

Computing on the shoulders of giants:

how existing knowledge is represented
and applied in bioinformatics

Benjamin Good

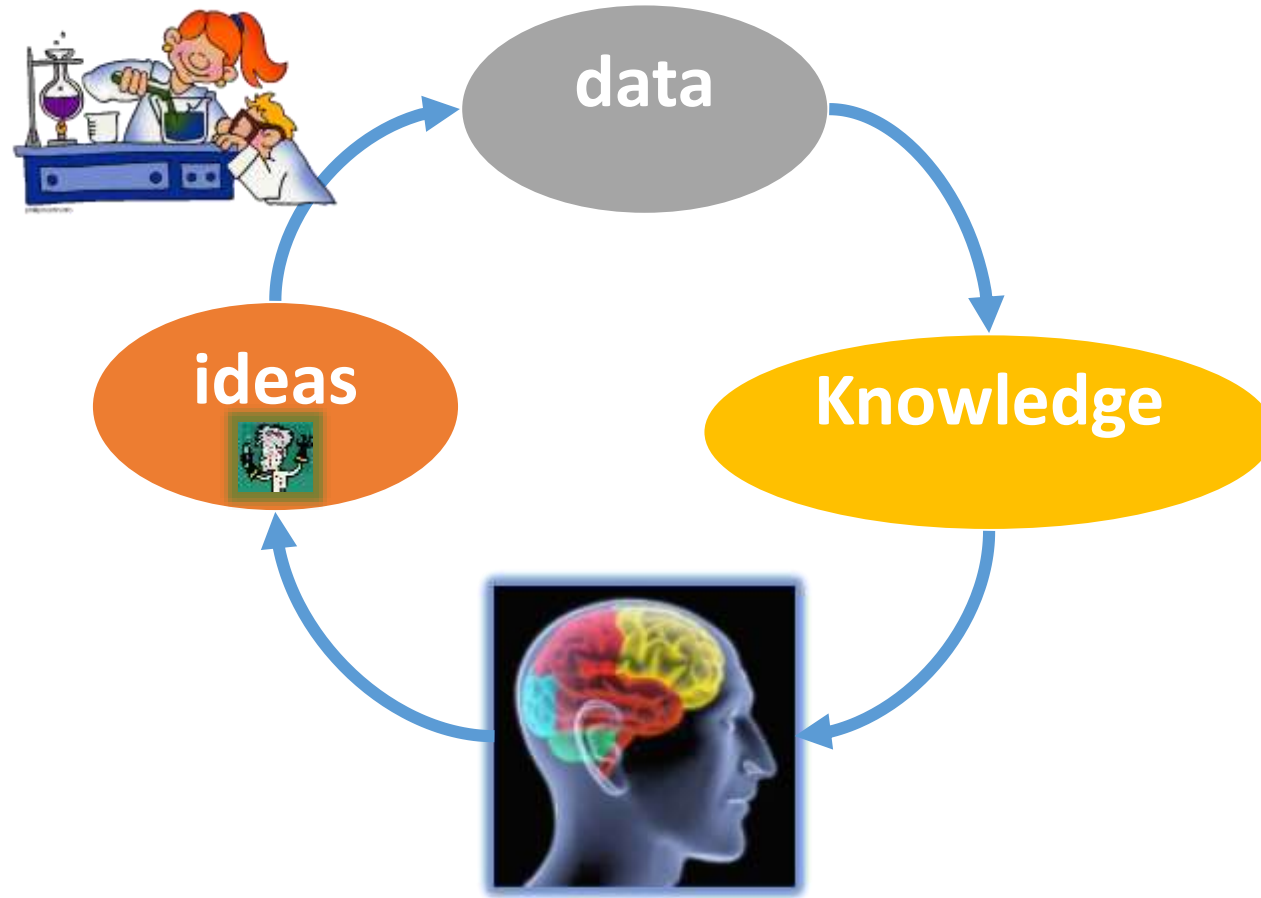
bgood@scripps.edu

Assistant Professor of the Department of
Molecular and Experimental Medicine

Specialty: artificial intelligence, crowdsourcing



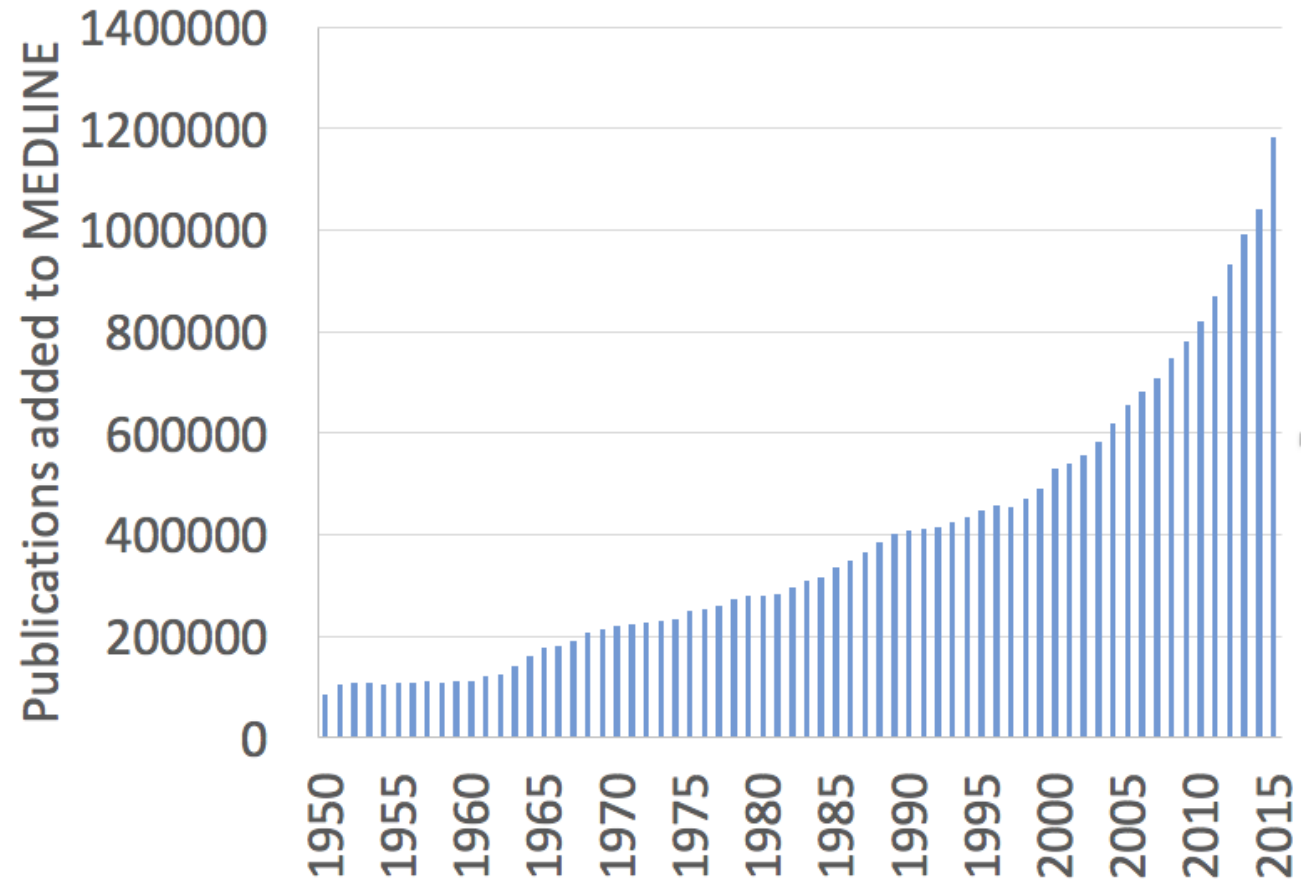
The more you can 'know' the better a scientist you can become



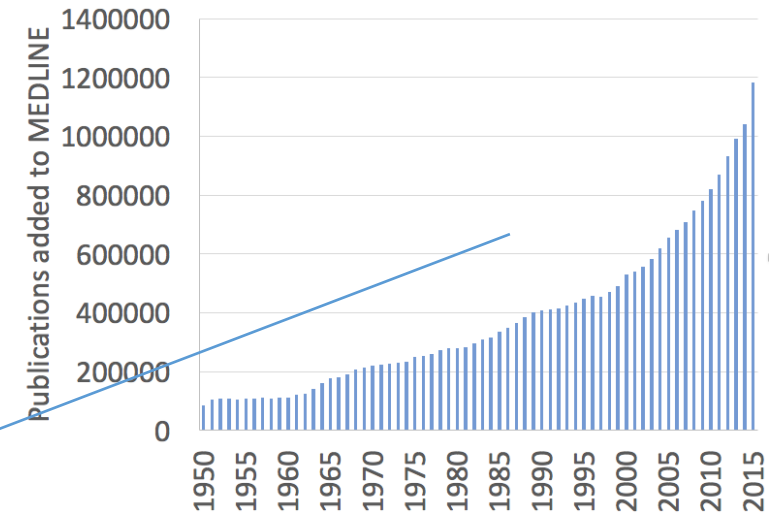
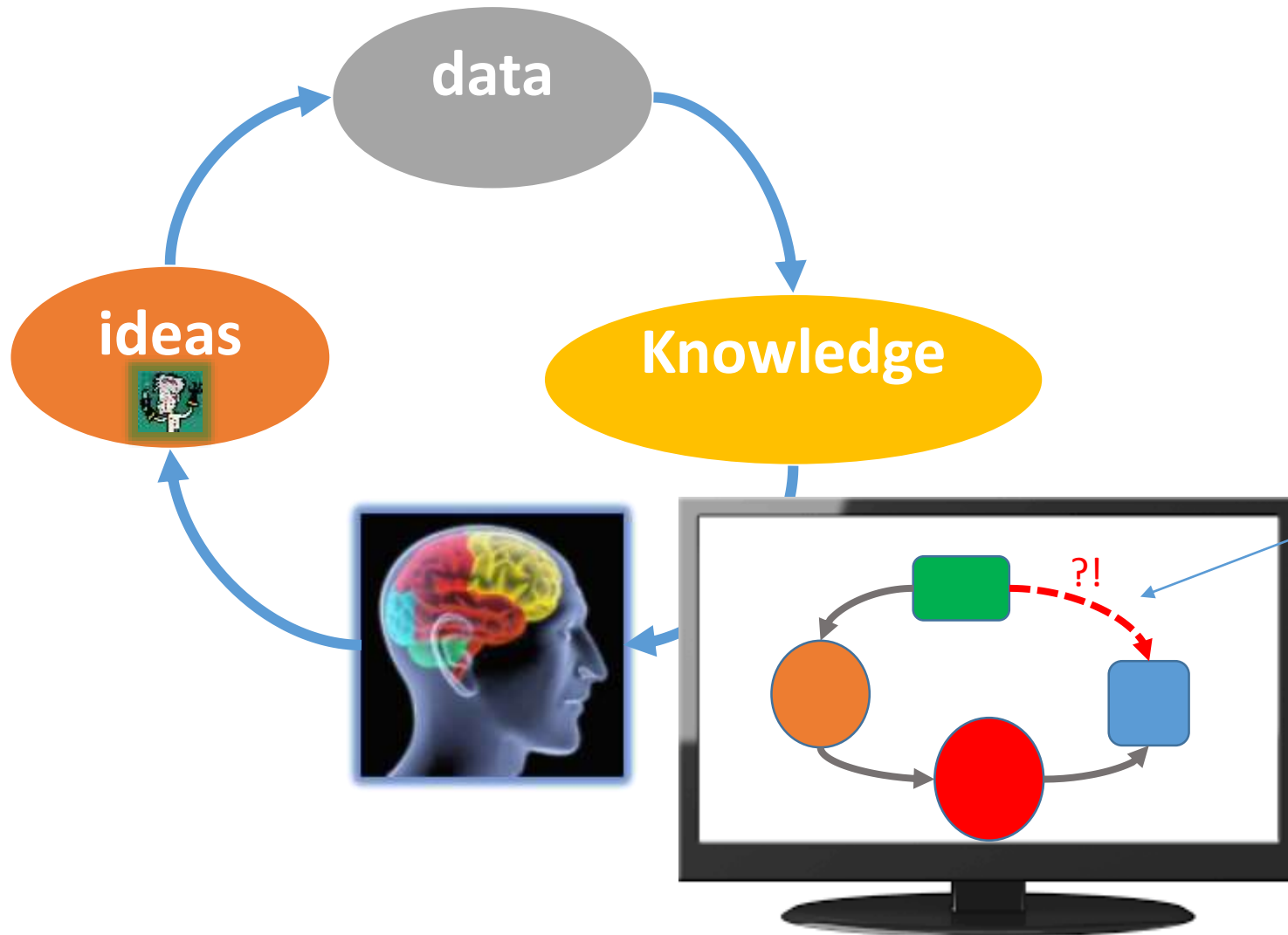
Too much to know

