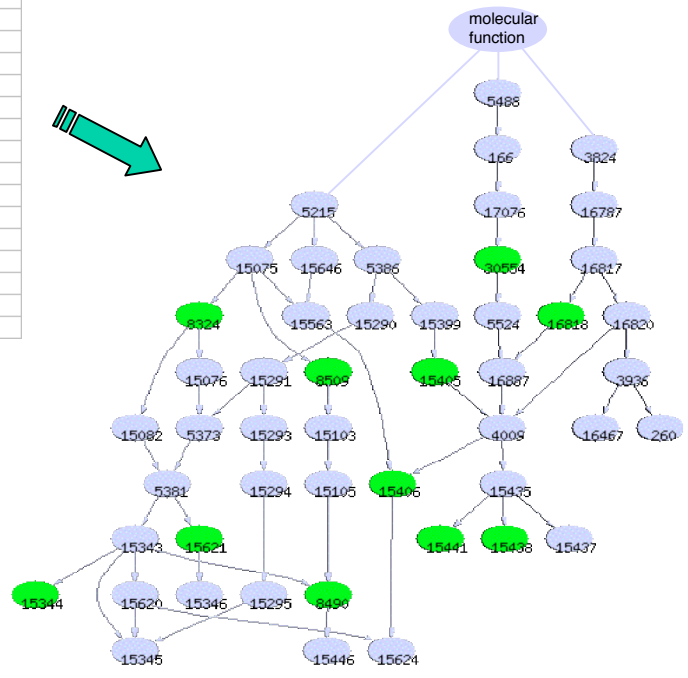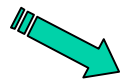transferase_transferring_phosphorus_containing_groups

One-stop-shopping for biological information

Digraphs are computable

protein_serine_threonine_kinase

receptor_signaling_protein

receptor_signaling_protein_serine_threonine_kinase

G_protein_coupled_receptor_kinase

NF_kappaB_inducing_kinase

MAP_kinase

MAP_kinase_kinase_kinase

MAP_kinase_1

MAP_kinase_2

JNK_kinase_kinase

# Goal

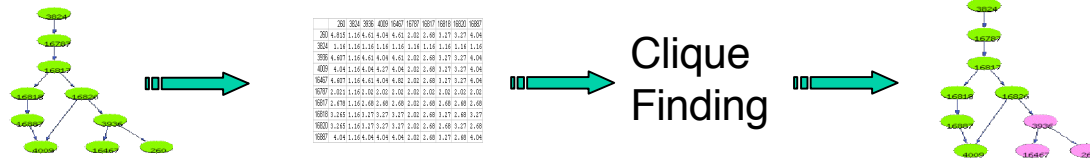The closer a node is to the root, the more general its biological classification, thus a greater amount of information is conveyed by higher level edges

The more common parent nodes shared the higher the degree of similarity

# Pair-wise similarity score between GO terms



$$W_p = \sum_{n=0}^{p} (wt)^n, \; p > 0; W_0 = 0$$

A weighting factor (*wt*) was assigned to each edge as a function of the depth (*n*) in the digraph, I chose a value of 0.815 to maximize (*wt6 – wt3*).

Determining the longest partial path shared by two nodes, *Wp* is the sum of weights for edges from root to level *p*.

$$C = \sum_{n=0}^{max-1} (wt)^n$$

A partial normalization scheme was applied to factor in the unevenness of the GO digraph.

