# Getting to know Blast2GO

## Functional annotation:
## from sequences to functional labels

# Outline

Concepts on Functional Annotation:

Biological Databases

Blast2GO annotation strategy

-------------------------------------------------------------------

The Blast2GO annotation framework:

Annotation steps,  Modulation of annotation intensity, Export/Import Functions, Sequence Selection,  Additional Tools

-------------------------------------------------------------------
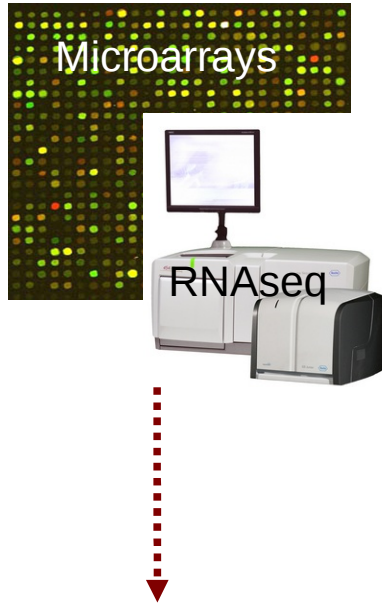
Data Mining: Visualization Techniques

Combined Graph, Charts and Pies

-------------------------------------------------------------------

Hands on: Blast2GO annotation excercise

# Why Blast2GO?

**Experiment with novel Sequences**

Microarrays

RNAseq

**Data-Analysis**
- preprocessing
- clustering
- normalization
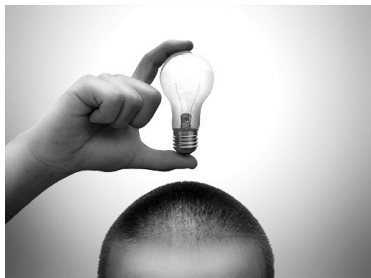- differential expression
- correlation

**Gene-List**

MNAT1
CTNNBL1
ENOX2
GTPBP1
RALY
TAGLN2
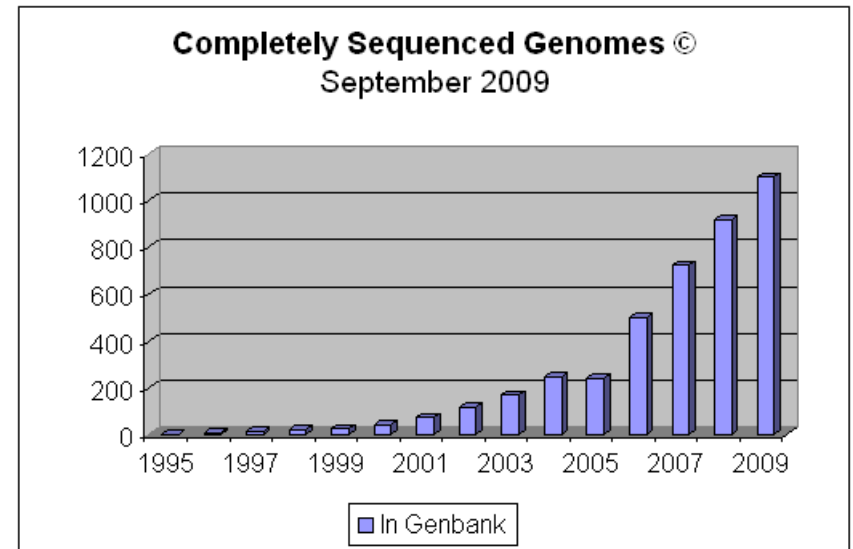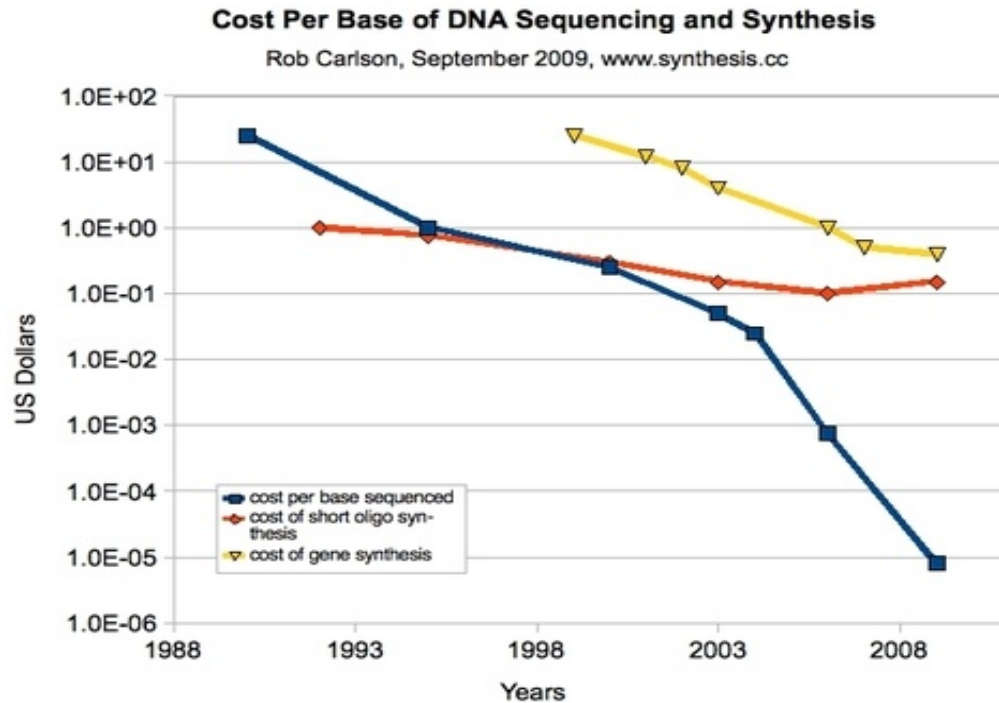RAB3A
PPP2R5A
MAPRE1
. . . . .
. . .