Knowledge

- PubMed lists > 1 million articles published each year (more than 2 per minute)
- **Your capacity to read and comprehend is limiting**



Knowledge representation



Goals for representing knowledge (outline)

- Make things (articles, genes, antibodies, etc.) easier to find

Controlled vocabularies (MeSH)
Ontologies (Gene Ontology)

- Answer questions

- Generate hypotheses

knowledge graphs on the Web:
the SPARQL query language

knowledge plus computation =
inference, the ABC model

Part 1: Medical Subject Headings (MeSH)

Finding what to read: controlled vocabularies for indexing PubMed

- What happens when you search PubMed?

The screenshot shows the PubMed website interface. At the top, there's a navigation bar with 'NCBI', 'Resources', and 'How To' links, along with a 'Sign in to NCBI' button. Below this is the 'PubMed.gov' logo and the text 'US National Library of Medicine National Institutes of Health'. A search bar contains the text 'fainting', with a dropdown menu showing 'PubMed'. To the right of the search bar are links for 'Create RSS', 'Create alert', and 'Advanced', and a 'Search' button. Below the search bar, there's a 'Help' link. On the left side, there's a sidebar with 'Article types' (Clinical Trial, Review, Customize ...), 'Text availability' (Abstract, Free full text, Full text), and 'PubMed Commons' (Reader comments, Trending articles). The main content area shows 'Search results' for 'fainting'. It includes a summary dropdown, '20 per page', and 'Sort by Most Recent'. Below this, it says 'Items: 1 to 20 of 18088'. A red arrow points from the search bar to the 'Search results' heading. The first result is 'Emergent pacemaker placement in a patient with Lyme carditis-induced complete heart block and ventricular asystole.' by Brownstein AJ, Gautam S, Bhatt P, Nanna M. The result includes the journal 'BMJ Case Rep. 2016 May 20;2016. pii: bcr2016214474. doi: 10.1136/bcr-2016-214474.', the PMID '27207985', and a link to 'Similar articles'. On the right side, there's a 'Results by year' bar chart and a 'PMc Images search for fainting' section.

NCBI Resources How To Sign in to NCBI

PubMed.gov US National Library of Medicine National Institutes of Health

PubMed fainting Search

Create RSS Create alert Advanced Help

Article types
Clinical Trial
Review
Customize ...

Text availability
Abstract
Free full text
Full text

PubMed Commons
Reader comments
Trending articles

Summary 20 per page Sort by Most Recent Send to: Filters: Manage Filters

Search results

Items: 1 to 20 of 18088 << First < Prev Page 1 of 905 Next > Last >>

1. [Emergent pacemaker placement in a patient with Lyme carditis-induced complete heart block and ventricular asystole.](#)
Brownstein AJ, Gautam S, Bhatt P, Nanna M.
BMJ Case Rep. 2016 May 20;2016. pii: bcr2016214474. doi: 10.1136/bcr-2016-214474.
PMID: 27207985
[Similar articles](#)

Results by year

Download CSV

PMc Images search for fainting

Summary ▾ 20 per page ▾ Sort by Most Recent ▾

Send to: ▾

Filters: [Manage Filters](#)

Search results

Items: 1 to 20 of 18088

<< First < Prev Page 1 of 905 Next > Last >>

Results by year ▾

PMC Images search for *fainting* ▾

Titles with your search terms ▾

Find related data ▾

Search details ▴

"syncope"[MeSH Terms] OR "syncope"
[All Fields] OR "fainting"[All
Fields]

Search

[See more...](#)

Summary ▾ 20 per page ▾ Sort by Most Recent ▾

Search results

Items: 1 to 20 of 976

[induced complete heart block and](#)

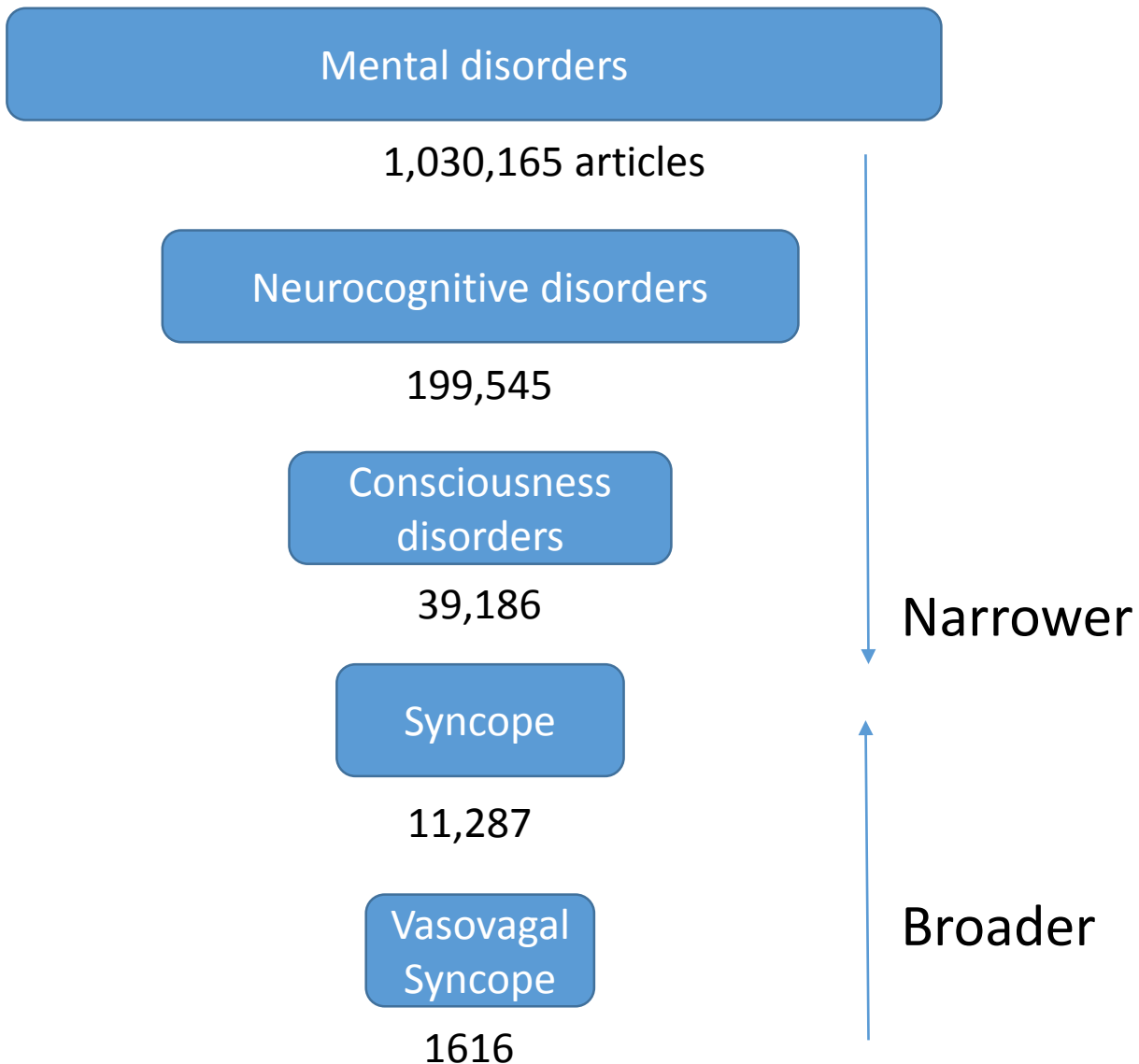
2016-214474.