$$Nf_p = \frac{W'_p}{W_p}$$

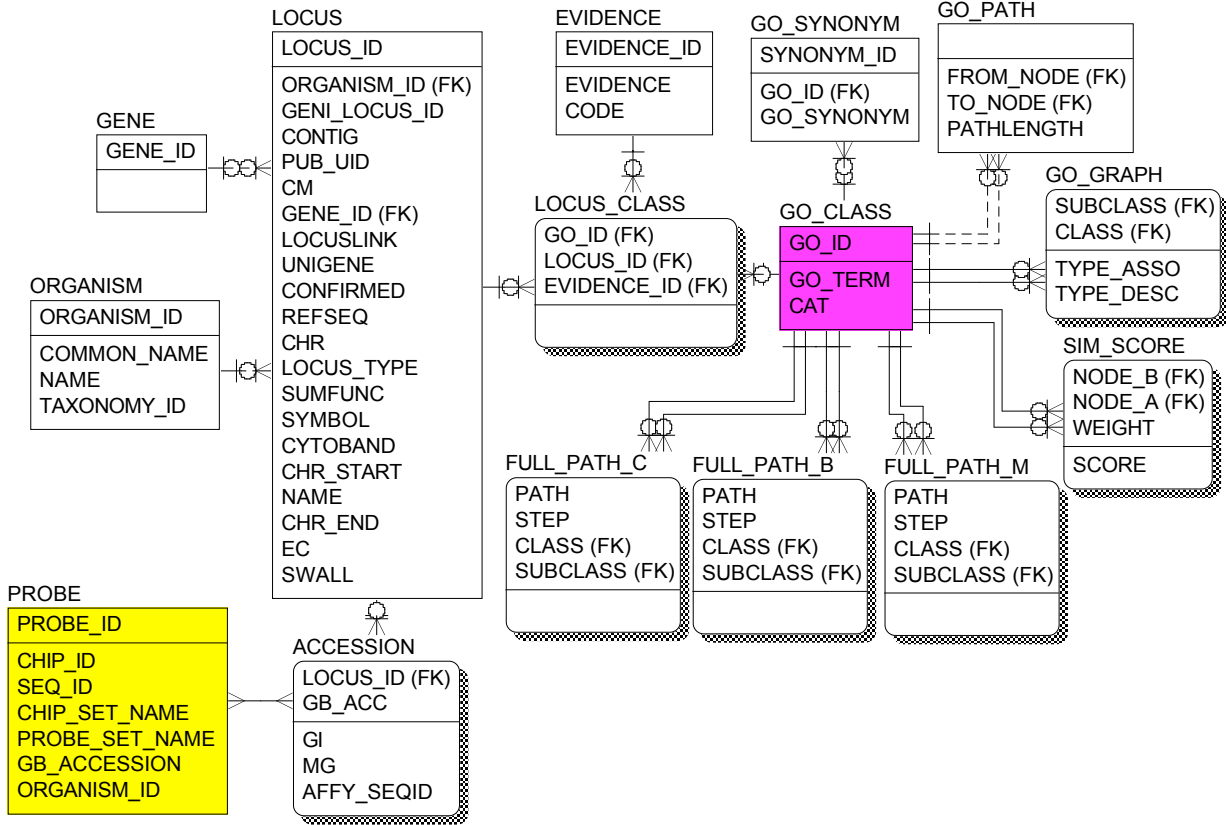Calculate the average length for all paths that go through the shared partial path (*p*), followed by the weight for a hypothetical path with *p* edges (*Wp*).

*Wp* is transformed to *W'p*, the mean of *Wp* and *C*.