**Functional Annotation**

**+**

**Functional interpretation**

FatiGO

FatiScan

GSEA
Gene Set Enrichment Analysis

snow

KEGG
Kyoto Encyclopedia of Genes and Genomes

**Functional Profiling**

# Why automated functional annotation?



We can not characterize experimentally all these
sequence data at the current growth rate:
We need automated annotation tools to achieve an equivalent throughput

# What is functional annotation?

The function is on the protein

But frequently we annotate nt-sequences

**Functional Label**

Cellular Rol
Expression
Interactions
.....

Controlled Vocabulary

High throughput

Accessible

# Functional Vocabularies


*the Gene Ontology*

**Molecular Function**
**Biological Process**
**Cellular Component**


Kyoto Encyclopedia of Genes and Genomes

**Metabolic pathways**


InterPro proSite

**Functional motifs**

**KEGG orthologues**

### Example proteins

P25024 High affinity interleukin-8 receptor A (IL-8R A) (IL-8 rec

More proteins

IPR000174 Interleukin-8 receptor
IPR000276 Rhodopsin-like GPCR superfamily
IPR001277 C-X-C chemokine receptor, type 4
IPR001355 Interleukin 8A receptor
ModBase
PDB Chain

# The Gene Ontology

✓ Project developed by the Gene Ontology Consortium

✓ Provides a controlled vocabulary to describe gene and gene product attributes in any organism

✓ Lastest version more than 22.000 terms

✓ Includes both the development of the Ontology and the maintenance of a Database of annotations

**http://www.geneontology.org**

# Gene Ontology

**The three categories of GO**

**Molecular Function**

the tasks performed by individual gene products; examples are *transcription factor* and *DNA helicase*

**Biological Process**

broad biological goals, such as *mitosis* or *purine metabolism*, that are accomplished by ordered assemblies of molecular functions

**Cellular Component**

subcellular structures, locations, and macromolecular complexes; examples include *nucleus*, *telomere*, and *origin recognition complex*

GO:0003673 : Gene_Ontology (65883)
  GO:0008150 : biological_process (44405)
    GO:0007610 : behavior (357)
    GO:0000004 : biological_process unknown (7877)
    GO:0009987 : cellular process (32672)
      GO:0007154 : cell communication (5384)
      GO:0008219 : cell death (744)
      GO:0030154 : cell differentiation (464)
      GO:0008151 : cell growth and/or maintenance (28802)
      GO:0006928 : cell motility (911)
      GO:0006944 : membrane fusion (257)
    GO:0016265 : death (793)
    GO:0007275 : development (4615)
    GO:0008371 : obsolete (1581)
    GO:0007582 : physiological processes (31124)
    GO:0016032 : viral life cycle (115)
  GO:0005575 : cellular_component (32869)
  GO:0003674 : molecular_function (53910)

# The Ontology

- Annotations are given to te **most specific** (low) level
- True path rule: annotation at a term implies **annotation to all its parent terms**
- Annotation is given with an **Evidence Code:**
  - **IDA**: inferred by direct assay
  - **TAS**: traceable author statement
  - **ISS**: infered by sequence similarity
  - **IEA**: electronic annotation
  - ….

More general

More specific

biological_process
GO:0008150

is a — is a

cellular process
GO:0009987

physiological process
GO:0007582

is a — is a

cellular physiological process
GO:0050875

is a

cell cycle
GO:0007049

is a — part of

mitotic cell cycle
GO:0000278

M phase
GO:0000279

part of — is a

M phase of mitotic cell cycle
GO:0000087

part of

mitosis
GO:0007067

GO-Description:
The division of the eukaryotic cell nucleus to produce two daughter nuclei that, usually, contain the identical chromosome complement to their mother.

# The Gene Ontology Database (GOA)

http://www.geneontology.org/GO.current.annotations.shtml

✓ There is a collaborating institution per organism to provide annotations

✓ Most of the GOA annotations come from UniProt

✓ Most of the annotations are electronic annotations

Some numbers of the actual version (jan11):
Terms:  31.794, mf(8.912), bp(20.092), cc(2.790)
Genes: 8.864.425
Annotations: 52.571.310

# The GO has a DAG structure



molecular_function - combined Graph

# KEGG pathways

## 1. Metabolism

### 1.1 Carbohydrate Metabolism
Glycolysis / Gluconeogenesis
Citrate cycle (TCA cycle)
Pentose phosphate pathway
Pentose and glucuronate interconversions
Fructose and mannose metabolism
Galactose metabolism
Ascorbate and aldarate metabolism
Starch and sucrose metabolism
Aminosugars metabolism
Nucleotide sugars metabolism
Pyruvate metabolism
Glyoxylate and dicarboxylate metabolism
Propanoate metabolism
Butanoate metabolism
C5-Branched dibasic acid metabolism
Inositol metabolism
Inositol phosphate metabolism