[ing in patients with syncope\].](#)

ia A, Broullón FJ, Álvarez-García N,

[.medcli.2016.03.041. [Epub ahead of

MeSH controlled vocabulary (AKA 'thesaurus')



- **Descriptor Unique ID:** D013575
- **Definition:** A transient loss of consciousness and postural tone caused by diminished blood flow to the brain...
- **Entry Terms:** Syncope, **Fainting**, Syncopal Vertigo, Presyncope, Drop Attack, Carotid Sinus Syncope,...
- **Relations to other terms**

MeSH: medical subject headings

- >27,000 descriptors
- >87,000 entry terms
- 16 hierarchical trees
- Constantly being revised

1. + Anatomy [A]
2. - Organisms [B]
 - [Eukaryota \[B01\]_+](#)
 - [Archaea \[B02\]_+](#)
 - [Bacteria \[B03\]_+](#)
 - [Viruses \[B04\]_+](#)
 - [Organism Forms \[B05\]_+](#)
3. + Diseases [C]
4. + Chemicals and Drugs [D]
5. + Analytical, Diagnostic and Therapeutic Techniques and Equipment [E]
6. + Psychiatry and Psychology [F]
7. + Phenomena and Processes [G]
8. + Disciplines and Occupations [H]
9. + Anthropology, Education, Sociology and Social Phenomena [I]
10. + Technology, Industry, Agriculture [J]
11. + Humanities [K]
12. + Information Science [L]
13. + Named Groups [M]
14. + Health Care [N]
15. + Publication Characteristics [V]
16. + Geographicals [Z]



Demo and play time

- View and explore the MeSH trees:
 - https://www.nlm.nih.gov/mesh/2016/mesh_browser/MeSHtree.html
- Use MeSH to query PubMed
 - Go to: <http://www.ncbi.nlm.nih.gov/mesh>
 - Search for the term 'fainting'
 - click 'Add to search builder'
 - click search PubMed
 - click back, search for other things..



Query demos

- **Query expansion**

- [Hand Bones \[Mesh\]](#)
- [Hand Bones \[Mesh:NoExp\]](#)

- **Boolean operators**

- cardiac hypertrophy and use rodents besides mice and rats in their experiments
 - [\("Cardiomegaly"\[Mesh\]\)](#)
 - [AND "Rodentia"\[Mesh\]](#)
 - [NOT "Mice"\[Mesh\] NOT "Rats"\[Mesh\]](#)

- **Article type filter**

- Review papers about cardiac hypertrophy
- [Cardiomegaly \[MeSH\] AND Review\[ptyp\]](#)
- Try with <http://www.ncbi.nlm.nih.gov/pubmed/advanced>

Questions about MeSH ?

- Good 3 minute tutorial video on practical use:
<http://www.youtube.com/watch?v=uyF8uQY9wys>

Part 2: Ontology

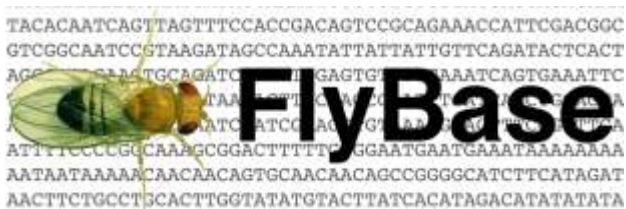
“Ontology”

- The word comes from philosophy:
 - “the branch of metaphysics dealing with the nature of being”
- In practice they are:
 - **A set of concepts, definitions and inter-relationships.**
 - (The dividing line between “controlled vocabulary”, “thesaurus”, “ontology” is hazy and not terribly important for practical purposes.)
- We have hundreds of ontologies in biology, e.g. see:
 - <http://www.obofoundry.org> (100+)
 - <http://biportal.bioontology.org> (500+)

The Gene Ontology

Started in 1999

As a collaboration between 3 Model Organism Databases



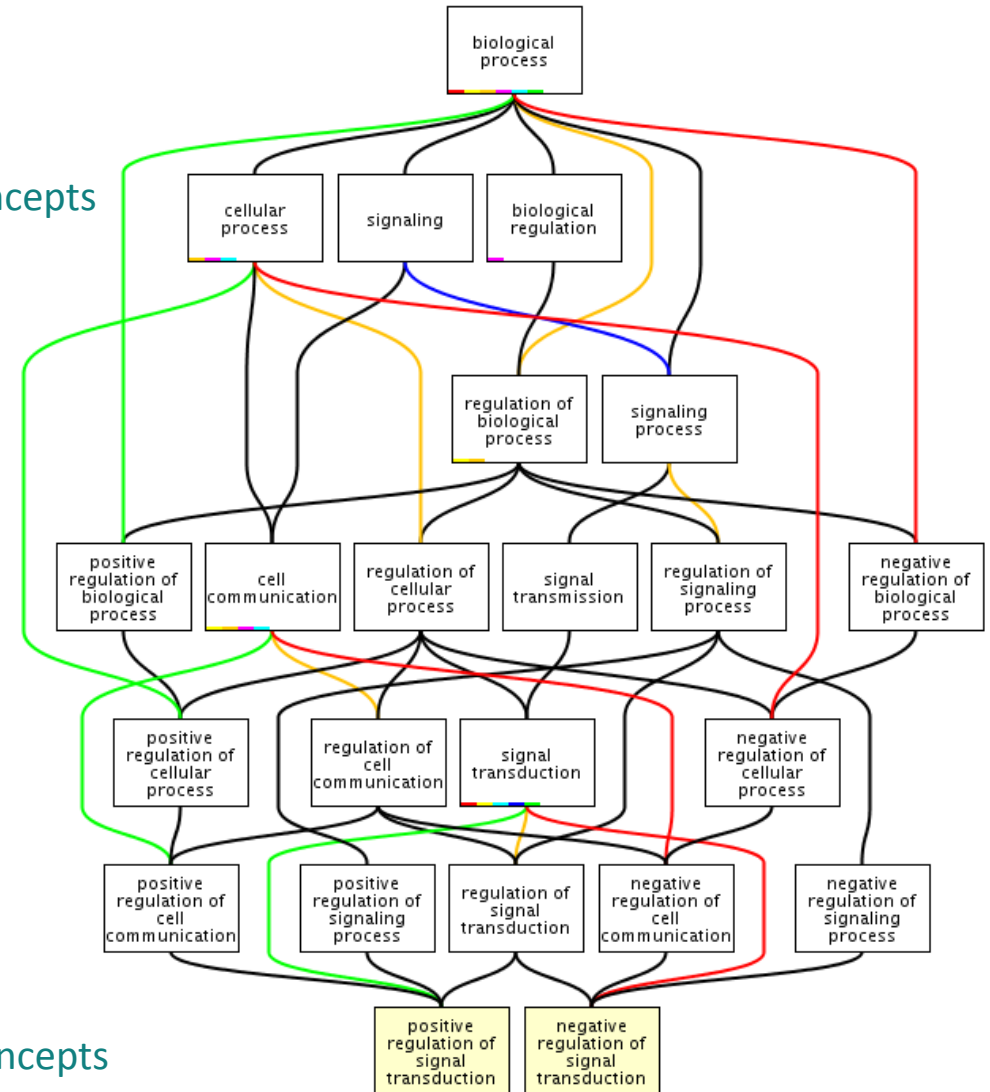
Saccharomyces
GENOME DATABASE



The Gene Ontology

- A way to capture biological knowledge for individual gene products in a computable form
- A set of concepts and their relationships to each other arranged as a hierarchy

Less specific concepts

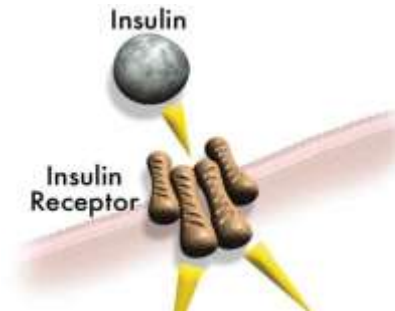


<http://www.ebi.ac.uk/QuickGO>

The GO branches

1. Molecular Function

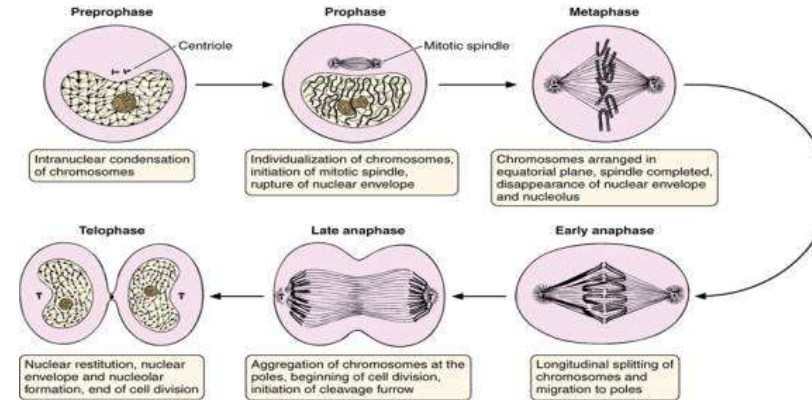
An elemental activity or task or job



- protein kinase activity
- insulin receptor activity

2. Biological Process

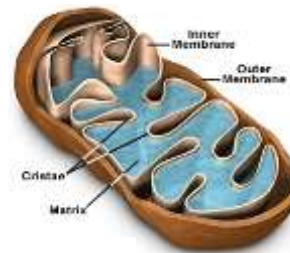
A commonly recognized series of events



- cell division

3. Cellular Component

Where a gene product is located



- mitochondrion
- mitochondrial matrix
- mitochondrial inner membrane

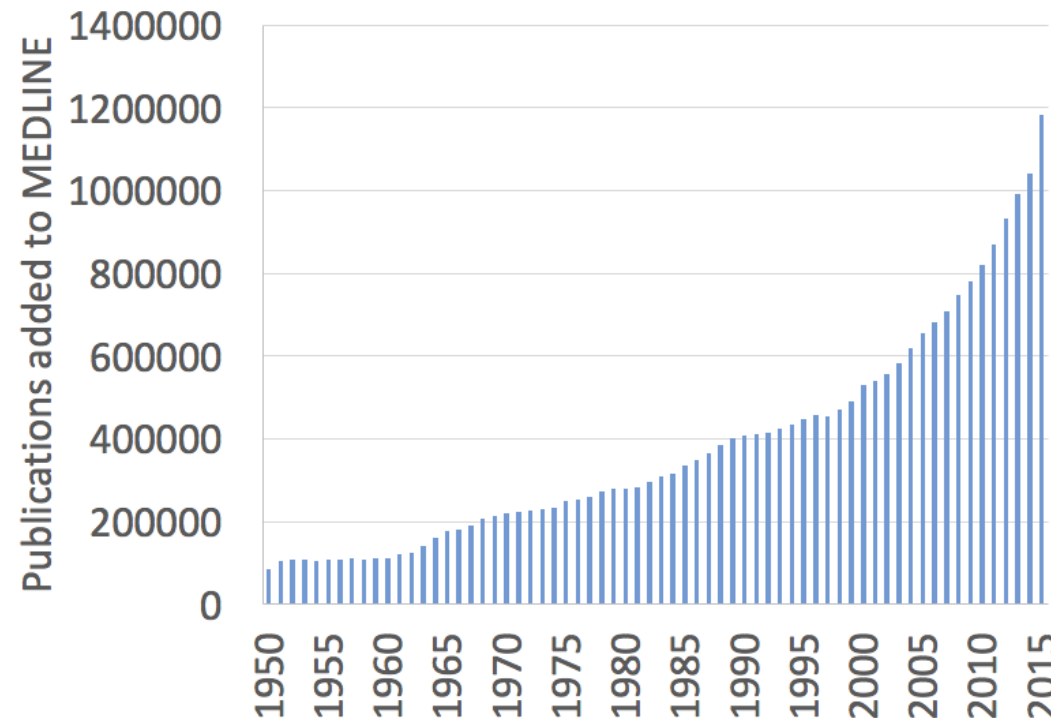
Building the GO (now covering more than 40,000 terms)

