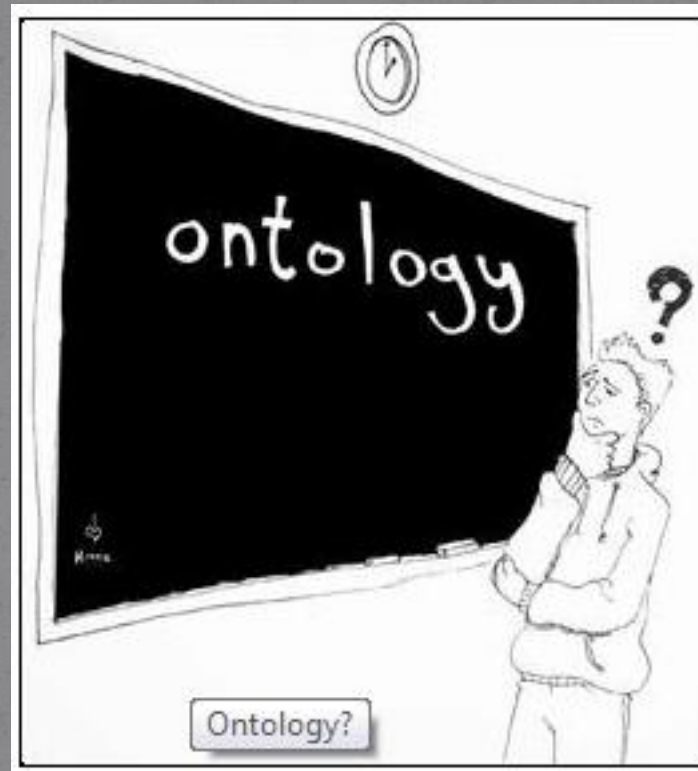


Introduction to ontologies and GO term analysis



- Pantelis Topalis – Bioinformatics Support Group

Metadata

- Data about data.
- Annotation is the process of linking metadata to data.
- Modern databases store and query metadata to facilitate more complicated queries.

Tons of data are available

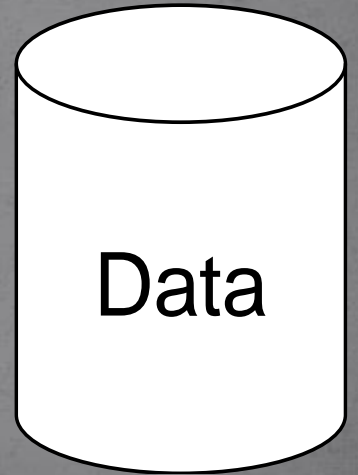
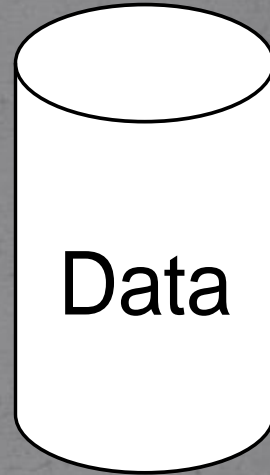
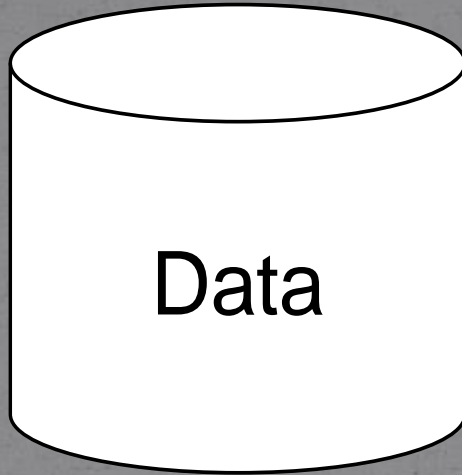
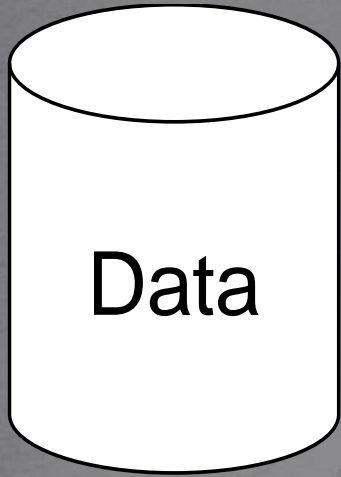
- High throughput methods yield large amount of data, often annotated in several different but related ways:
 - Annotations related to tissue and developmental stage where a gene is expressed.
 - Geographical distribution of alleles.
 - Resistance to drugs.
 - Ability to transmit a disease.

The modern tower of Babel (I)

- Currently data produced by a research project are stored in a homemade database specifically designed to serve the purpose of the project.
 - ✓ Innovative Vector Control Consortium - IVCC
 - ✓ The WHO/Gates Foundation Vector Biology and Control Project - VBC
 - ✓ African Network on Vector Research – ANVR
 - ✓ Malaria Atlas