### 1.2 Energy Metabolism
Oxidative phosphorylation
Photosynthesis
Photosynthesis - antenna proteins
Carbon fixation
Reductive carboxylate cycle (CO2 fixation)
Methane metabolism
Nitrogen metabolism
Sulfur metabolism

### 1.3 Lipid Metabolism
Fatty acid biosynthesis
Fatty acid elongation in mitochondria
Fatty acid metabolism
Synthesis and degradation of ketone bodies
Biosynthesis of steroids
Bile acid biosynthesis

## Current Statistics

**KEGG Release 48.0+/10-01, Oct 08**

**KEGG PATHWAY**   90,787 pathways generated from 251 reference pathways

**KEGG BRITE**   17,388 hierarchies generated from 54 reference hierarchies

**KEGG GENES**   11,213 KO groups
4,016,008 genes in 99 eukaryotes + 708 bacteria + 52 archaea

**KEGG LIGAND**   15,221 compounds, 7,844 drugs, 10,966 glycans, 7,780 reactions, 10,154 reactant pairs

**http://www.genome.jp/kegg/**

VII International Course of Massive Data Analysis

# KEGG pathways

# InterPro

http://www.ebi.ac.uk/interpro/databases.html

- ✓ Collection of databases with functional annotation of protein motifs
- ✓ Functional vocabulary at UniProt
- ✓ There is an equivalence table between GO and InterPro

# InterPro Scan

# Functional assignment

# Function transfer

# Annotation by similarity: concerns

$GO_1$, $GO_2$, $GO_3$, $GO_4$

**HIT**

$GO_1$, $GO_2$, $GO_3$, $GO_4$

**QUERY**

- ✓ Level of homology (~ from 40-60% is possible)
- ✓ The overlap query and hit sequences (not much a problem)
- ✓ The domain or structure function association
- ✓ The paralog problem: genes with similar sequences might have different functional specifications
- ✓ The evidence for the original annotation
- ✓ Balance between quality and quantity: depends on the use

# Blast2GO

- ✓ Suite for functional annotation and data mining on functional data
  - o Considerations for **annotation**
    - Simlarity
    - Length of the overlap
    - Percentage of hit sequence spanned by the overlap
    - Evidence original annotation
    - Blast hits and motif hits
    - Refinement by additional methods
  - o Visualization:
    - Annotation charts
    - **Knowledge discovery on the DAG**
- ✓ Desktop Java application

# Scientific Impact



Blast2GO citations per year



Scope Papers Citing Blast2GO

# Blast2GO Annotation strategy

# Blast2GO Annotation Strategy



| | |
|---|---|
| Sq1 | go1,go2 |
| | go1,go3 |
| | go3 |
| | go1 |
| Sq2 | go8 |
| | go8 |
| | go8 |
| Sq3 | go2 |
| | go2 |
| | go2 |
| | go2,go4 |
| Sq4 | |

**Refinement**

InterPro
Annex
GOSlim
Manual

| | |
|---|---|
| Sq1 | go1 go2, go3, GO11 |
| Sq2 | go8, GO12, GO13 |
| Sq3 | go2,go4 |
| Sq4 | GO15 |

# Blast2GO annotation rule

**Quality of annotation source**

| EC | weight |
|----|--------|
| IC | 1 |
| TAS | 1 |
| IDA | 1 |
| IMP | 0.9 |
| IGI | 0.9 |
| IPI | 0.9 |
| ISS | 0.8 |
| IEP | 0.8 |
| NAS | 0.7 |
| IEA | 0.7 |
| ND | 0.5 |
| NR | 0.5 |
| RCA | 0.5 |

Evidence codes

**Possibility of abstraction**

**Similarity requirement**

$$sim = \frac{\sum positives_{hsp}}{\sum alignmentlength_{hsp}}$$

**[max.(sim x ECw) + (#terms-1 x GOw)**

# Blast2GO annotation rule

Lowest term above threshold

Similarity requirement

$$sim = \frac{\sum positives_{hsp}}{\sum alignmentlength_{hsp}}$$

Quality of annotation source

| EC | weight |
|----|--------|
| IC | 1 |
| TAS | 1 |
| IDA | 1 |
| IMP | 0.9 |
| IGI | 0.9 |
| IPI | 0.9 |
| ISS | 0.8 |
| IEP | 0.8 |
| NAS | 0.7 |
| IEA | 0.7 |
| ND | 0.5 |
| NR | 0.5 |
| RCA | 0.5 |

Evidence codes

Possibility of abstraction

Recall vs. Precision

**Lowest.node [max.(sim x ECw) + (#terms-1 x GOw) >= threshold]**

VII International Course of Massive Data Analysis

# Blast2GO annotation rule

**Lowest.node [max.(sim x ECw) + (#terms-1 x GOw) >= threshold]**

- When I have a GO with ECw = 1 and I do not allow abstraction (GOw = 0), then the
  Annotation Score = %similarity

- If the ECw < 1 my similarity requirement is higher to obtain the same Annotation Score

- If I allow abstraction GOw > 0, then with less similarity I can obtain the required Annotation Score at a parent node

# Outline

Concepts on Functional Annotation:

Biological Databases

Blast2GO annotation strategy

-------------------------------------------------------------

The Blast2GO annotation framework:

Annotation steps,  Modulation of annotation intensity, Export/Import Functions, Sequence Selection,  Additional Tools