- GO editorial team based at the European Bioinformatics Institute

- Submission via <http://www.ebi.ac.uk/ontology/>

- Submissions via <http://www.ebi.ac.uk/ontology/>

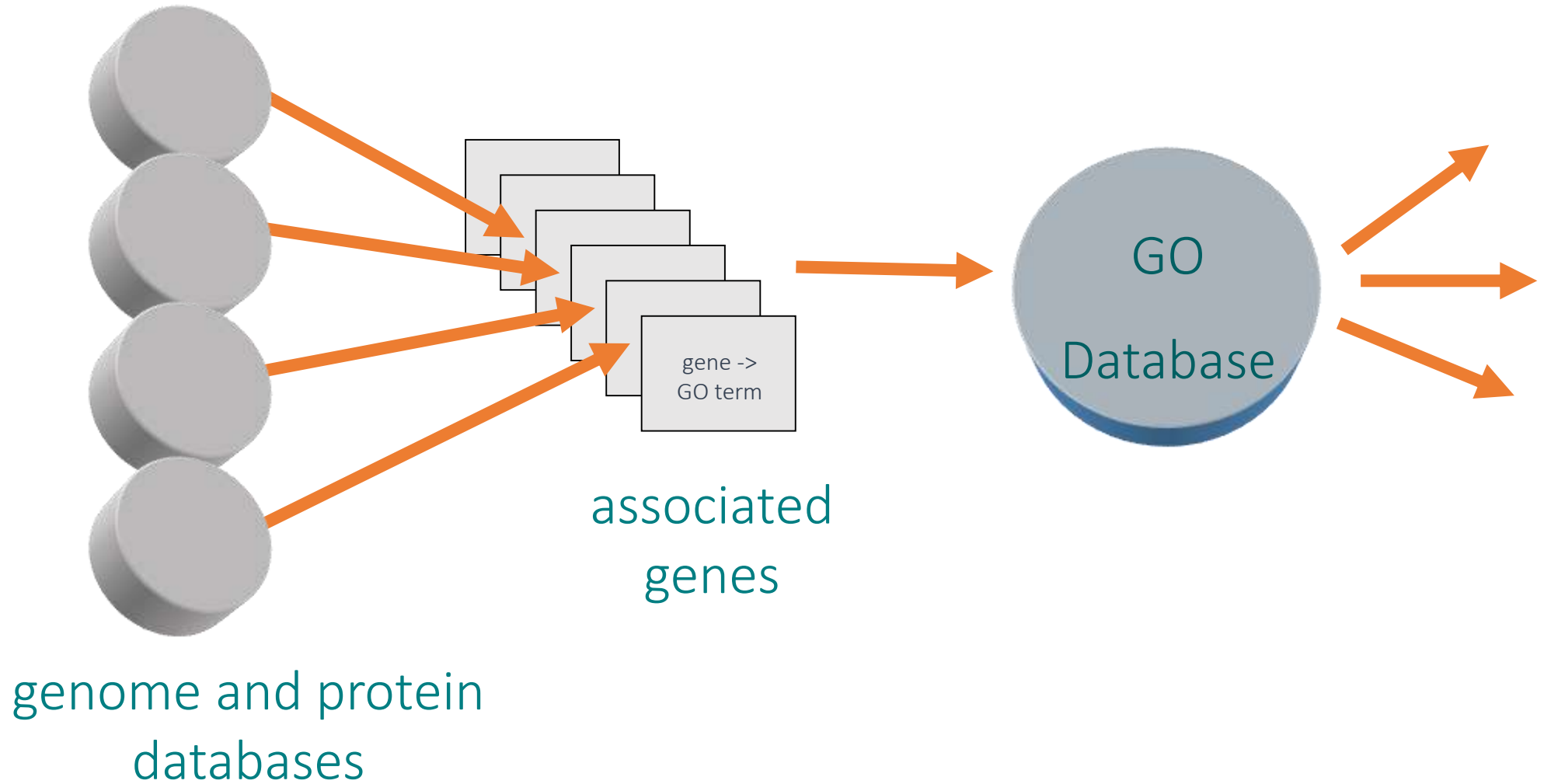
- In principal, an editors make t



<http://www.ebi.ac.uk/ontology/>

ontology, but the GO

Using the GO to describe gene products



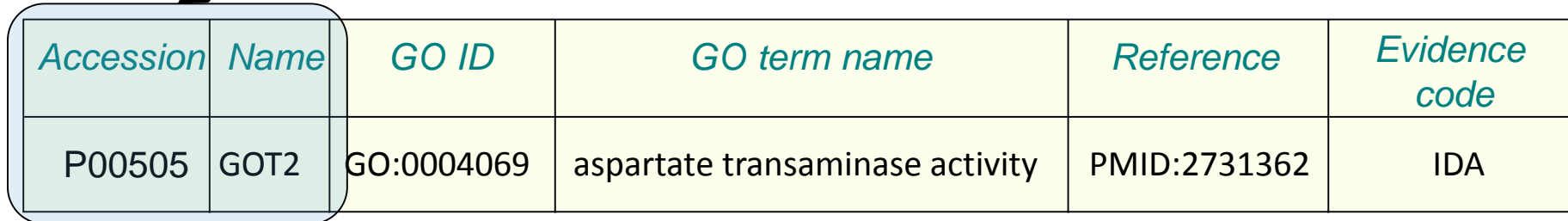
Contributors

PomBase



A GO annotation is ...

...a statement that a gene product;




<i>Accession</i>	<i>Name</i>	<i>GO ID</i>	<i>GO term name</i>	<i>Reference</i>	<i>Evidence code</i>
P00505	GOT2	GO:0004069	aspartate transaminase activity	PMID:2731362	IDA

A GO annotation is ...

...a statement that a gene product;

1. has a particular **molecular function**
or is involved in a particular **biological process**
or is located within a certain **cellular component**




<i>Accession</i>	<i>Name</i>	<i>GO ID</i>	<i>GO term name</i>	<i>Reference</i>	<i>Evidence code</i>
P00505	GOT2	GO:0004069	aspartate transaminase activity	PMID:2731362	IDA

A GO annotation is ...

...a statement that a gene product;

1. has a particular **molecular function**
or is involved in a particular **biological process**
or is located within a certain **cellular component**
2. as described in a particular reference




<i>Accession</i>	<i>Name</i>	<i>GO ID</i>	<i>GO term name</i>	<i>Reference</i>	<i>Evidence code</i>
P00505	GOT2	GO:0004069	aspartate transaminase activity	PMID:2731362	IDA

A GO annotation is ...

...a statement that a gene product;

1. has a particular **molecular function**
or is involved in a particular **biological process**
or is located within a certain **cellular component**
2. as described in a particular reference
3. according to a particular method



<i>Accession</i>	<i>Name</i>	<i>GO ID</i>	<i>GO term name</i>	<i>Reference</i>	<i>Evidence code</i>
P00505	GOT2	GO:0004069	aspartate transaminase activity	PMID:2731362	IDA

Kinds of evidence for GO annotations by curators



Experimental data

[Inferred from Experiment \(EXP\)](#)
[Inferred from Direct Assay \(IDA\)](#)
[Inferred from Physical Interaction \(IPI\)](#)
[Inferred from Mutant Phenotype \(IMP\)](#)
[Inferred from Genetic Interaction \(IGI\)](#)
[Inferred from Expression Pattern \(IEP\)](#)



Computational analysis

[Inferred from Sequence or structural Similarity \(ISS\)](#)
[Inferred from Sequence Orthology \(ISO\)](#)
[Inferred from Sequence Alignment \(ISA\)](#)
[Inferred from Sequence Model \(ISM\)](#)
[Inferred from Genomic Context \(IGC\)](#)
[Inferred from Biological aspect of Ancestor \(IBA\)](#)
[Inferred from Biological aspect of Descendant \(IBD\)](#)
[Inferred from Key Residues \(IKR\)](#)
[Inferred from Rapid Divergence \(IRD\)](#)
[Inferred from Reviewed Computational Analysis \(RCA\)](#)



Author statements/ curator inference

[Traceable Author Statement \(TAS\)](#)
[Non-traceable Author Statement \(NAS\)](#)
[Inferred by Curator \(IC\)](#)
[No biological Data available \(ND\) evidence code](#)

Inferred from Electronic Annotation (IEA)

The one evidence code used for completely automated annotation

Manual annotations

- Time-consuming process producing lower numbers of annotations (~2,800 taxons covered)
- More specific GO terms
- Manual annotation is essential for creating predictions



Aleksandra
Shypitsyna



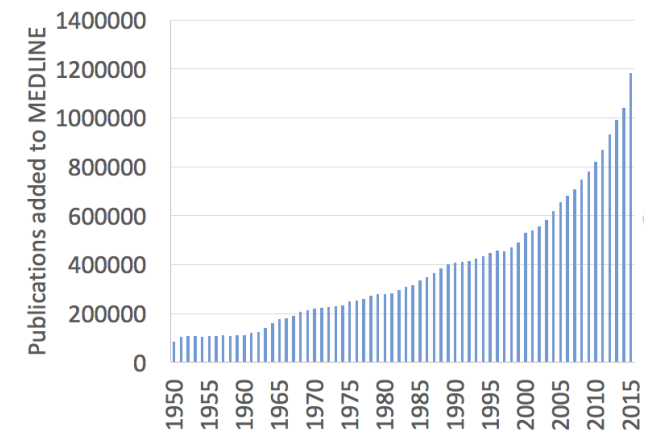
Elena
Speretta



Alex
Holmes

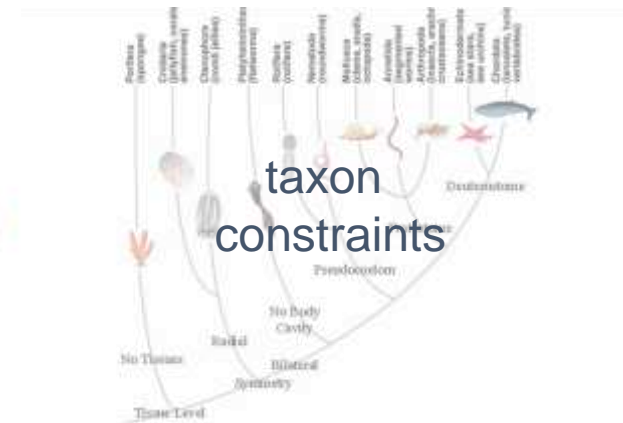


Tony
Sawford



Electronic Annotations (IEA)