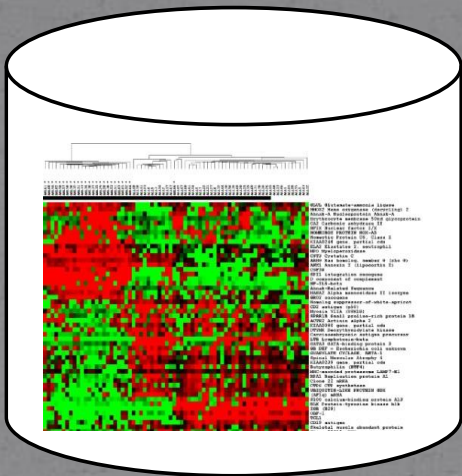
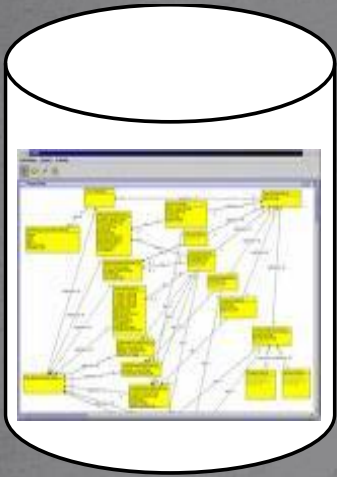
. The modern tower of Babel (II)

- Not a single format / standard to describe the same datatypes.
 - Genbank vs EMBL
 - Generic Feature Format GFF version 2 vs version 3
 - GVF (Genome Variation Format) vs VCF (Variant Call Format)
 - OBO vs OWL

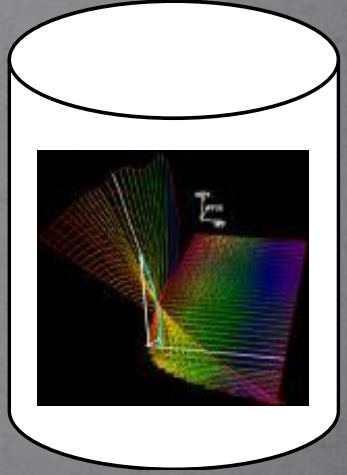
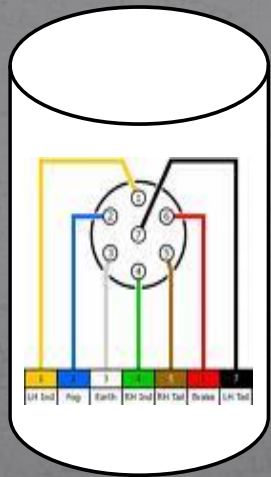
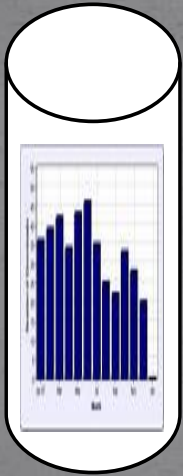
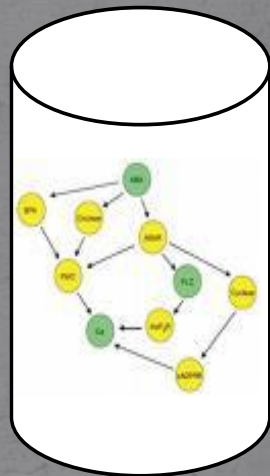


Biological data separated in “silos”

- Lab / pathological data.
- Data related to medical records.
 - Clinical data.
 - Patient history
 - Imaging data
- Gene expression datasets (microarrays, RNAseq).
- Proteomics data (mass spectrometry).
- Genomic Variation data.



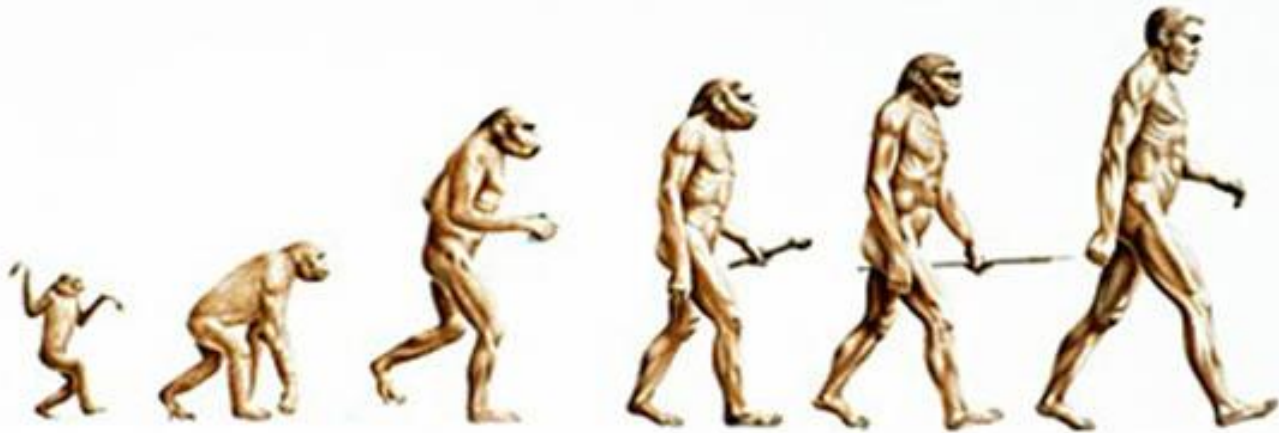
A vertical cylinder containing a column of binary code (0s and 1s) on a dark background. The code is arranged in a regular, grid-like pattern, with each character being a single bit.



- We need to **link** all these databases and to make them **interoperable**.
- Need to develop IT tools and decision support systems to take advantage and retrieve data from all those sources.

Ontology

- Branch of philosophy focusing to everything that exists.
- (IT) Ontology is a typical network of knowledge representation.
- In biomedical domain , ontologies are ways to represent and organize meta-data.