-------------------------------------------------------------

Data Mining: Visualization Techniques

Combined Graph, Charts and Pies

-------------------------------------------------------------

Hands on: Blast2GO annotation excercise

# Start Blast2GO

**www.blast2go.org**

- Desktop application

- Java WebStart

- Internet connection

# Java and Java Web Start

- SUNs Java Runtime Environment (1.6)

- Java Web Start, a technology to stay always up to date

- Activate the Java Console for debugging

- Create a desktop short cut

- Define the memory B2G can use

# Input data
(in FASTA format, AA or nt)

>my_favourite_species_seq1 | still unknown
gtgatggaaaagaaaagttttgttatcgtcgacgcatatgggtttcttttcgcgcgtattatgcgctgcctggattaagcacctcatacaattttcctgtaggaggtgtatatggttt
tataaacatacttttgaaacatctctcttccacgatgcagattatttagttgtggtatttgattcggggtcgaaaaattttcgtcacactatgtattccgaatacaaaactaatcgc
cctaaagcaccagaggatctgtcactacaatgtgctccgctacgtgaggctgttgaagcgtttaatattgtaagtgaagaagtgcttaactacgaagcagacgacgtaata
gctacactctgtacaaaatatgcatctagtaatgttggagtgagaatactgtcagcagataaggatttactacaactcctaaatgataatgttcaagtttacgaccctataaaa
agcagatacctcaccaatgaatacgttttagaaaaatttggtgtttcatcagataagttgcatattgatacggttgcatcgagttataatgagaaaattattctcagctaagctgt
acaccgtttattacacactcgaaaggccgttag
>my_favourite_species_seq2 | no clue
ttgttagctaaaaaggaagactttcacacctttggtaatggtgttggctctgctggaacaggtggagttgtagtttctgcatccatgttgtctgcggatttttcaaatcttagagaag
agatagcagcggttagtacggctggtgcagattggttacacattgatgtgatggatgggtgcttcgtccccagtttgactatgggtcctgtggtgatttccggcattaggaaatg
tacaaatatgtttcttgatgtgcatttgatgattaatcgcccaggcgatcatctgaagagtgtggtagatgctggagctgataagatagagcacattcgcaagatgatagagg
aaagctcatcaaccgcgaaaatcgctgttgatggtggtgtttcaacggataatgcccgggctgttatcgaggcaggtgcgaatatactcgttgttggaacggcgctgtttgct
gctgacgatatgagtaaagttgtaagaactttaaaatcattttaa
>my_favourite_species_seq3 | just sequenced
gtgggactgctcatccctgtaggcagggtggctattttttgtgtaaaggcagtctttcatagtcttgtaccgccatactatctatggataactacaaagcagttttttgaggtgtggt
ttttctctcttcctatagtagcagttacatctttgtttacgggaggcgcgttagcccttcaggatacctcgtgggaagcgctaaagtatcagggtaatggagttttttactcctgca
agatgtaatagagggtctggtaaaagctgtatcgtttgggctggtaatttcgctagttgggtgttacaacgggtatcactgtgagataggcgcaaggggtgtaggaacagcg
acaacaaaaacttcggtagcagcttctatgctcataattttgttaaactatataattactgtttttttacgcgta
>my_favourite_species_seq4 | we will see soon...
atgtacgctgtatctctttcaaatttgcatgtctctttcaacaacaaggaggttttgaaaggtgttgacttggacatagcatggggggattccctggttatactgggagaatctggt
agtggaaagtctgtactaacaaaggttgtattgggtctaatagtgccccaagagggaagtgttactgtagatggcaccaatattcttgagaataggcagggcatcaagaat
tttagtgttttgtttcaaaactgtgcgttatttgacagtcttacgatttgggaaaatgtagtattcaatttccgtaggaggcttcgtttagataaggataatgccaaggctttggcttta
cggggattggagcttgtgggattggacgccagtgtaatgaacgtgtatcctgtggagctatcaggcgggatgaaaaagcgcgtagctttggcaagagctattataggtagt
cccaaaattctaattttggatgagccaacttcgggattggatcctataatgtcttcagtggt

# Blast2GO Application



(1) Blast

(2) Mapping

(3) Annotation

Any operation will only affect to selected sequences!!!!

Main Sequence Table

Application statistics

Blast results

Application messages

Graph visualisation

# The First Check



Click on the green arrow to check you can connect to DB
A GO graph should appear

# Database configuration



**FOR TODAY: mem20**

Open port 3306 (mysql) for outgoing connections at your institute

Configure/check personal firewalls

Actual settings can be found at www.blast2go.org

# Load Sequences

# Run BLAST search



BLAST against NCBI or locally
Choose different DBs

In combination with URL

Limit to query-hit overlap

Recommended to save as XML

Text mining on BLAST hit description

# Choose other DB at NCBI



Set at blast2go.properties file

# BLAST Results

# Blast Distribution Charts



Evaluate the similarity of
your sequences with public DBs

# Single Sequence Menu

# Mapping Results

# Resources for mapping

Gene Ontology Database

NCBI data-files:
    gene2accession (4 079 414 entries)
    gene_info (1 635 614 entries)

Protein Information Resource (PIR):
    Non-Redundant Reference Protein Database including
    PSD, UniProt, Swiss-Prot, TrEMBL, RefSeq, GenPept
    and PDB

BLAST Hit IDs ACCs/GIs Gene-Symbols → Mapping Resources → GO terms EC sim %

# Annotation Menu

BLAST based annotation



Other Annotation modes

Validation and Annex

# Annotation



Allows to set a minimum percentage of the HIT sequence
which should be expand by the QUERY sequence