- Quick way of producing large numbers of annotations
- Annotations use less-specific GO terms
- Only source of annotation for ~438,000 non-model organism species



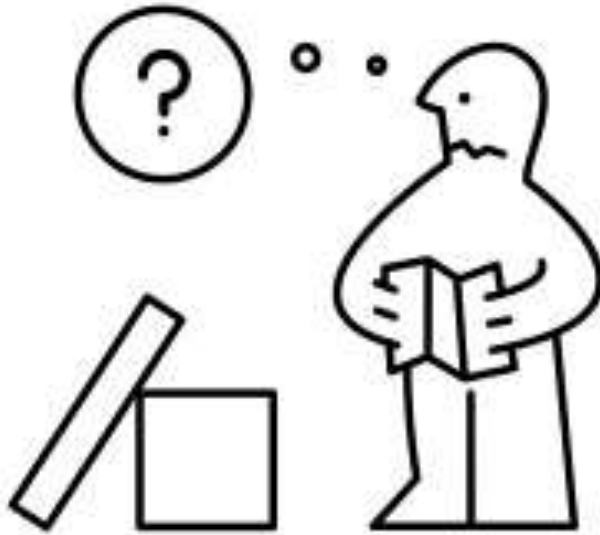
A public resource of data and tools

Number of annotations in UniProt-GOA database (March 2016)

Electronic annotations	269,207,317
Manual annotations*	2,752,604

* Includes manual annotations integrated from external model organism and specialist groups

<https://www.ebi.ac.uk/QuickGO/>
<http://www.ebi.ac.uk/GOA>



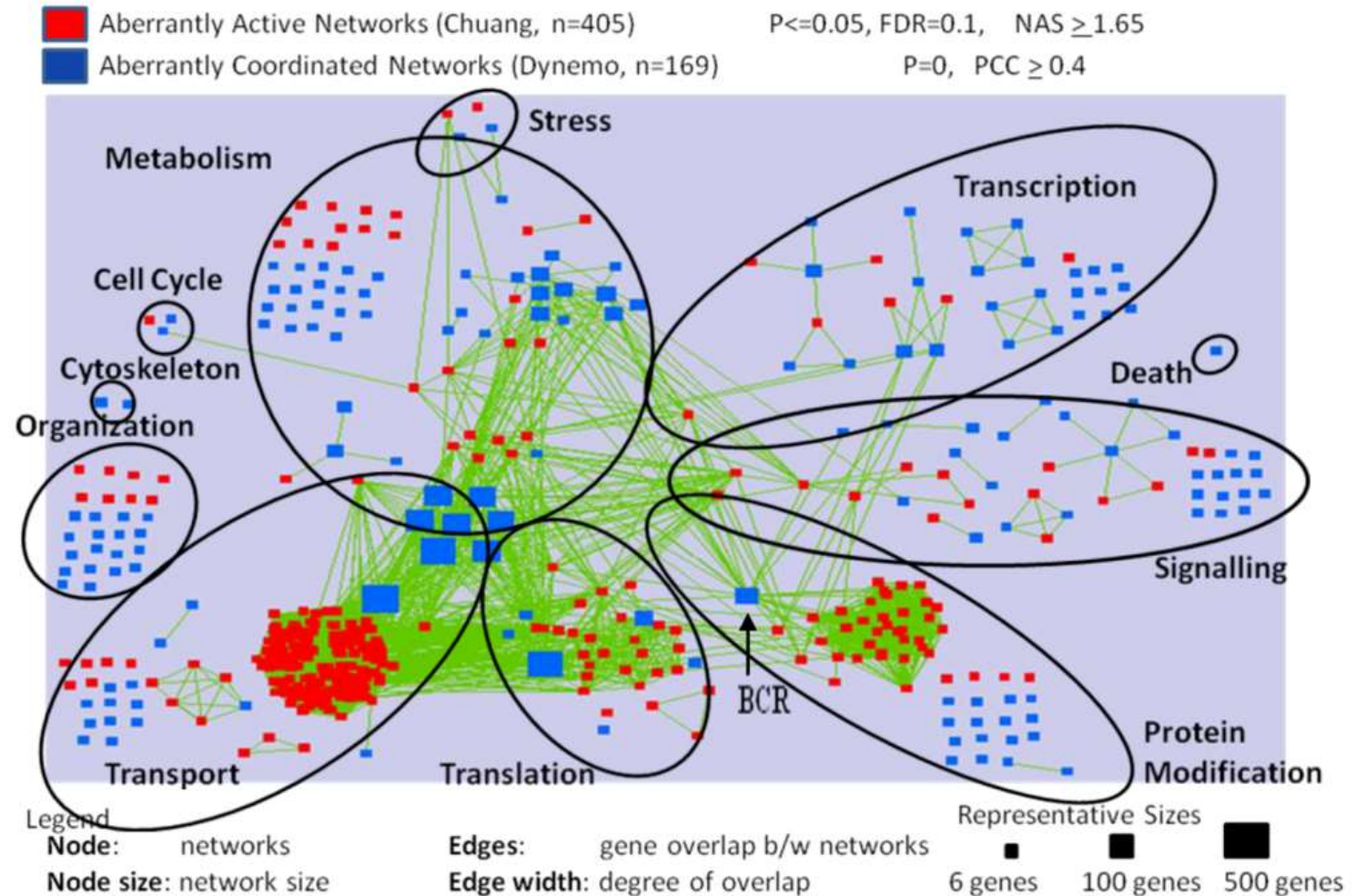
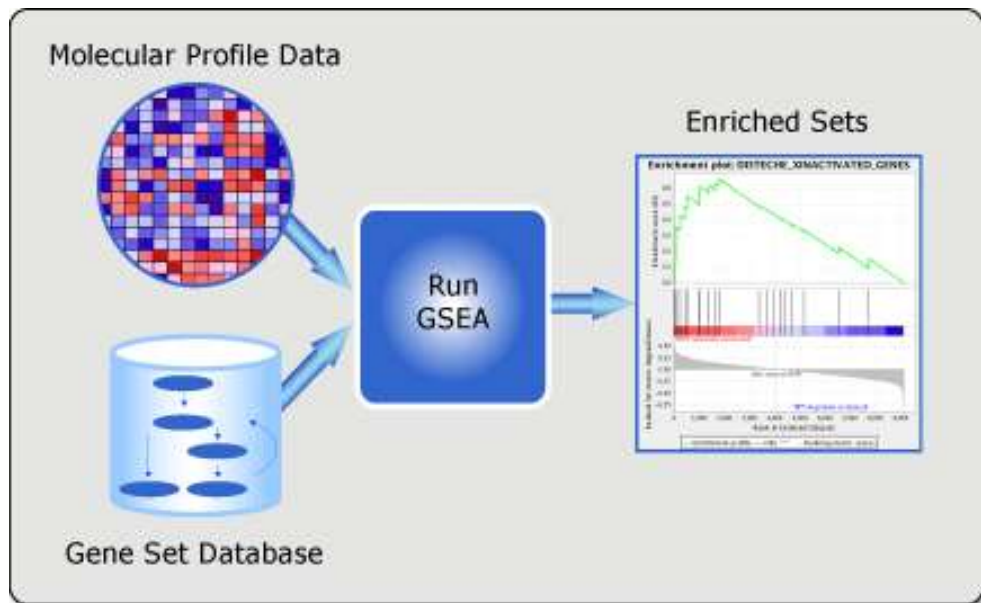


Using AMIGO2:

<http://amigo.geneontology.org>

- Find the Gene Ontology term for [Nucleus](#)
- Find its child term [Pronucleus](#)
- Find a C. Elegans gene associated with this term and find the PubMed id of the reference supporting the annotation
- Repeat for a human gene, what is the evidence for the annotation?

Gene Set Enrichment Analysis (previously covered)

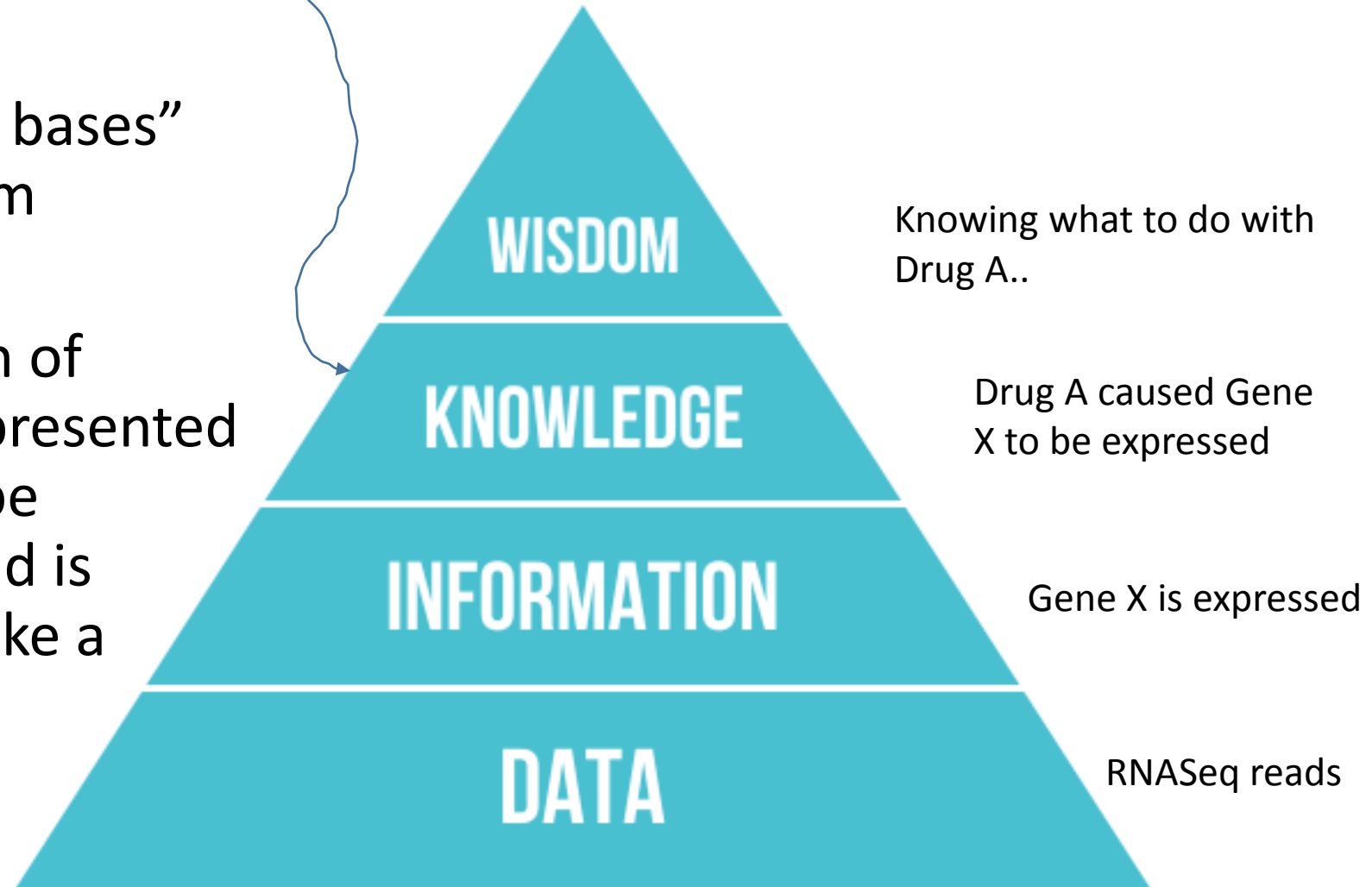


Questions about GO or other ontologies?