Λ
Ι
Σ
Τ
Α

Ο
Ρ
Ω
Ν

Λ
Ε
Ξ
Ι
Κ
Ο

Τ
Α
Ξ
Ι
Ν
Ο
Μ
Η
Σ
Η

Λ
Ε
Ξ
Ι
Λ
Ο
Γ
Ι
Ο

Θ
Η
Σ
Α
Υ
Ρ
Ο
Σ

Ο
Ν
Τ
Ο
Λ
Ο
Γ
ΙΑ

Basic elements of an ontology

- Terms unambiguously defined.
- Synonyms.
- Connections with the other terms of the ontology via logically defined relations.

Ontology Tree Editor

- ← I oostasis
- ⊖ ← I physiological process of malaria vector
 - ⊕ ← I behavioural process
 - ⊕ ← I chorio formation
 - ← I circulation
 - ⊕ ← I developmental process
 - ⊕ ← I distention of midgut
 - ⊕ ← I egg laying
 - ⊕ ← I endocrine system process
 - ⊕ ← I excretion
 - ← I fertilization (sensu Anophelinae)
 - ← I formation of ovarian follicles
 - ⊕ ← I formation of peritrophic matrix
 - ⊕ ← I growth
 - ⊕ ← I immune system process
 - ⊕ ← I muscular system process
 - ⊕ ← I nervous system process
 - ⊕ ← I nutritional process
 - ⊕ ← I previtellogenic development
 - ⊕ ← I regulation of biological process
 - ← I release of 20-hydroxyecdysone
 - ⊕ ← I reproduction
 - ⊕ ← I respiration
 - ⊕ ← I response to stimulus
 - ← I rRNA synthesis in oocyte and nurse cells
 - ← I saliva secretion
 - ← I secretion of peritrophic matrix in larvae
 - ⊕ ← I sensory perception

Text Editor

ID IDOMAL:0002033

Namespace malaria_ontology

Name formation of peritrophic matrix

Definition * Comment Cross Products

Definition

The formation of a membrane which is deposited around the ingested blood mass.

Dbxrefs

ISBN:0-412-40180-0

Xrefs Synonyms * Subsets

formation of peritrophic membrane

Scope: *Related Synonym*

Xrefs

● ISBN:0-12-473276

Reference and application ontologies

- Reference ontologies describe a scientific domain and provide the basis for the
- Application ontologies which aim to serve the needs of a specific scientific community.

Scientific vs operational ontologies

	Upper level ontologies	Domain specific ontologies
Scientific ontologies	BFO , Dolce , SUMO	GO, FMA, IDO
Operational ontologies	FOAF	Amazon.com, Library of Congress

The tower of Babel revisited

- There are many ways to create an ontology but the number of available ontologies is not solving the “data isolation” problem by itself.
- If every community or group creates an ontology using a different jargon interoperability is not promoted.
 - Drug database based on brand names only.
 - Metabolic resistance vs Detoxification process
- Scientific ontologies have to be linked together.

Ways to link ontologies together

- Use common reference ontologies while developing application ontologies for related fields.
 - Disease ontology (DO)
 - Infectious disease ontology (IDO)
 - Malaria ontology (IDOMAL)
 - Dengue Fever ontology (IDODEN)
- Reference ontologies contain generic terms which remain more or less the same, while application ontologies are more flexible to cover the needs as our knowledge is increased.
 - Neither DO or IDO have changed when a fifth serotype of Dengue virus had discovered and added to IDODEN.

Ways to link ontologies (II)

- Common use of relations in all ontologies related to a specific domain.
- OBO Foundry is an international consortium of ontology developers for the biomedical domain has issued a set of rules to facilitate ontology interoperability.

OBO Foundry rules

- Every term has to be related with another term via an is a relation. (is a completeness).
- Every term has only one is a parent.
- There is a predefined set of relations that has to be used in all biomedical ontologies.

Relations available for use in OBO Foundry

is_a	has_agent
part_of	instance_of
integral_part_of	realizes
proper_part_of	inheres_in
located_in	bearer_of
contained_in	has_quality
adjacent_to	has_function
transformation_of	has_role
derives_from	has_disposition
preceded_by	has_participant

Ontological realism

- Describe the scientific domain according to our best knowledge at the present time.
- Add terms and relations in order to describe what is true in reality and not some simplified computer-based model.
- An ontology is the product of collaboration between specialist from several different fields.