This helps to avoid the problem of cis-annotation

# Annotation Result

# Graph Visualization

# Annotation Charts

# Annotation Charts



Commonly, level 5 is the most abundant specificity level in the Gene Ontology

# Additional Annotation: ANNEX

Recovers implicit biological process and cellular component GO terms based on molecular function annotations



**Myhre *et al*, Bioinformatics 2006**

# Additional Annotation: InterProScan

Runs InterProScan searches at the EBI through Blast2GO

| Your email address: | | |
|---|---|---|
| **Choose applications to run:** | | |
| BlastProDom: | ☐ | |
| FPrintScan: | ☑ | |
| HMM-PIR: | ☑ | |
| HMM-Pfam: | ☑ | |
| HMM-Smart: | ☑ | |
| HMM-Tigr: | ☑ | |
| ProfileScan: | ☑ | |
| PatternScan: | ☑ | |
| SuperFamily: | ☑ | |
| Gene3D: | ☑ | |
| HMM-Panther: | ☑ | |
| SignalP: | ☑ | |
| TM-HMM: | ☑ | |

| Annotation | Analysis | Statistics | Sele |
|---|---|---|---|
| **Run Annotation Step** | | | |
| **Set Evidence Code Weights** | | | |
| **Reset Annotation** | | | |
| **Validate Annotations** | | | |
| **Remove 1.Level Annotations** | | | |
| **Run ANNEX (Annotation Augmentation)** | | | |
| **InterProScan** ▶ | | | |
| **Enzyme Code and KEGG** ▶ | | | |
| **GO-Slim** ▶ | | | |

- Run InterProScan (online)
- Stop InterProScan
- Import InterProScan Results (XMLs)
- Reset InterProScan Results
- Merge InterProScan GOs to Annotation

Results are stored at your computer as XML files. You can upload them later

Once you have completed your InterPro annotation, results can be transformed to GO terms and merged to Blast annotation

# InterProScan Results

# Additional Annotation: GOSlim



GOSlim is a reduction of the Gene Ontology to a more reduced vocabulary → Helps to summarize information

After GOSlim transformation sequences get YELLOW



**Different GOSlims available at Blast2GO**

# Enzyme annotation and Kegg Maps

## GO → Enzyme Codes → KEGG maps

# Additional Annotation: Manual Curation



You can modify manually annotation of particular sequences

If you click in this box, curated sequences get purple

# Export Results



Saves the complete B2G project (heavy)

Export annotation results in different formats

# Export formats

## .annot

```
C04018C10   GO:0004707      mitogen-activated protein kinase 3
C04018C10   EC:2.7.11.24
C04018A12   GO:0016798      class iv chitinase
C04018A12   GO:0000272
```

Also for import!

## GeneSpring Format

| | | | |
|---|---|---|---|
| C04013E10 | response to water deprivation; regulation c | nucleus; | transcription factor activity; |
| C04013A12 | translation; | ribosome; plastid; | structural constituent of ribosome; |
| C04013C12 | galactose metabolic process; | plastid; | aldose 1-epimerase activity; carbohydrate binding; |

## GoStat

```
C04018C10                   4707,9409,6979,10200,5524,169
C04018A12                   16798,272,44248
C04018C12                   4869,12505,8233
```

## By Seq

| | | | |
|---|---|---|---|
| C04018A02 | glyoxalase i | GO:0004462 | F:l |
| C04018C02 | metallothionein-like protein | GO:0046872 | F:r |
| C04018G02 | protein phosphatase | GO:0008287 | C:| |

# More export formats

## Export Sequence Table

| Seq. Name | Seq. Description | Seq. Length | #Hits | min. eValue | mean Similarity | #GOs | GOs | Enzyme Codes | InterProScan |
|---|---|---|---|---|---|---|---|---|---|
| C04018C12 | cysteine proteinase inhibitor | 663 | 20 | 25 | 80.00% | 3 | F:GO:0004869; C:GO:0012505; F: | | IPR000010; IPR01 |
| C04018E12 | protein phosphatase 2c | 663 | 20 | 77 | 85.00% | 2 | N:GO:0015071; F:GO:0003824 | | IPR001932; IPR01 |
| C04018G12 | alpha beta fold family protein | 578 | 20 | 84 | 79.00% | 4 | F:GO:0016787; C:GO:0005739; C: | | noIPR |
| C04018A02 | glyoxalase i | 600 | 20 | 64 | 74.00% | 2 | P:GO:0005975; F | EC:4.4.1.5 | IPR004360; noIPR |
| C04018C02 | metallothionein-like protein | 625 | 18 | 14 | 74.00% | 1 | F:GO:0046872 | | IPR000347 |
| C04018E02 | haemolysin-iii related familyex | 612 | 20 | 32 | 72.00% | 1 | C:GO:0016020 | | noIPR |
| C04018G02 | protein phosphataseexpressed | 645 | 20 | 97 | 81.00% | 5 | C:GO:0008287; N:GO:0015071; P: | | no IPS match |
| C04018C04 | phosphoglycerate bisphospho | 780 | 20 | 63 | 66.00% | 2 | P:GO:0008152; F:GO:0003824 | | IPR001345; IPR01 |
| C04018E04 | polyubiquitin | 707 | 20 | 115 | 99.00% | 2 | P:GO:0006464; C:GO:0005622 | | IPR000626; IPR01 |
| C04018G04 | meiotic recombination 11 | 575 | 20 | 45 | 89.00% | 21 | C:GO:0019013; P:GO:0007126; F: | | IPR003701; IPR00 |
| C04018A06 | late embryogenesis-abundant | 648 | 20 | 43 | 68.00% | 2 | P:GO:0009737; P:GO:0009409 | | no IPS match |