Part 3: Knowledge graphs

Knowledge Graphs

- Also called “knowledge bases” to distinguish them from databases.
- An integrated collection of assertions or claims represented in something that can be visualized as a graph and is technically very much like a database.



Example knowledge graphs

- Wikidata: The structured equivalent of Wikipedia
 - <http://wikidata.org>
- UniProt Knowledge Base: Manually curated Protein knowledge base
 - <http://www.uniprot.org/uniprot/>
- Microsoft Knowledge Graph (“Satori”)
- Google Knowledge Graph

Example: “Google Knowledge Graph” (GKG)

Vemurafenib
405,000 results
1 infobox
1 node in GKG

The image shows a Google search for "vemurafenib". The search results page displays several links, including Wikipedia, Genentech's patient information, the NCI Drug Dictionary, a NEJM article, and the FDA website. On the right side, a detailed infobox for "Vemurafenib" is shown. Red arrows indicate the mapping from search results to the infobox sections:

- From the Wikipedia result to the "Chemotherapy" section.
- From the Genentech result to the "SIDE EFFECTS", "INTERACTIONS", and "WARNINGS" sections.
- From the NCI Drug Dictionary result to the "Brands" and "Availability" sections.
- From the NEJM article result to the "May treat" section.
- From the FDA website result to the "Related medications" section.

The infobox for Vemurafenib includes the following information:

- Vemurafenib** (Common brands: Zelboraf)
- Chemotherapy**: It can treat melanoma.
- SIDE EFFECTS**, **INTERACTIONS**, **WARNINGS**
- Brands**: Zelboraf
- Availability**: Prescription needed
- Pregnancy**: Consult a doctor
- Alcohol**: No known interactions with light drinking
- Drug class**: BRAF kinase inhibitor antineoplastic agent
- May treat**: Melanoma (The most serious type of skin cancer), Common
- Related medications**: Romidepsin (Istodax), Paclitaxel (Abraxane), Pentostatin

Why Knowledge Graphs?

- Answer explicit questions
- Uncover implicit relations

Vemurafenib ⓘ
Common brands: Zelboraf

Chemotherapy
It can treat melanoma.

[SIDE EFFECTS](#) [INTERACTIONS](#) [WARNINGS](#)


Brands: Zelboraf
Availability: Prescription needed
Pregnancy: Consult a doctor
Alcohol: No known interactions with light drinking
Drug class: BRAF kinase inhibitor antineoplastic agent

May treat

Melanoma
The most serious type of skin cancer.

Common

B-Raf proto-oncogene, serine/threonine kinase



Available structures

PDB Ortholog search: [PDB](#) [RCSB](#)
List of PDB id codes [\[show\]](#)

Identifiers

Aliases [BRAF](#), [B-RAF1](#), [BRAF1](#), [NS7](#), [RAFB1](#)

External IDs [MGI: 88190](#) [HomoloGene: 3197](#) [GeneCards: 673](#)

Breast cancer

[ABOUT](#) [SYMPTOMS](#) [TREATMENTS](#)



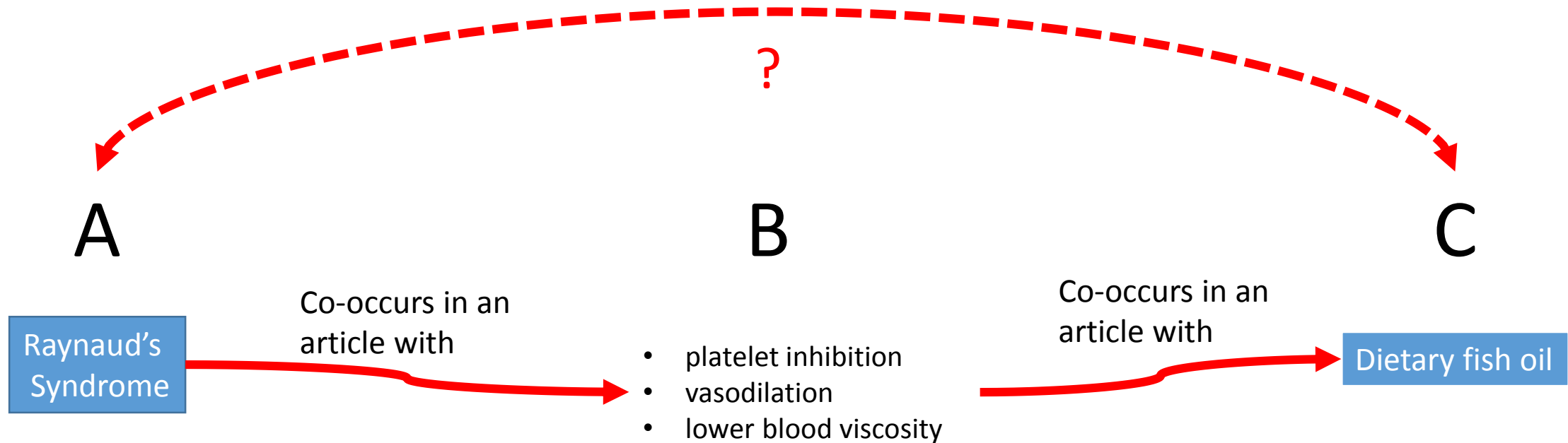
A cancer that forms in the cells of the breasts.

Common
More than 200,000 US cases per year

-  Treatable by a medical professional
-  Requires a medical diagnosis
-  Lab tests or imaging always required

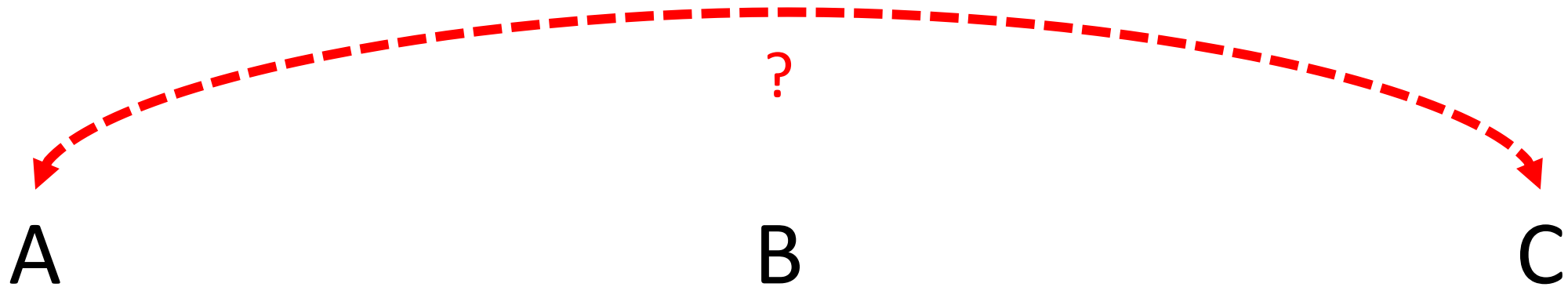
Breast cancer can occur in women and rarely in men.
Symptoms of breast cancer include a lump in the breast, bloody discharge from the nipple, and changes in the shape or texture of the nipple or breast.
Treatment depends on the stage of cancer. It may consist of chemotherapy, radiation, and surgery.

Implicit relations for hypothesis generation ABC model



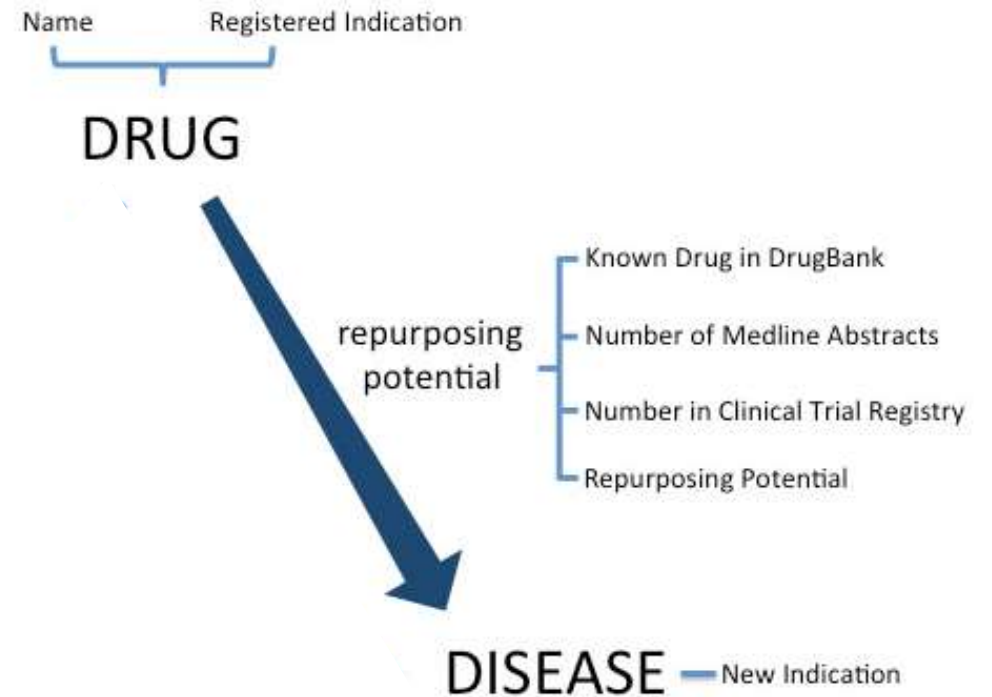
Open Discovery and Closed Discovery

- Open, you don't know what C or B is (e.g. disease -> ?drug)
- Closed, you know what C is and are looking for B (e.g. disease – why? – drug)



Example question: drug repurposing

- For a given drug, what diseases might it be used to treat?