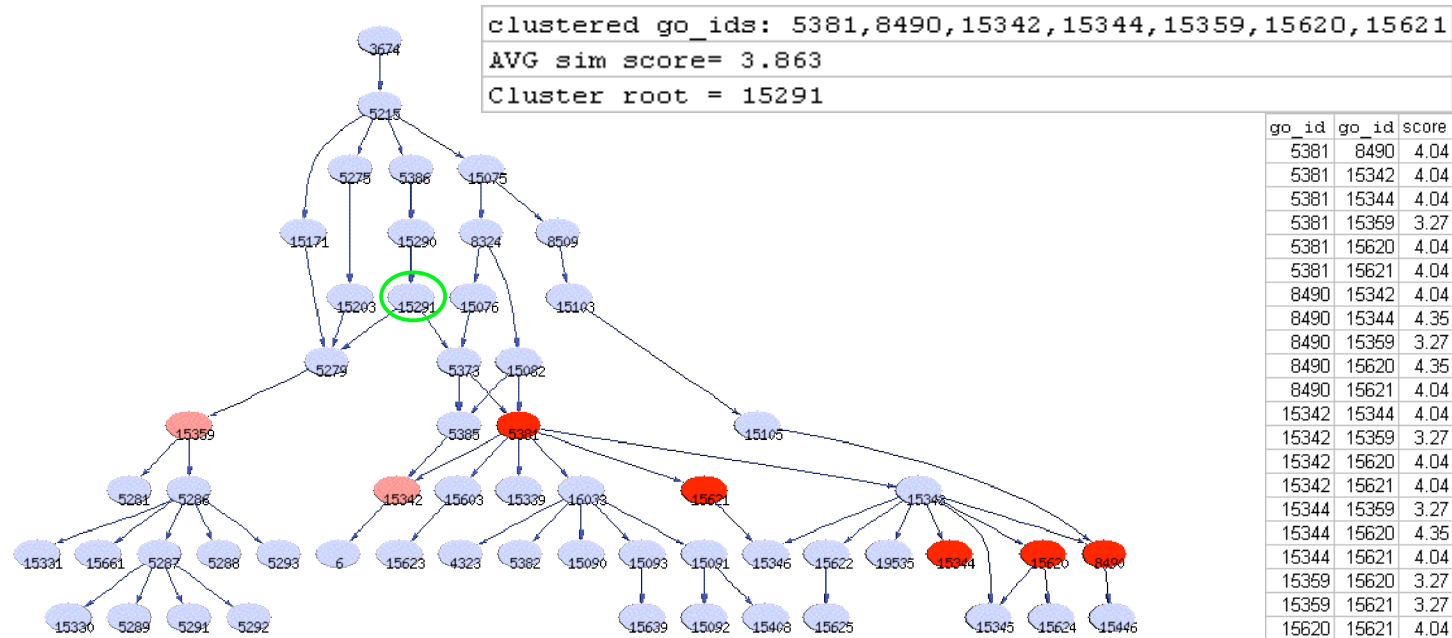$$W_m = Nf_p \sum_{n=0}^{m} (wt)^n, \; m > 0$$

The normalization factor (*Nfp*) is the ratio of *W'p* and *Wp*

The value for a partial path with m edges (*Wm*) is normalized by applying *Nfp*.

# Annotation database schema

**GENE**
| GENE_ID |
| --- |
|  |

**ORGANISM**
| ORGANISM_ID |
| --- |
| COMMON_NAME<br>NAME<br>TAXONOMY_ID |

**PROBE**
| PROBE_ID |
| --- |
| CHIP_ID<br>SEQ_ID<br>CHIP_SET_NAME<br>PROBE_SET_NAME<br>GB_ACCESSION<br>ORGANISM_ID |

**LOCUS**
| LOCUS_ID |
| --- |
| ORGANISM_ID (FK)<br>GENI_LOCUS_ID<br>CONTIG<br>PUB_UID<br>CM<br>GENE_ID (FK)<br>LOCUSLINK<br>UNIGENE<br>CONFIRMED<br>REFSEQ<br>CHR<br>LOCUS_TYPE<br>SUMFUNC<br>SYMBOL<br>CYTOBAND<br>CHR_START<br>NAME<br>CHR_END<br>EC<br>SWALL |

**ACCESSION**
| LOCUS_ID (FK)<br>GB_ACC |
| --- |
| GI<br>MG<br>AFFY_SEQID |

**EVIDENCE**
| EVIDENCE_ID |
| --- |
| EVIDENCE<br>CODE |

**LOCUS_CLASS**
| GO_ID (FK)<br>LOCUS_ID (FK)<br>EVIDENCE_ID (FK) |
| --- |
|  |

**GO_CLASS**
| GO_ID |
| --- |
| GO_TERM<br>CAT |

**GO_SYNONYM**
| SYNONYM_ID |
| --- |
| GO_ID (FK)<br>GO_SYNONYM |

**GO_PATH**
|  |
| --- |
| FROM_NODE (FK)<br>TO_NODE (FK)<br>PATHLENGTH |

**GO_GRAPH**
| SUBCLASS (FK)<br>CLASS (FK) |
| --- |
| TYPE_ASSO<br>TYPE_DESC |

**SIM_SCORE**
| NODE_B (FK)<br>NODE_A (FK)<br>WEIGHT |
| --- |
| SCORE |

**FULL_PATH_C**
| PATH<br>STEP<br>CLASS (FK)<br>SUBCLASS (FK) |
| --- |
|  |

**FULL_PATH_B**
| PATH<br>STEP<br>CLASS (FK)<br>SUBCLASS (FK) |
| --- |
|  |

**FULL_PATH_M**
| PATH<br>STEP<br>CLASS (FK)<br>SUBCLASS (FK) |
| --- |
|  |

# Spike-in experiment



clustered go_ids: 5381,8490,15342,15344,15359,15620,15621
AVG sim score= 3.863
Cluster root = 15291