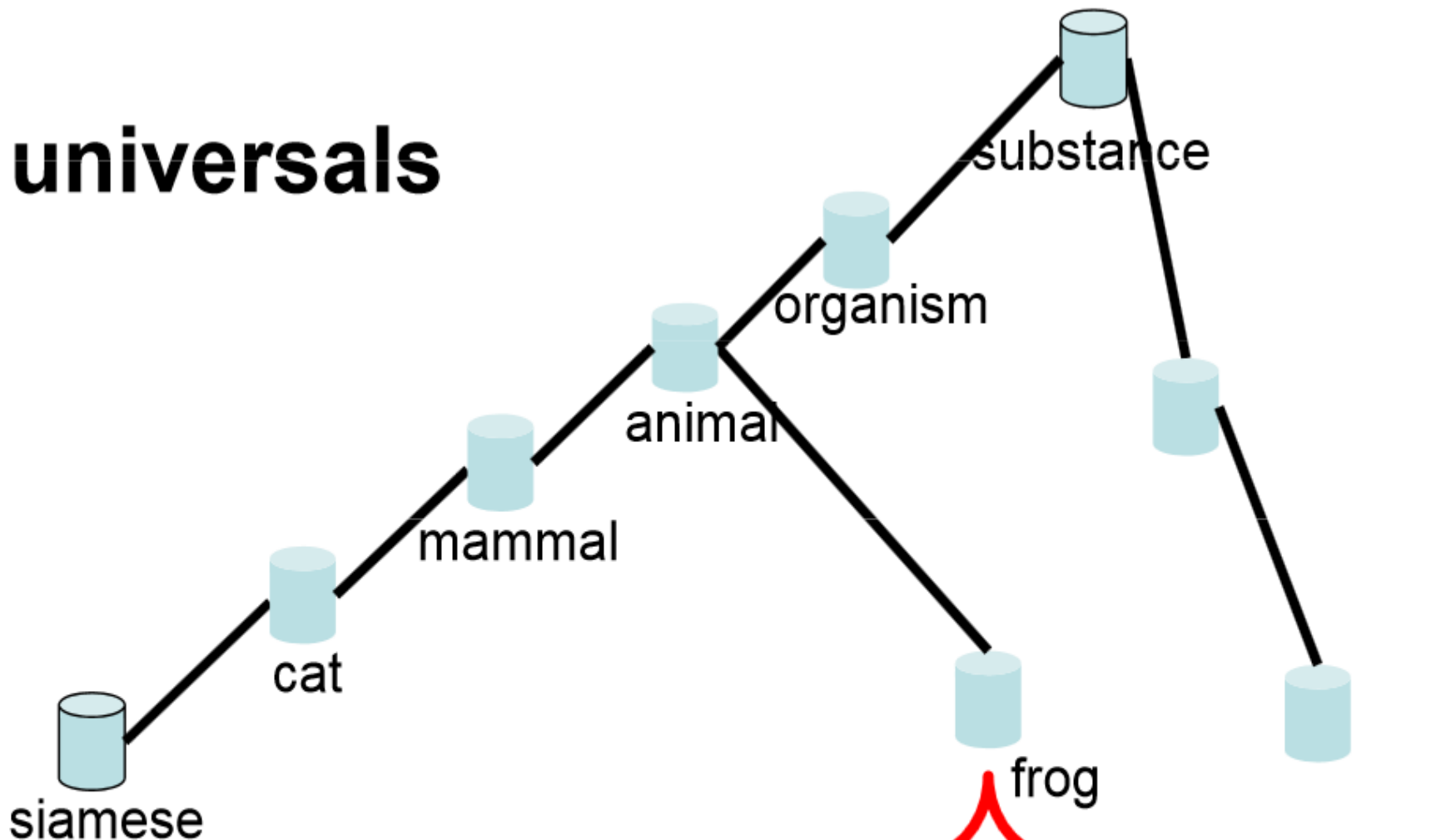
Basic elements of an ontology

- Terms unambiguously defined.
- Synonyms.
- Connections with the other terms of the ontology via logically defined relations.

Some definitions

- Entity: Everything that exists: An object, a process, a function, space, text, computer software.
- Universal entities versus instances.
- Human beings versus specific individuals.
- An ontology focuses on universals where as a database focuses on instances.

