



International course of

Massive Data Analysis

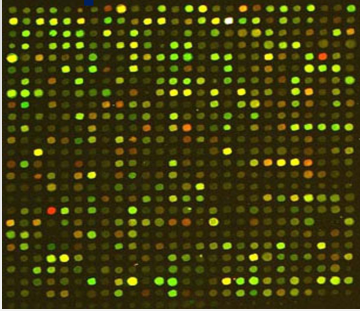


Functional annotation with Blast2GO

Ana Conesa
Stefan Götz

Why Blast2GO?

Experiment



Data-Analysis



GEPAS
GENE EXPRESSION PATTERN ANALYSIS SUITE

Gene-List

MNAT1
CTNBL1
ENOX2
GTPBP1
RALY
TAGLN2
RAB3A
PPP2R5A
MAPRE1
.....
....

Functional Annotation

+



Functional Interpretation



Functional Profiling



What does Blast2GO do?



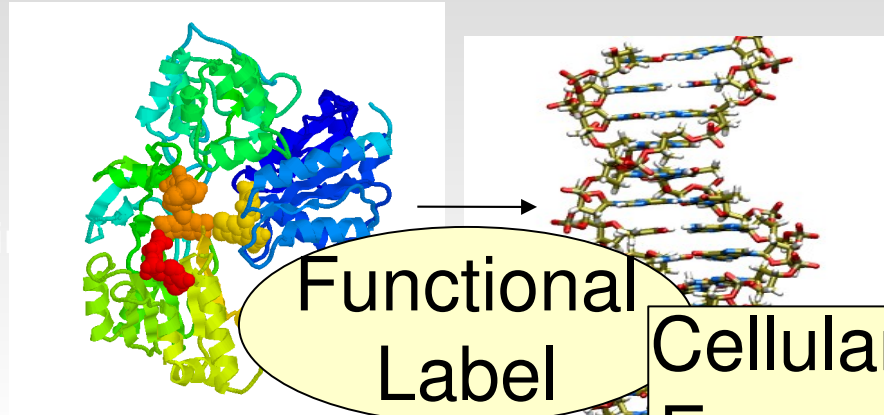
Generates annotations

Visualization of functional annotations



What is functional annotation?

The function
on the protein



But frequently we
annotate nt sequences

Cellular Role
Expression
Interactions
.....