## Export BestHit Data

| Sequence name | Sequence desc. | Sequence length | Hit desc. | Hit ACC | E-Value | Similarity | Score | Alignment length | Positives |
|---|---|---|---|---|---|---|---|---|---|
| C04018C10 | mitogen-activated protein | 717 | gi|122894104|gb|ABM6769 | ABM67698 | 1.35E-123 | 99 | 445.28 | 222 | 221 |
| C04018E10 | ---NA--- | 706 | gi|157356307|emb|CAO624 | CAO62459 | 2.69E-036 | 83 | 155.22 | 119 | 99 |
| C04018G10 | protein | 620 | gi|114153154|gb|ABI52743. | ABI52743 | 7.47E-015 | 63 | 83.57 | 90 | 57 |
| C04018A12 | class iv chitinase | 715 | gi|3608477|gb|AAC35981.1 | AAC35981 | 1.45E-061 | 78 | 239.2 | 171 | 134 |
| C04018C12 | cysteine proteinase inhil | 663 | gi|8099682|gb|AAF72202.1| | AAF72202 | 9.33E-025 | 83 | 116.7 | 99 | 83 |
| C04018E12 | protein phosphatase 2c | 663 | gi|46277128|gb|AAS86762. | AAS86762 | 2.76E-077 | 91 | 291.2 | 180 | 164 |
| C04018G12 | alpha beta fold family pr | 578 | gi|147865769|emb|CAN832 | CAN83251 | 1.67E-084 | 94 | 314.69 | 179 | 169 |
| C04018A02 | glyoxalase i | 600 | gi|2213425|emb|CAB09799 | CAB09799 | 2.16E-064 | 81 | 248.05 | 114 | 93 |
| C04018C02 | metallothionein-like prote | 625 | gi|3308980|dbj|BAA31561.1 | BAA31561 | 2.23E-014 | 100 | 82.03 | 40 | 40 |

# Sequence Selection



Sequence Selection tool to obtain a selection based on annotation status

# Sequence Selection

Select | Tools | View | Info

Select by Color
Select by Sequence Name or ID
Select by Sequence Description
Select by Function (GO-Terms or GO-IDs)

Invert Selection
Delete Selected Sequences

Order Sequences by Selection
Restore Inicial Table Order

**By Name/Description**

From File:
Select/Unselect: ✔
GO-Terms/GO-IDs: ✔
Exact match: ✔
Case sensitive: ✔
Include GO parents: ✔
Functions:

**By Function**

From File:
Select/Unselect: ✔
Exact match: ✔
Case sensitive: ✔
Sequences Names/IDs:

# View Menu



Functions to switch between displaying IDs or descriptions for GO annotation or InterPro results

# Other Tools

| Tools | View | Info |
|---|---|---|

**Invert selection**

**Delete Sequence Selection** — Permits to reduce the project size

**Run Blast-Description-Annotator (BDA)**

**Recover original Best-Blast-Hit descriptions** — Manipulation of sequence desc.

**Add .dat to existing Project (beta)**

**Add .annot annotations to the sequences of a existing Project (beta)** — Merging .annot and .dat projects

**Search loaded annotations in another annotation set**

**Calculate dissimilarity/homogenity (GoetzScore) of selected sequences**

**Start JAVA memory monitor** — Get more out of your memory

**Force to free unused memory**

**Clear properties cache**

**Import PIR Mapping to a local B2G-DB** — Check when connection problems

**DB configuration**

| | |
|---|---|
| DB Host | 193.144.127.204 |
| DB Name | b2g_apr |
| DB User | blast2go |
| DB Password | blast4it |

# Outline

Concepts on Functional Annotation:

Biological Databases

Blast2GO annotation strategy

--------------------------------------------------------------------

The Blast2GO annotation framework:

Annotation steps,  Modulation of annotation intensity, Export/Import Functions, Sequence Selection,  Additional Tools

--------------------------------------------------------------------

## Data Mining: Visualization Techniques

Combined Graph, Charts and Pies

--------------------------------------------------------------------

## Hands on: Blast2GO annotation excercise

# Data Mining on the DAG

- ✓  When working with large datasets, annotation results need to be summarized

- ✓  The DAG provides visualization of annotation data within its biological context

- ✓  In Blast2GO --> *Combined Graph* Function

# Combined Graph

Each term has a number of sequences associated



Nodes can be coloured to indicate relevance

Each term is displayed around its biological context

Node shape to differentiate between direct and indirect annotation

# Combined Graph

# Combined Graph

Let's paint the DAG of the dataset of 1000 sequences



Too many nodes!!!