universals



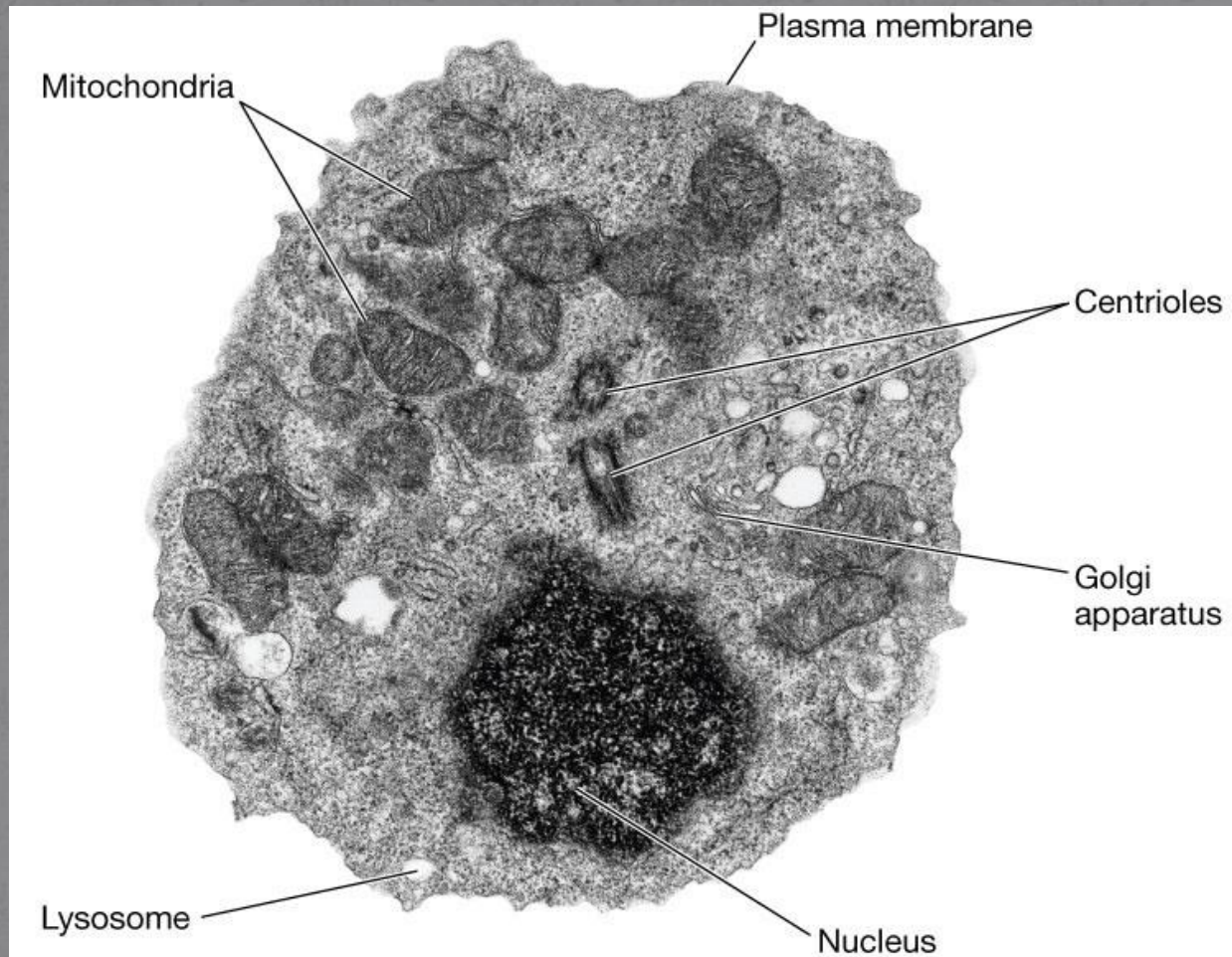
instances



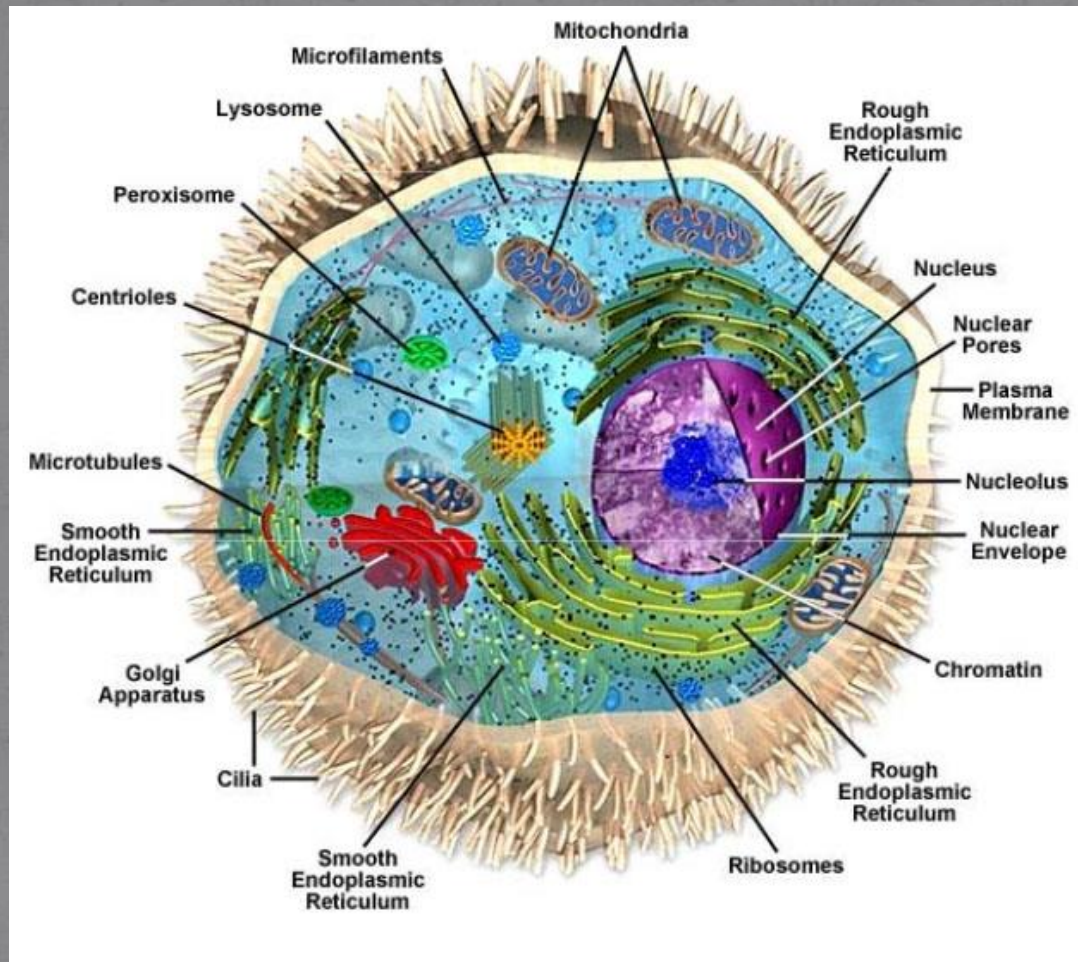
In a scientific ontology

- Every term represents both universal entities and their instances.
- Every term represents only one universal entity.
- Therefore we can define ontology as a representation of universals.