Controlled
Vocabulary

High
throughput

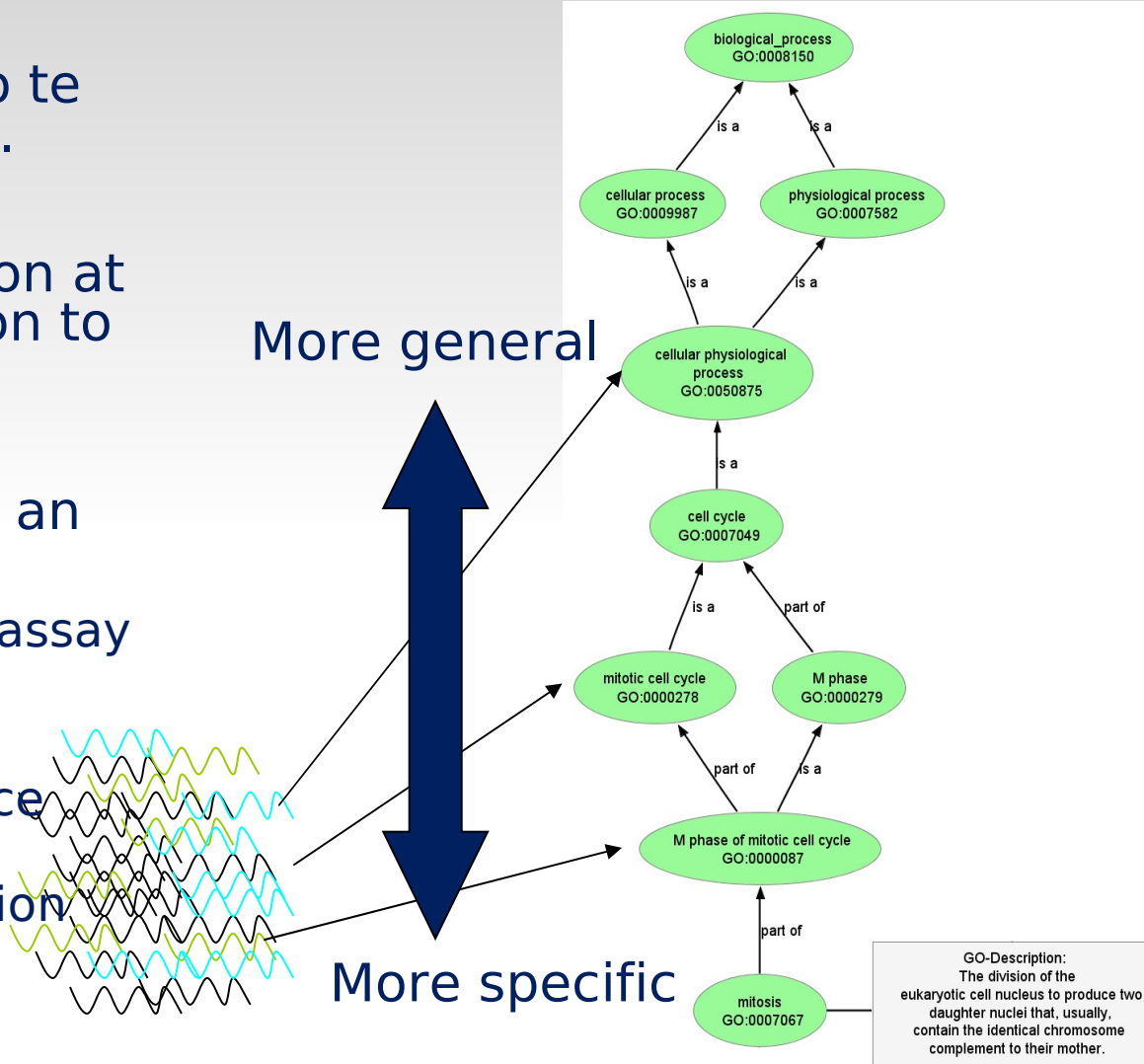
Accessible

The Gene Ontology

- ✓ Project developed by the **Gene Ontology Consortium**
- ✓ Provides a **controlled vocabulary** to describe gene and gene product attributes in **any organism**
- ✓ Includes both the development of the **Ontology** and the maintenance of a **Database** of annotations

The Ontology

- ✓ Annotations are given to the **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**: inferred by sequence similarity
 - **IEA**: electronic annotation
 -