| go_id | go_id | score |
|---|---|---|
| 5381 | 8490 | 4.04 |
| 5381 | 15342 | 4.04 |
| 5381 | 15344 | 4.04 |
| 5381 | 15359 | 3.27 |
| 5381 | 15620 | 4.04 |
| 5381 | 15621 | 4.04 |
| 8490 | 15342 | 4.04 |
| 8490 | 15344 | 4.35 |
| 8490 | 15359 | 3.27 |
| 8490 | 15620 | 4.35 |
| 8490 | 15621 | 4.04 |
| 15342 | 15344 | 4.04 |
| 15342 | 15359 | 3.27 |
| 15342 | 15620 | 4.04 |
| 15342 | 15621 | 4.04 |
| 15344 | 15359 | 3.27 |
| 15344 | 15620 | 4.35 |
| 15344 | 15621 | 4.04 |
| 15359 | 15620 | 3.27 |
| 15359 | 15621 | 3.27 |
| 15620 | 15621 | 4.04 |

Five related GO nodes with GOids 5381, 8490, 15344, 15620, and 15621; labeled **red**; were spiked into a randomly selected pool of 20 nodes and subjected to GO clustering. The similarity analysis successfully re-created the set of related GO nodes. Column 1and 2 in the table shows a pair of GO nodes and column 3 shows the pair-wise similarity scores. Nodes colored **pink** (15342, 15359) are from the randomly selected 20 Go nodes and were clustered with the spiked GO nodes. **Green circle** indicates the cluster root (15291), which is the lowest level common ancestor node.

# RA stimulated MPRO cell differentiation time-series experiment

Transgenic Myeloid Progenitor (MPRO) cells transgenic for the dominant negative Retinoic Acid (RA) receptor were induced to differentiate into Neutrophils with high doses of RA.

Gene expression at 0, 1, 2, 4, and 8 hours post RA induction was analyzed with Affymetrix U74Av2 mouse microarray.
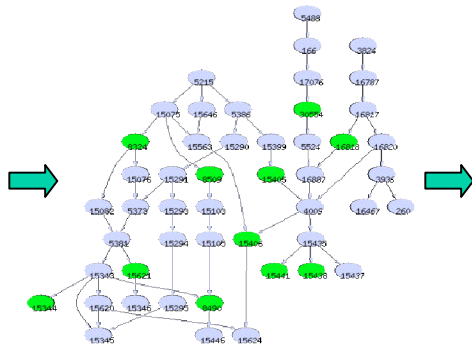
Genes showing significant changes in their expression level across a series of time points are modulated by retinoic acid stimulation and cell differentiation.

We arbitrarily took the top 80 genes based on the F-score ranking.