Need way to find relevant information

# Node information content



**Accumulated by node (Sequence Count)**

**Incomming information (Node Score)**

$$\sum_{g\in desc(g')} seq(g) * \alpha^{\ dist\ (g,\ g')}$$

The node score that reflects the amount of direct information at the node

# Node score



$$\sum_{g \in desc(g')} seq(g) * \alpha^{dist\,(g,\,g')}$$

$$\alpha = 0.6$$

NodeScore (GO1) = $\mathbf{1}$ * $0.6^{0}$ = $\mathbf{1}$

NodeScore (GO2) = $\mathbf{3}$ * $0.6^{0}$ = $\mathbf{3}$

NodeScore (GO3) = $\mathbf{1}$ * $0.6^{1}$ + $\mathbf{3}$ * $0.6^{1}$ = 0.6 + 1.8 = $\mathbf{2.4}$

NodeScore (GO4) = $\mathbf{1}$ * $0.6^{2}$ + $\mathbf{3}$ * $0.6^{2}$ + $\mathbf{1}$ * $0.6^{0}$ = 0.36 + 1.08 + 1 = $\mathbf{2.5}$

# Filtered Graph



# Filtered Nodes

Transition nodes

Direct annotations

# Compacting Graphs by GOSlim

# Show node content

# Saving Options



Save as picture and as txt

| Level | GO ID | Term | Type | #Seqs | Graph Score | Sequences |
|---|---|---|---|---|---|---|
| 4 | GO:0016052 | carbohydrate catabolic process | biological_process | 36 | 4.67 | C02009A12, C04019G12, |
| 7 | GO:0043687 | post-translational protein modification | biological_process | 62 | 34.82 | C04016C08, C08010E08, |
| 3 | GO:0016043 | cellular component organization and | biological_process | 242 | 51.93 | C04018F11, C18004G08, |
| 5 | GO:0051252 | regulation of RNA metabolic process | biological_process | 25 | 13.77 | C04016E04, C04013G11, |
| 4 | GO:0006725 | aromatic compound metabolic proce | biological_process | 44 | 28.64 | C08012A08, C02016F08, |
| 4 | GO:0046907 | intracellular transport | biological_process | 38 | 24.57 | C18004G08, C02009C02, |
| 8 | GO:0006094 | gluconeogenesis | biological_process | 21 | 21 | C04013C12, C07009E02, |
| 4 | GO:0006519 | amino acid and derivative metabolic | biological_process | 78 | 34.82 | C02016F08, C04013E11, |
| 3 | GO:0009719 | response to endogenous stimulus | biological_process | 48 | 20.37 | C18004D02, C18002H02, |
| 5 | GO:0007047 | cell wall organization and biogenesis | biological_process | 15 | 5.75 | C02015B04, C04018D06, |
| 4 | GO:0044248 | cellular catabolic process | biological_process | 75 | 23.13 | C08012A08, C08011C08, |

# Graph Charts

# Graph Charts

· Sequence Distribution/GO as Bar-Chart

· Sequence Distribution/GO as Level-Pie (level selection)