A photograph is a representation of an instance



We can't photograph a universal but we can create it's representation



Periodic Table of the Elements

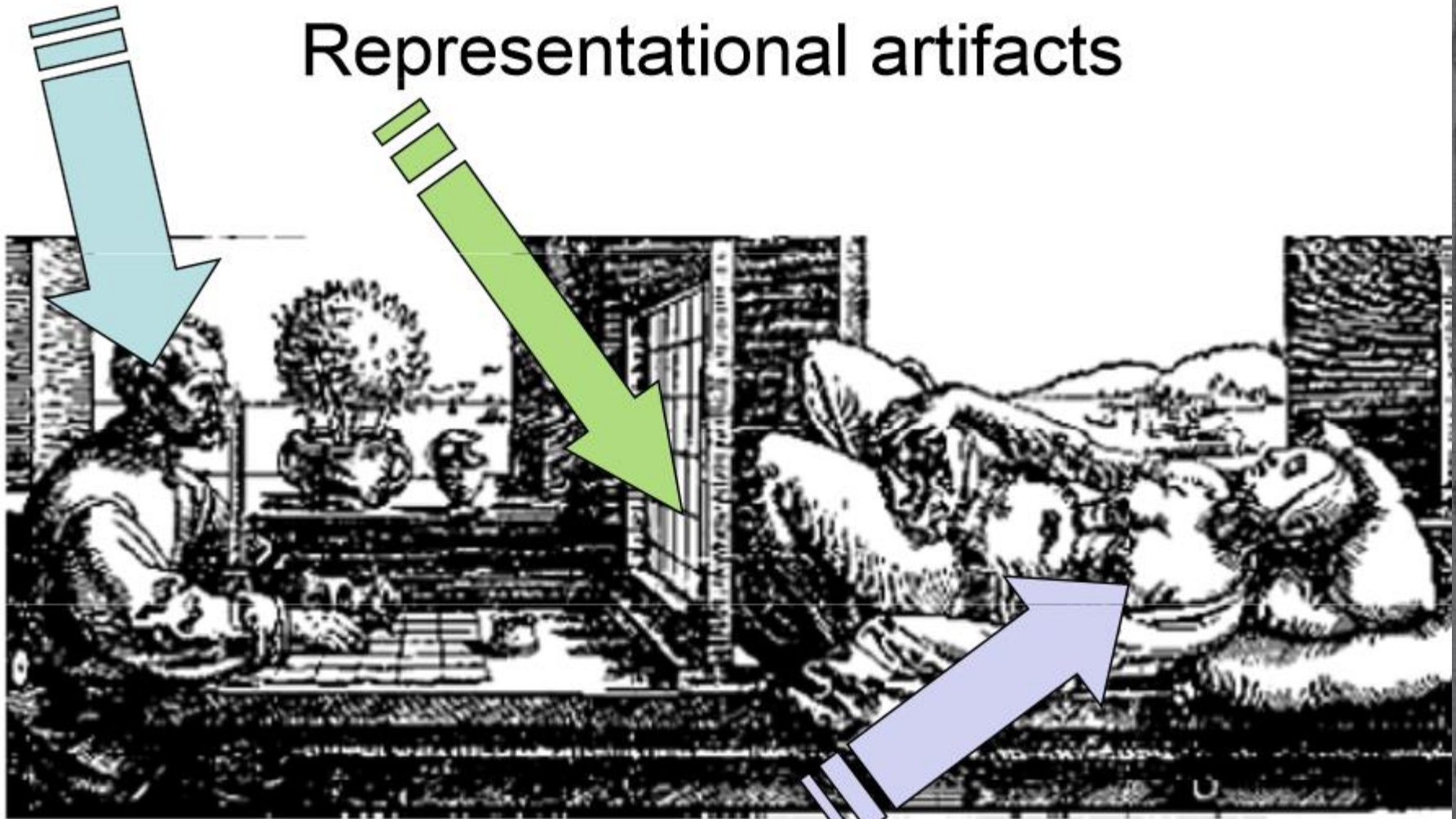
- hydrogen
- alkali metals
- alkali earth metals
- transition metals
- poor metals
- nonmetals
- noble gases
- rare earth metals

1 H																	2 He
3 Li	4 Be											5 B	6 C	7 N	8 O	9 F	10 Ne
11 Na	12 Mg											13 Al	14 Si	15 P	16 S	17 Cl	18 Ar
19 K	20 Ca	21 Sc	22 Ti	23 V	24 Cr	25 Mn	26 Fe	27 Co	28 Ni	29 Cu	30 Zn	31 Ga	32 Ge	33 As	34 Se	35 Br	36 Kr
37 Rb	38 Sr	39 Y	40 Zr	41 Nb	42 Mo	43 Tc	44 Ru	45 Rh	46 Pd	47 Ag	48 Cd	49 In	50 Sn	51 Sb	52 Te	53 I	54 Xe
55 Cs	56 Ba	57 La	72 Hf	73 Ta	74 W	75 Re	76 Os	77 Ir	78 Pt	79 Au	80 Hg	81 Tl	82 Pb	83 Bi	84 Po	85 At	86 Rn
87 Fr	88 Ra	89 Ac	104 Unq	105 Unp	106 Unh	107 Uns	108 Uno	109 Une	110 Unn								

58 Ce	59 Pr	60 Nd	61 Pm	62 Sm	63 Eu	64 Gd	65 Tb	66 Dy	67 Ho	68 Er	69 Tm	70 Yb	71 Lu
90 Th	91 Pa	92 U	93 Np	94 Pu	95 Am	96 Cm	97 Bk	98 Cf	99 Es	100 Fm	101 Md	102 No	103 Lr

Cognitive representations

Representational artifacts



Reality







Linguistic problems

- Sometimes the same word/term is used for a universal entity and a subset of instances generating confusion.
- Ebola virus can cause death.
- Ebola virus is transmitted uncontrollably in west Africa.

Scientific ontology

- An artifact to represent universals and the relations that relate them

Gene Ontology (GO)

- A structured representation of the qualities / properties of genes available to all interested to the universal biological reality.