# GO clustering

# GO clustering on Leukocyte differentiation time-series experiment

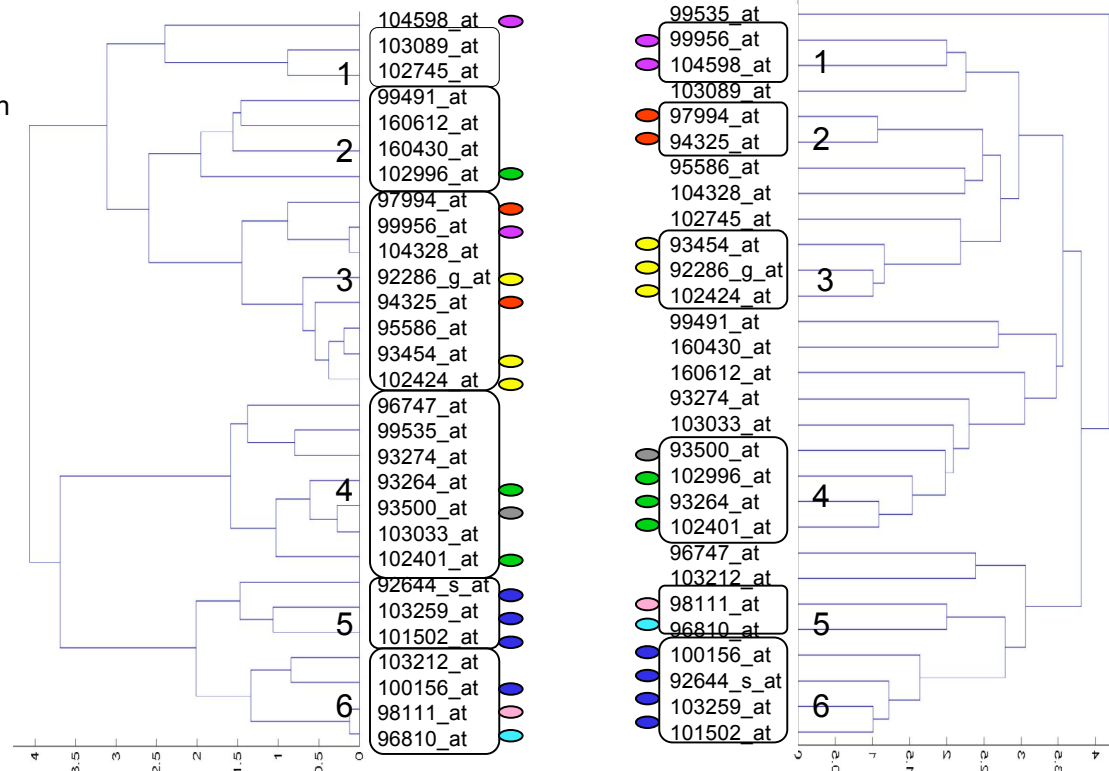| Rank | Title | Score | probe sets | Genes | X1 | n1 | X2 | n2 | Bootstrap p-val (10000) btstrap | Enrichment |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | defense response | 3.403 | 93454_at<br>102424_at<br>92286_g_at<br>102401_at<br>103033_at<br>102745_at | lymphocyte antigen 68<br>small inducible cytokine A3<br>interleukin 4<br>interferon regulatory factor 1<br>complement component 4<br>T-cell receptor gamma | 6 | 29 | 221 | 3163 | 0.0148 | Yes |
| 2 | transcription regulation | 4.04 | 102996_at<br>103259_at<br>102401_at<br>100156_at<br>92644_s_at<br>94325_at<br>93264_at<br>97994_at<br>101502_at | eleven-nineteen lysine-rich leukemia gene<br>growth factor independent 1<br>interferon regulatory factor 1<br>mini chromosome maintenance deficient 5<br>myeloblastosis oncogene<br>pre B-cell leukemia transcription factor 1<br>sterol regulatory element binding factor 1<br>transcription factor 7, T-cell specific<br>TG interacting factor | 9 | 29 | 486 | 3163 | 0.0287 | Yes |
| 3 | steroid metabolism, steroid biosynthesis | 3.862 | 93264_at<br>94325_at | sterol regulatory element binding factor 1<br>pre B-cell leukemia transcription factor 1 | 2 | 29 | 38 | 3163 | 0.0478 | Yes |
| 4 | cell cycle control | 3.05 | 99956_at<br>92644_s_at | kit oncogene,<br>myeloblastosis oncogene | 2 | 29 | 113 | 3163 | 0.2757 | No |
| 5 | cytoskeleton organization and biogenesis | 3.592 | 103212_at<br>96747_at | dynein, axon, heavy chain 11, Wnt1<br>responsive Cdc42 homolog | 2 | 29 | 119 | 3163 | 0.2972 | No |
| 6 | protein modification | 3.265 | 93274_at<br>99956_at<br>104598_at | CDC-like kinase,<br>kit oncogene,<br>PTP non-receptor type 16 | 3 | 29 | 282 | 3163 | 0.4871 | No |

# GO-guided expression clustering

# GO guided clustering on Leukocyte differentiation time-series experiment

Gene clusters where correlations between biological function and expression profile are both evident were identified by GO guided clustering.



Hierarchical clustering

GO-guided hierarchical clustering

# Acknowledgements

John Martin

Melissa Cline

David Finkelstein

Tarif Awad

Michael Stewart

Michael Siani-Rose

David Kulp

# Thank you!