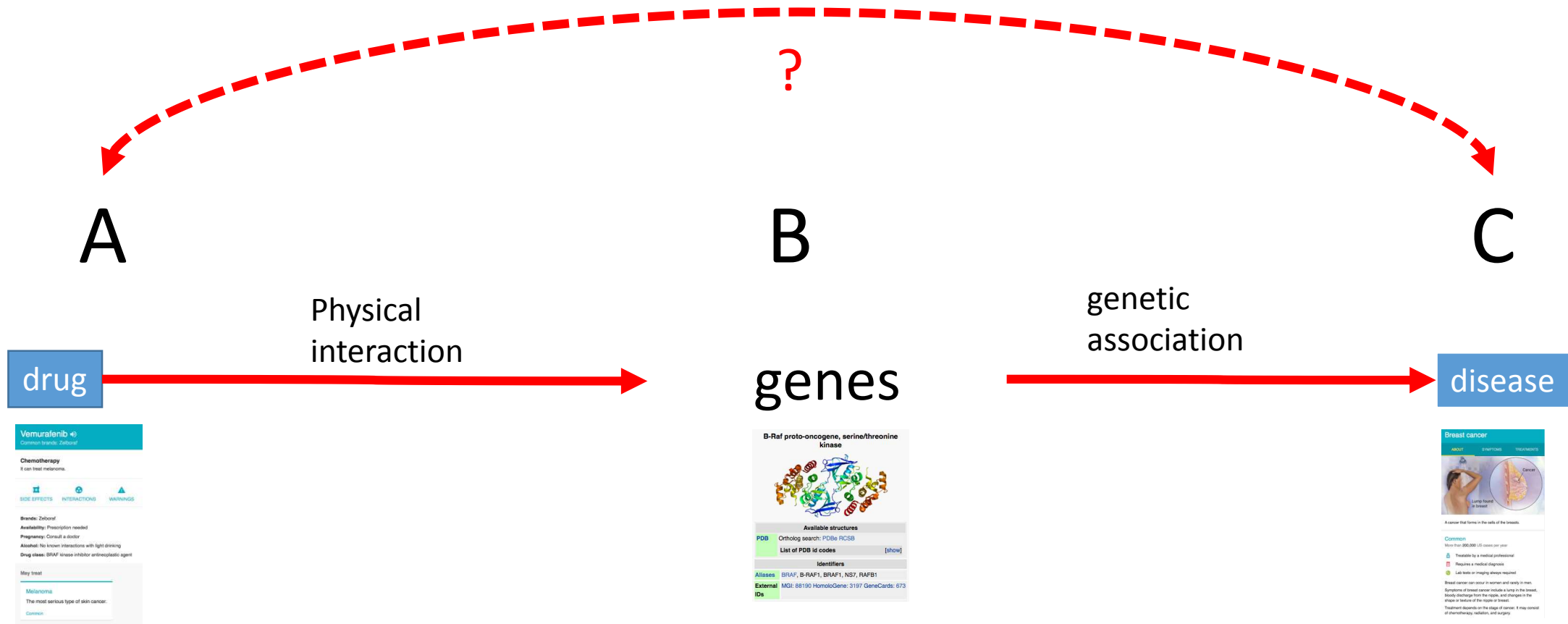


'RE:fine drugs': an interactive dashboard to access drug repurposing opportunities.

<http://www.ncbi.nlm.nih.gov/pubmed/27189611>

Implicit relations for hypothesis generation

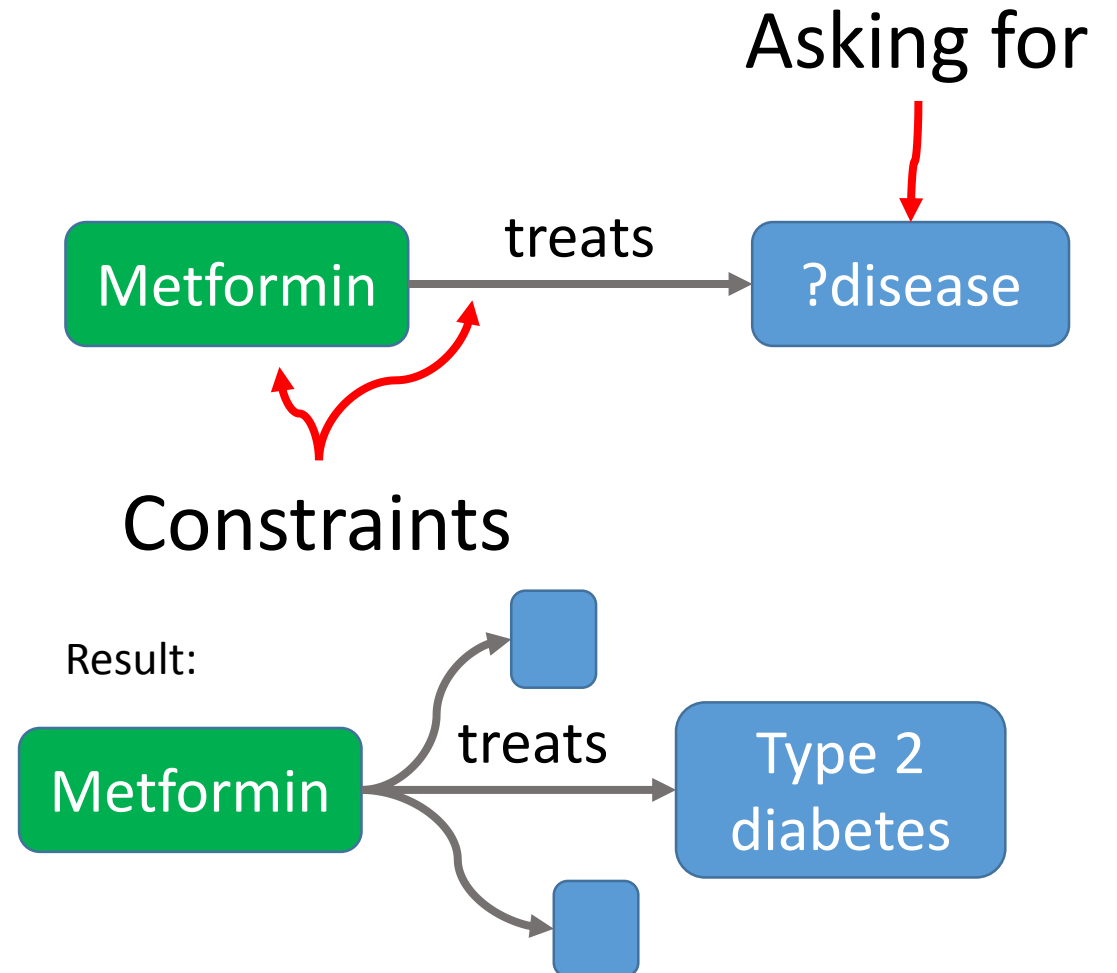
ABC model for drug repurposing



Questions on ABC model ?

Querying a knowledge graph with SPARQL

- “SPARQL protocol and RDF query language”
- RDF: Resource Description Framework (common standard for storing knowledge graphs)
- **A SPARQL query = a partially completed **graph****
 - ?’s show what you are looking for
 - rest constrains the search



#what disease does the drug treat ?

```
SELECT ?disease WHERE {  
  wd:Q19484 wdt:P2175 ?disease .  
}
```

Metformin's unique id: Q19484

Treats property id: P2175



disease

Q wd:Q3025883

Q wd:Q6717002

HANDS ON

<https://query.wikidata.org/>

Metformin's unique id: Q19484
Treats property id: P2175

#what disease does the drug treat ?

```
SELECT ?disease ?diseaseLabel WHERE {
```

```
wd:Q19484 wdt:P2175 ?disease .
```

#add the labels

```
SERVICE wikibase:label {
```

```
bd:serviceParam wikibase:language "en" .
```

```
}
```

```
}
```

Metformin

treats

?disease

disease

Q wd:Q3025883

Q wd:Q6717002

diseaseLabel

diabetes mellitus type 2

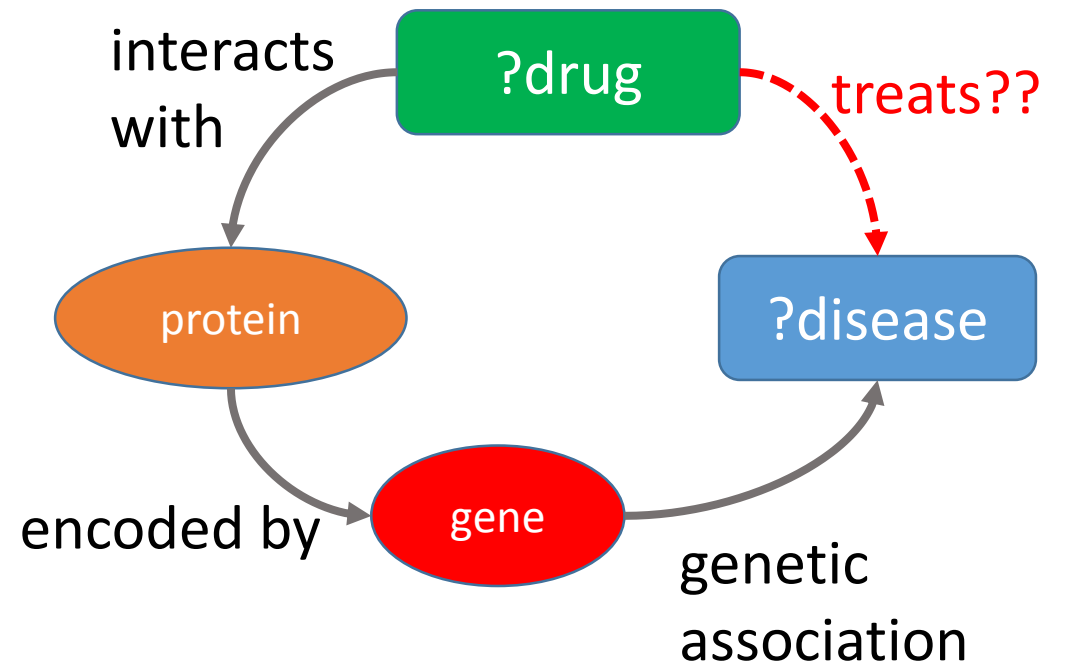
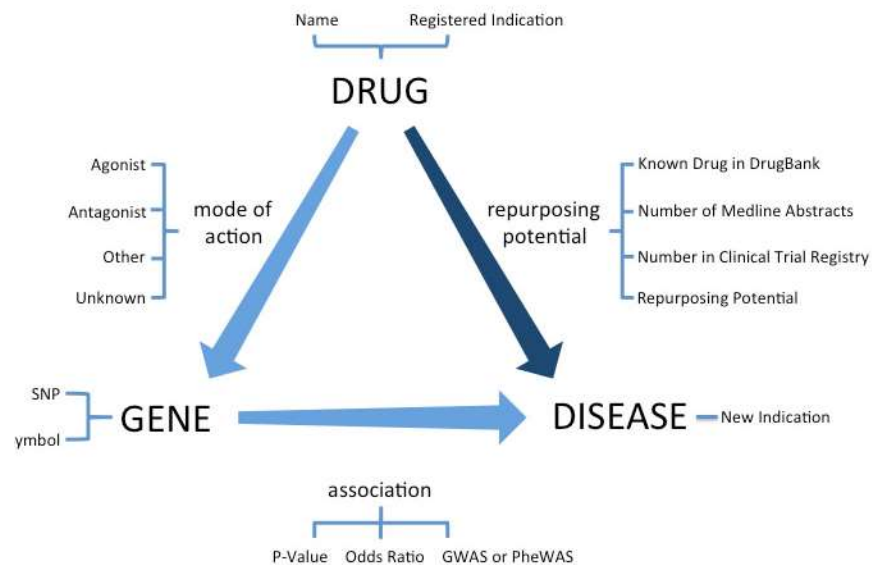
MODY 2