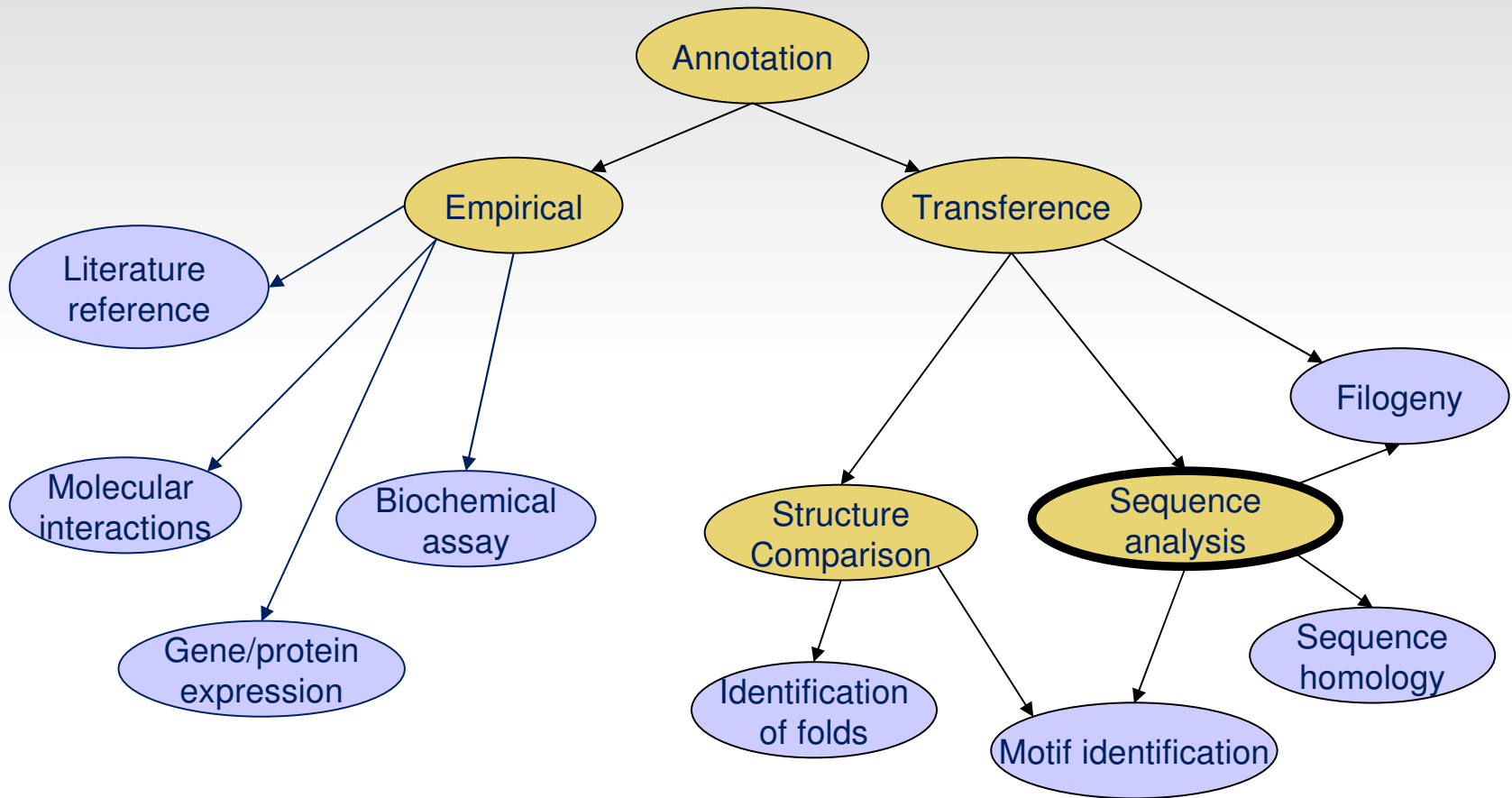


The Gene Ontology Database (GOA)

<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**

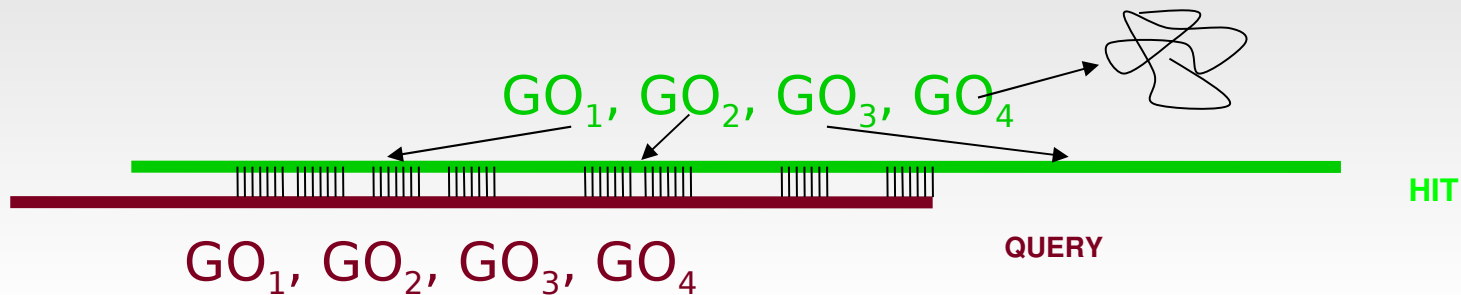
Functional assignment



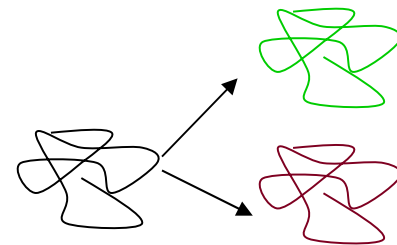
Automatic annotation

- ✓ GO annotations can be created by **comparison** to annotated sequences
- ✓ To achieve enough coverage, high-throughput, **automatic** annotation is required
- ✓ The most effective (also error prone) automatic annotation method is transfer from **sequence similarity**

Concerns in functional transfer by similarity



- ✓ Level of **homology** (~ from 40-60% is possible)
- ✓ The **overlap** query and hit sequences
- ✓ 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