GO is not an ontology, but 3 independent ontologies together

**cellular
component**

**molecular
function**

**biological
process**

No connections between the 3 main components of GO

cellular
component

molecular
function

biological
process

Continuant

Independent
Continuant

cell component

Dependent
Continuant

molecular function

Occurrent

biological process

**cellular
component**

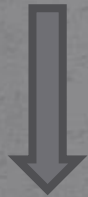
**biological
process**

**molecular
function**

**cellular
biological
process**

**organism-level
biological
process**

**molecular
function**



molecule

**cellular
biological
process**



**cellular
component**

**organism-level
biological
process**



organism

**molecular
process**

**cellular
process**

**organism-
level
biological
process**

**molecular
function**

**cellular
function**

**organism-
level
biological
function**

molecule

**cellular
component**

organism

Continuant

Independent
Continuant

cell component

Dependent
Continuant

molecular function

Occurrent

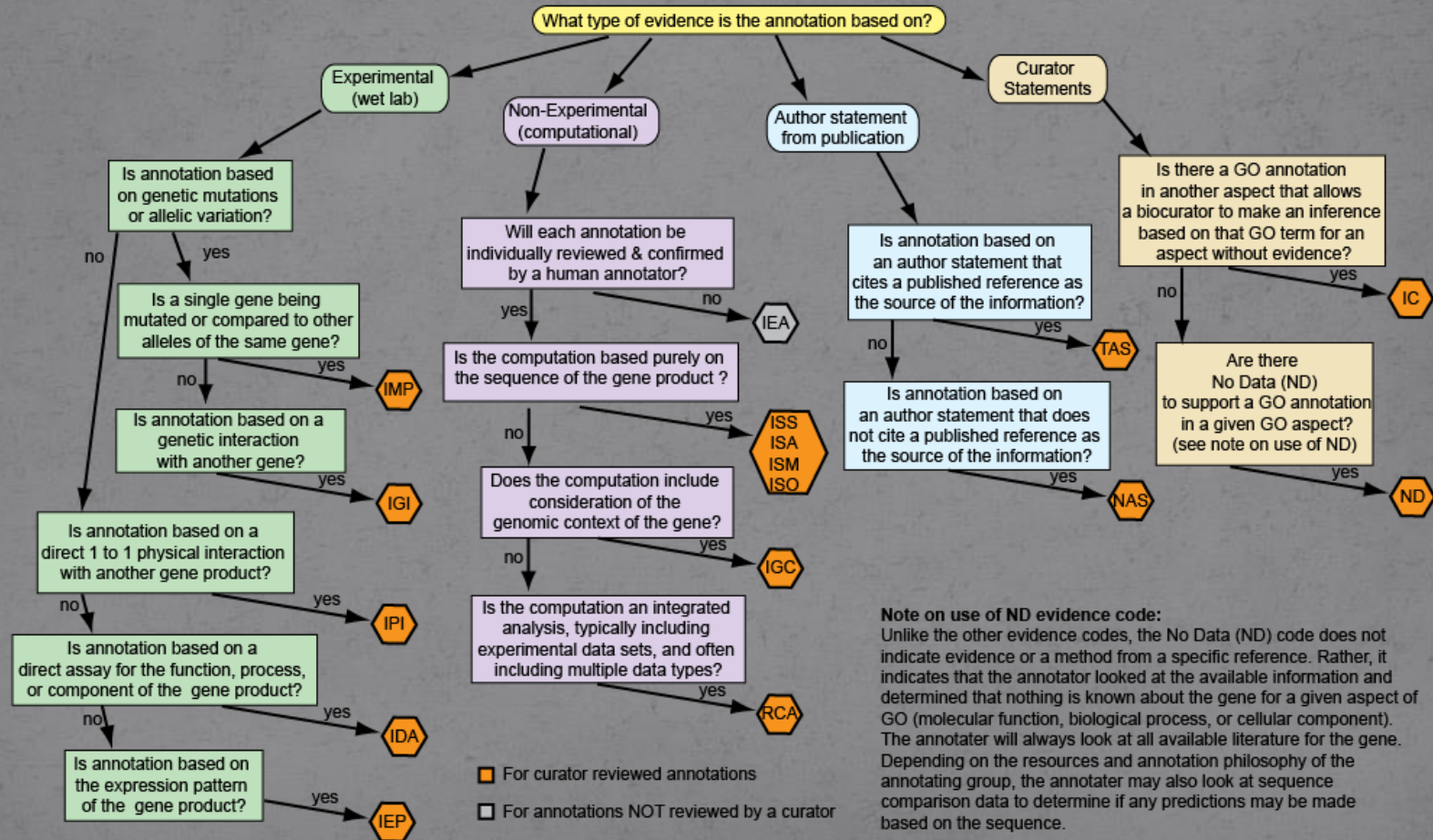
biological process

GO curation

- GO curators develop the ontology and also provide annotations on behalf of the model organism databases.

Evidence codes

GO Evidence Code Decision Tree



GO usage

- The Gene Ontology (GO) provides core biological knowledge representation for modern biologists, whether computationally or experimentally based.
- GO annotations provide largely species-neutral, comprehensive statements about what gene products do.
- The general user of GO resources often misses fundamental distinctions about GO structures, GO annotations, and what can and can not be extrapolated from GO resources.

Know the Source of the GO Annotations You Use

- The GO site has the most comprehensive and current sets of annotations.
- Annotations are also contained within external tools and applications, although these may not be updated often.

Understand the Scope of GO Annotations

- The GO annotation stream focuses on the capture of the knowledge about the functional activities of specific proteins, the larger biological process (such as photosynthesis) as part of which these specific functions collectively act, and the cellular locality where all this occurs.