· Sequence Distribution/GO as Multilevel-Pie (#score or #seq cutoff)
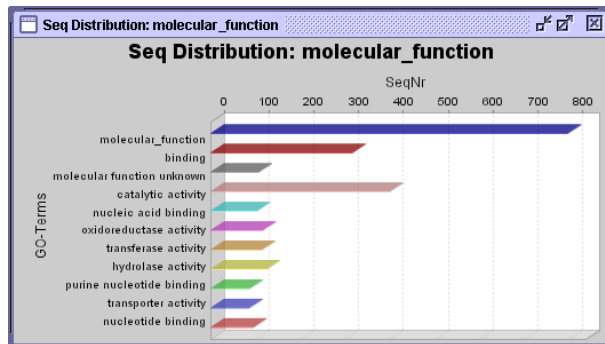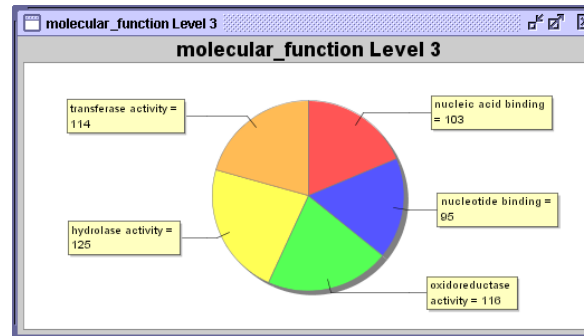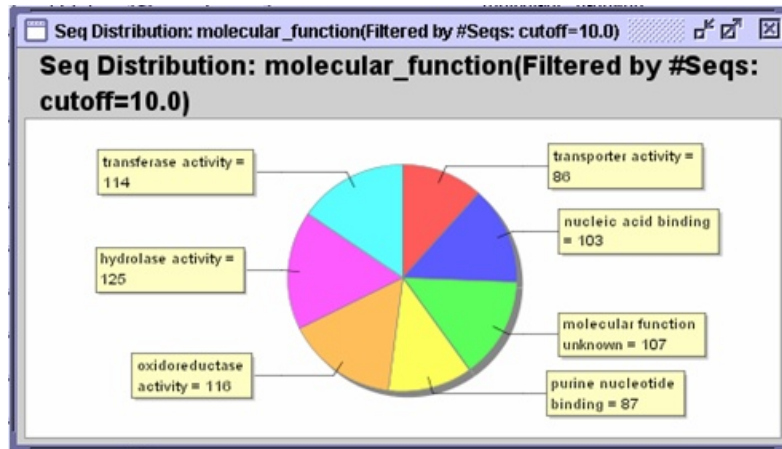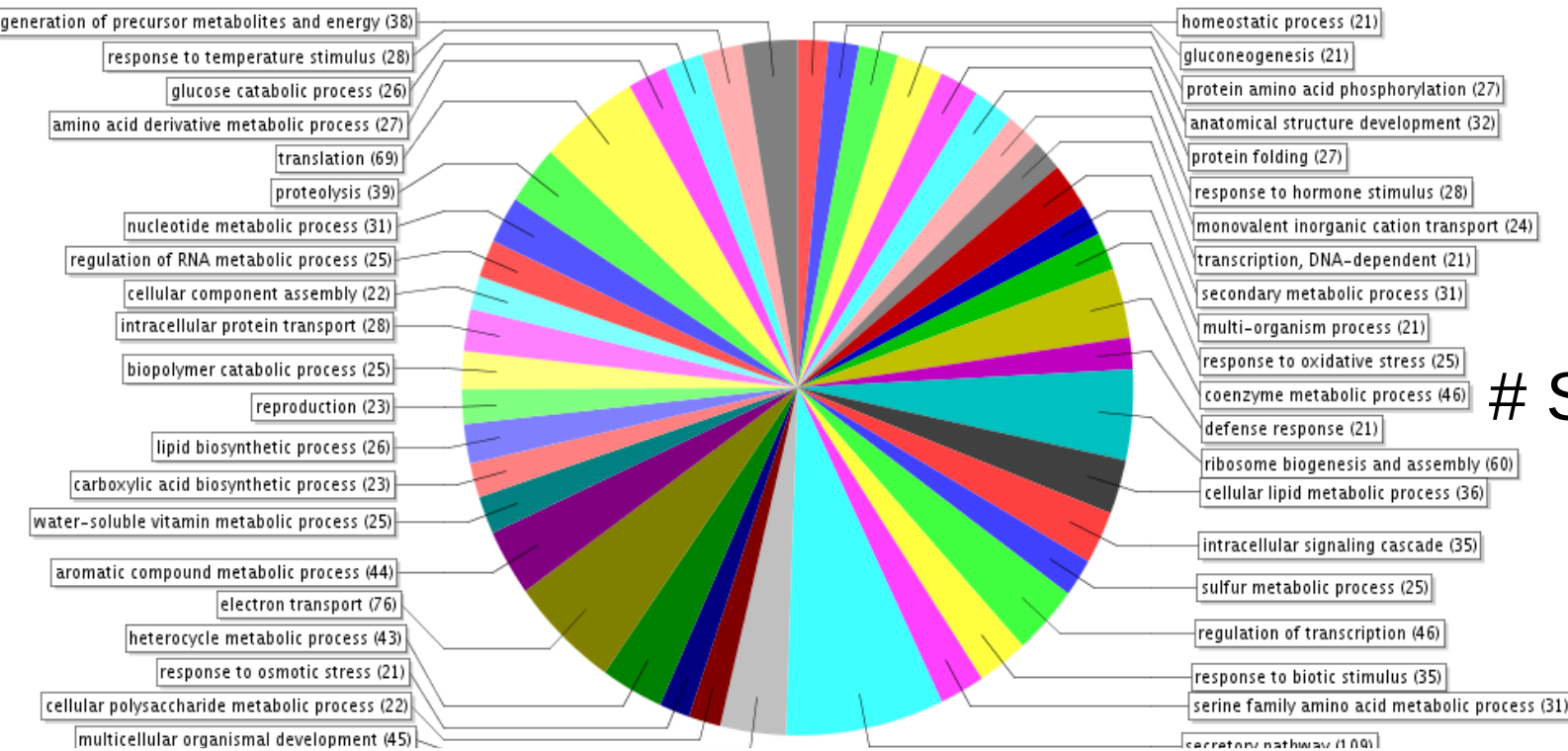
**M-Pies**

# Seq filter = 20

Handy to summarize functional content

GOSlim

# Colouring yourself the DAG

| Tree Type | ● Process | ○ Function | ○ Compon… |
|---|---|---|---|

| | |
|---|---|
| Seq Filter | 14 |
| Node Information | Hide |
| Mode of Graph-Colouring | byDesc |
| Score alpha | 0.6 |
| Node Score Filter | 14 |
| Graph Title Text | Combined Graph |

The byDesc option in the Graph-Colouring allows you to colour the DAG nodes according to an additional value

| | | |
|---|---|---|
| GO:0005792 | GO:0005792 | 1.00 |
| GO:0006412 | GO:0006412 | 0.81 |
| GO:0003735 | GO:0003735 | 0.71 |
| GO:0016705 | GO:0016705 | 0.65 |
| GO:0005840 | GO:0005840 | 0.65 |
| GO:0005506 | GO:0005506 | 0.64 |
| GO:0006631 | GO:0006631 | 0.61 |
| GO:0020037 | GO.. | |

The "special" .annot file:
3 columns
GO name, GO ID, Value

Scale between 0 and 1 used to colour the graph

# Graph Visualization

- ✓ DAGs are interesting for browsing functional annotation but can be too large

- ✓ With filtering and prunning options you can create more navegable DAGs

- ✓ Pies are good to compact information: try out levels

- ✓ GOSlim compacts to more equivaent terms than filtering the GO

# HANDS ON B2G

Go to the on-line course material

Blast, map and annotate several few sequences in Blast2GO by loading the 10 test sequences (within the file menu).

Generate some singel-Seq GO graphs to review annotation. (right mouse click on sequence table)
(http://www.blast2go.org → Start → 1024MB)

Annotated 1100 Citrus-Unigenes (nt) with Blast2GO. Analyse the annotation results. Generate a Combined Graph after a GoSlim-Reduction and try to export a handy graph as PDF.