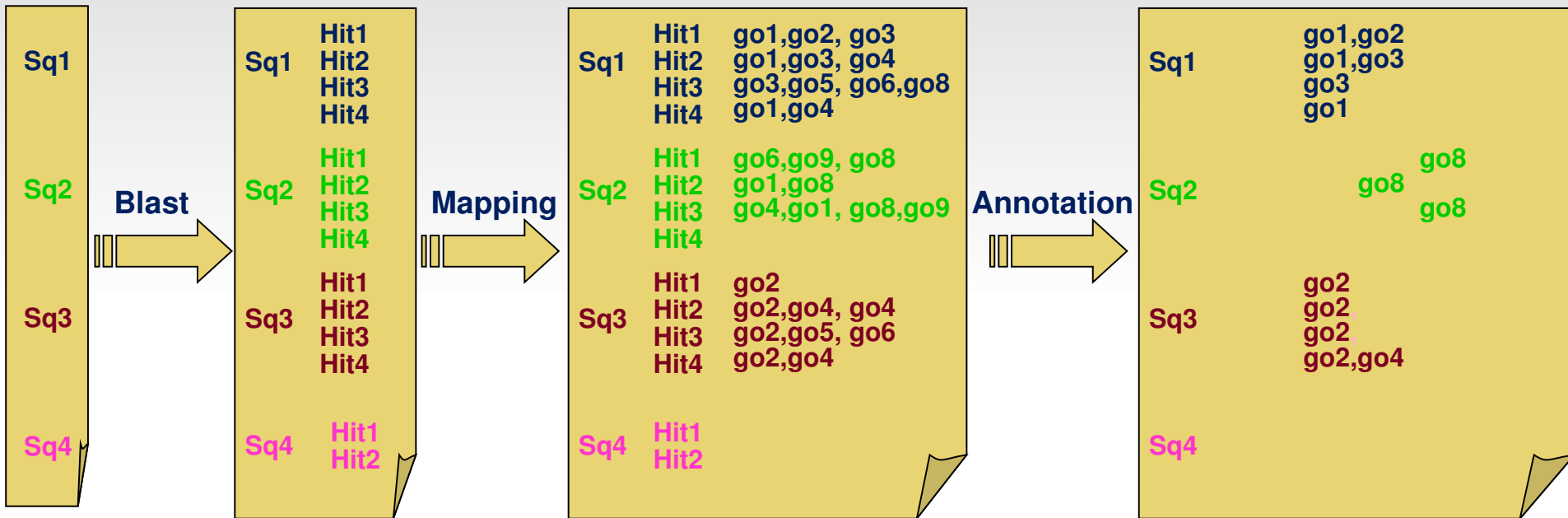


Blast2GO

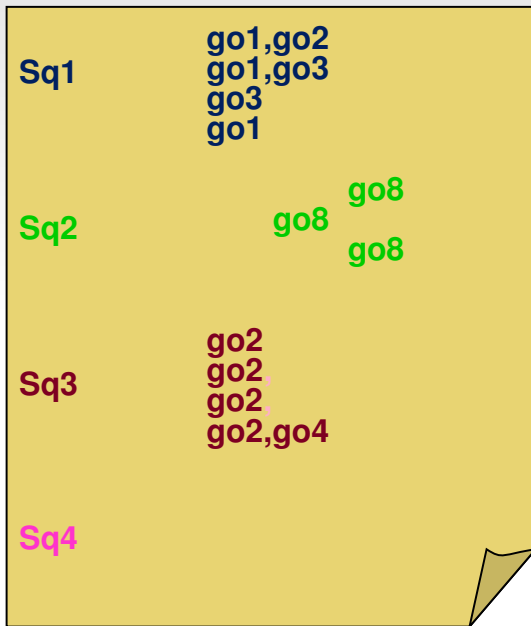
- ✓ Suite for functional annotation and data mining on functional data
 - Considerations for **annotation**
 - Similarity
 - Length of the overlap
 - Percentage of hit sequence spanned by the overlap
 - Evidence original annotation
 - Blast hits and motif hits
 - Refinement by additional methods
 - Visualization: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



Blast2GO Annotation strategy



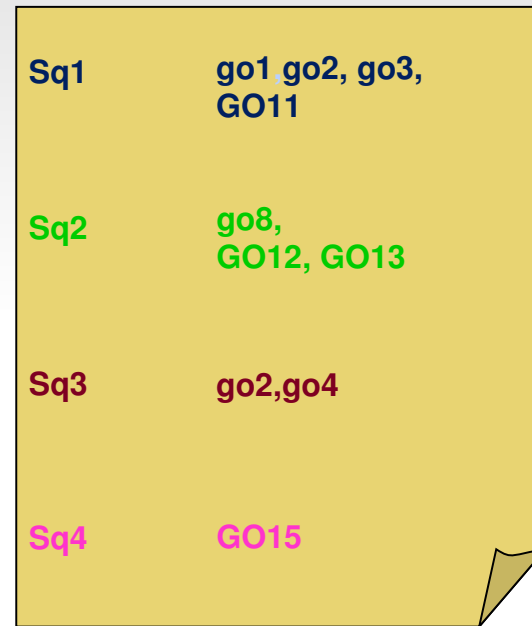
Blast2GO Annotation Strategy



Refinement



InterPro
Annex
GOSlim
Manual



Blast2GO annotation rule

Lowest term above threshold

Similarity requirement

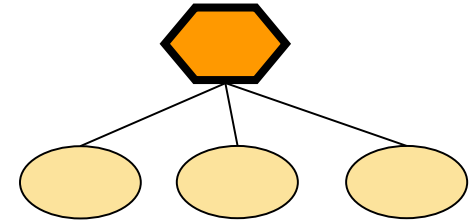
$$sim = \frac{\sum positives_{hsp}}{\sum alignmentlength_{hsp}}$$

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) + (#GO-1 x GOw) >= threshold]

Blast2GO annotation rule

Lowest.node [(max.sim x ECw) + (#GO-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