Consider Differences in Evidence Codes

- There are 21 evidence codes currently in use to document how data are obtained, and GO may adopt an evidence code classification in the near future. Some evidence codes indicate different classes of experiments such as “inferred by direct assay.” Some indicate different approaches to prediction from comparative analysis such as “inferred from sequence orthology.”

Probe Completeness of GO Annotations

- While the GO annotation corpus generally has excellent broad coverage of available knowledge about gene products via structurally based annotations from sources such as InterPro, the completeness of annotations derived from biomedical literature is uneven because there are not enough GO curators to keep them current on all fronts.
- The lack of a GO annotation does not mean that a gene product does not perform a particular function or act in a particular role.

Understand the Complexity of the GO Structure

- The GO structure, relations, and terminology are modified every day by GO ontology editors. New relations between terms are added as the GO refines the representation of biological knowledge.
- There is a need to know the version of the ontology which is being used by our analysis tool.

Choose Analysis Tools Carefully

- Hundreds of GO-focused applications and tools are available.
- Term enrichment analysis, a common use of GO resources, is incorporated into many different applications and analysis tool sites. Different implementations/algorithms may give different results.

AmiGO 2

More information on quick search [?](#)

Get Started with Grebe



Use the Grebe Search Wizard to **get started** in exploring the Gene Ontology data.

Advanced Search



Interactively **search** the Gene Ontology data for annotations, gene products, and terms using a powerful search syntax and filters.

GOOSE



Use GOOSE to query a legacy GO database with **SQL** or edit one of the templates.

Term Enrichment Service



Powered by [PANTHER](#)

Statistics



View the most recent **statistics** about the Gene Ontology data on the main site.

And Much More...



Many **more tools** are available from the software list, such as alternate searching modes, Visualize, non-JavaScript pages.

<http://www.amigo.org>

Term annotations

What are the direct and indirect annotations to term ? [Go »](#)

What are the direct and indirect annotations for organism (scientific name) to term ? [Go »](#)

Gene product annotations

What are all the annotations for gene product [Go »](#)

Gene products associated with terms

What are the gene products annotated to term but *not* term ? [Go »](#)

What are the gene products annotated to term and term ? [Go »](#)

Protein family

What are the annotations associated with the protein family ? [Go »](#)

AmiGO 2

More information on quick search [?](#)

Get Started with Grebe



Use the Grebe Search Wizard to **get started** in exploring the Gene Ontology data.

Advanced Search



Interactively **search** the Gene Ontology data for annotations, gene products, and terms using a powerful search syntax and filters.

GOOSE



Use GOOSE to query a legacy GO database with **SQL** or edit one of the templates.

Term Enrichment Service



Powered by [PANTHER](#)

Statistics



View the most recent **statistics** about the Gene Ontology data on the main site.

And Much More...



Many **more tools** are available from the software list, such as alternate searching modes, Visualize, non-JavaScript pages.

<http://www.amigo.org>

Gene IDs

Gene IDs...

Species

H. sapiens ▼

Ontology

biological process ▼

Correction

Use Bonferroni correction

Resource

PANTHER ▼

Results viewer

AmiGO ▼

Submit

Species

H. sapiens

H. sapiens

M. musculus

R. norvegicus

G. gallus

D. rerio

D. melanogaster

C. elegans

D. discoideum

S. pombe

S. cerevisiae

A. thaliana

E. coli

C. albicans

Ontology

biological process

biological process

molecular function

cellular component

biological process (experimental only)

molecular function (experimental only)

cellular component (experimental only)

Gene product annotations

What are all the annotations for gene product

Go »

Gene products associated with

What are the gene products annotated to term

What are the gene products annotated to term

Protein family

What are the annotations associated with the

Notch4 (MGI:MGI:107471/Mus musculus)
Notch2 (MGI:MGI:97364/Mus musculus)
Notch1 (MGI:MGI:97363/Mus musculus)
NOTCH1 (UniProtKB:F1MSM3/Bos taurus)
NOTCH2 (UniProtKB:F1MPE9/Bos taurus)
NOTCH4 (UniProtKB:G3X812/Bos taurus)
NOTCH3 (UniProtKB:E1BPT8/Bos taurus)
Human Notch 3 (UniProtKB:Q14962/Homo sapiens)
NOTCH2 (UniProtKB:F1P236/Gallus gallus)
NOTCH1 (UniProtKB:F1NZ70/Gallus gallus)



Free-text filtering

Your search is pinned to these filters

+ document_category: annotation

User filters

+ bioentity: UniProtKB:Q14962

- Source
- Assigned by
- Ontology (aspect)
- Evidence type
- PANTHER family
- Qualifier
- Taxon
- Direct annotation
- Inferred annotation
 - anatomical structure morphogenesis (1)
 - biological_process (1)
 - cell part (1)
 - cellular_component (1)
 - developmental process (1)
 - integral component of membrane (1)
 - integral component of plasma membrane (1)
 - membrane part (1)
 - plasma membrane part (1)
- Annotation extension

Found entities

Total: 2; showing 1-2 Results count

<input type="checkbox"/>	Gene/product	Gene/product name	Qualifier	Direct annotation	Annotation extension	Source	Taxon	Evidence	Evidence with	PANTHER family	Isoform	Reference
<input type="checkbox"/>	Human Notch 3	Notch 3 protein		anatomical structure morphogenesis		UniProtKB	Homo sapiens	TAS				PMID:7698746
<input type="checkbox"/>	Human Notch 3	Notch 3 protein		integral component of plasma membrane		UniProtKB	Homo sapiens	TAS				PMID:7835890

<http://geneontology.org/page/guide-go-evidence-codes>

Thank you!