HANDS ON

<http://tinyurl.com/gwd6pep>

Example question: drug repurposing

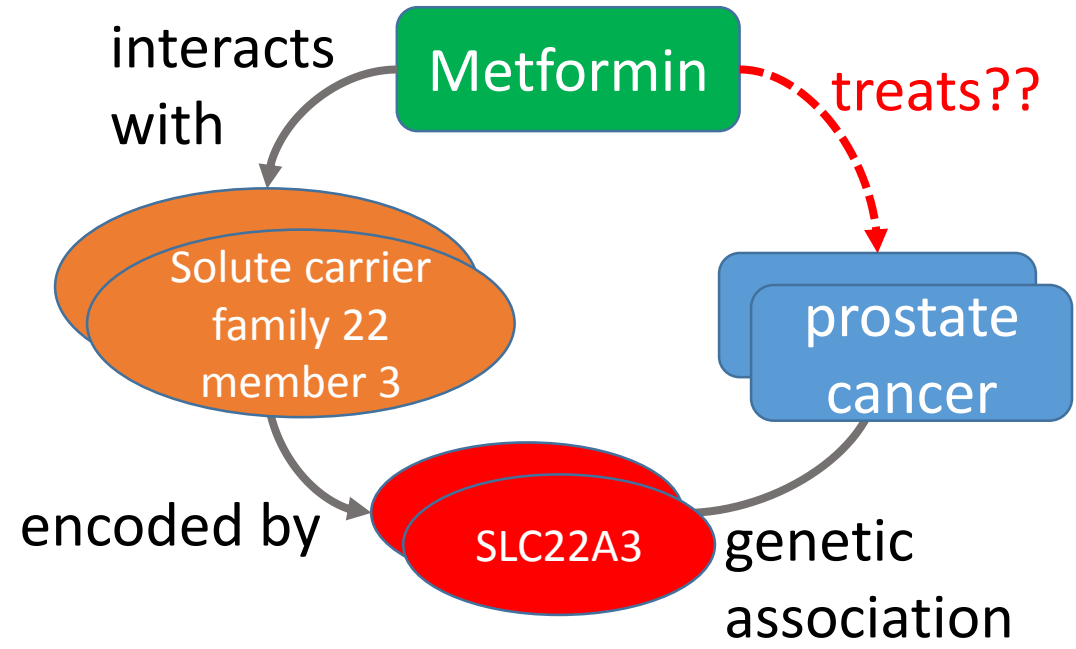
“For a given drug, what diseases might it be used to treat?”



Example question: repurposing Metformin

```
SELECT ?gene ?geneLabel ?disease ?diseaseLabel WHERE {  
  wd:Q19484 wdt:P129 ?gene_product . # Metformin interacts with a gene_product  
  ?gene_product wdt:P702 ?gene . # gene_product is encoded by a gene  
  ?gene wdt:P2293 ?disease . # gene is genetically associated with a disease  
  # add labels  
  SERVICE wikibase:label {  
    bd:serviceParam wikibase:language "en" .  
  }  
}
```

geneLabel	diseaseLabel
solute carrier family 22 (organic cation transporter), member 3	hepatitis C
solute carrier family 22 (organic cation transporter), member 3	prostate cancer
solute carrier family 22 (organic cation transporter), member 3	colorectal cancer
solute carrier family 22 (organic cation transporter), member 2	nephropathy



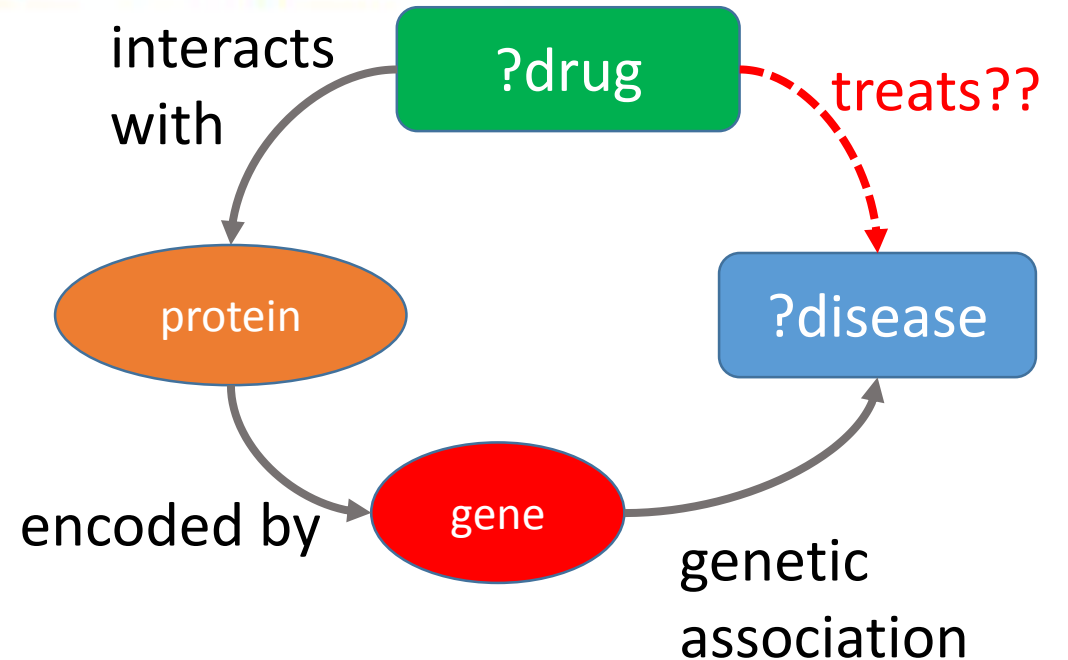
<http://tinyurl.com/zem3oxz>

Aside

- “Validating drug repurposing signals using electronic health records: a case study of metformin associated with reduced cancer mortality”
- <https://jamia.oxfordjournals.org/content/22/1/179>

Example question: repurposing all drugs

```
SELECT ?drug ?drugLabel ?gene ?geneLabel ?disease ?diseaseLabel WHERE {  
  ?drug wdt:P129 ?gene_product . # drug interacts with a gene_product  
  ?gene_product wdt:P702 ?gene . # gene_product is encoded by a gene  
  ?gene wdt:P2293 ?disease . # gene is genetically associated with a disease  
# add labels  
  SERVICE wikibase:label {  
    bd:serviceParam wikibase:language "en" .  
  }  
}  
limit 1000
```



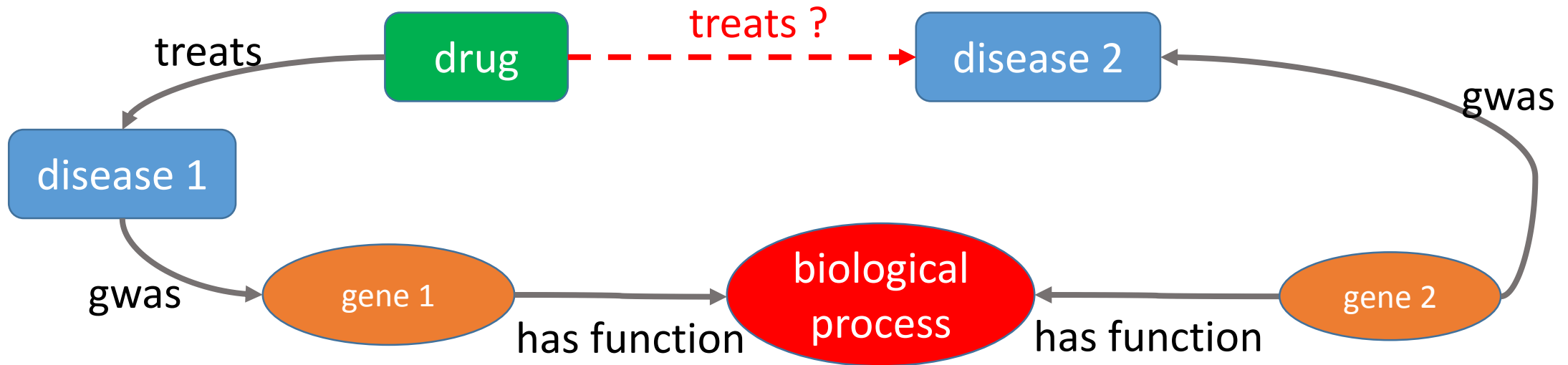
<http://tinyurl.com/hwm9388>

Adding constraints

- Find drugs that may treat disease
 - according to the drug->gene->disease model
 - constrained to focus on cancers
 - **?disease wdt:P279* wd:Q12078 .**
 - limited to genes related to cell proliferation
 - **?gene_product wdt:P682 ?biological_process**
 - **?biological_process wdt:P279* wd:Q14818032**
- <http://tinyurl.com/j222k6g>

Other patterns?

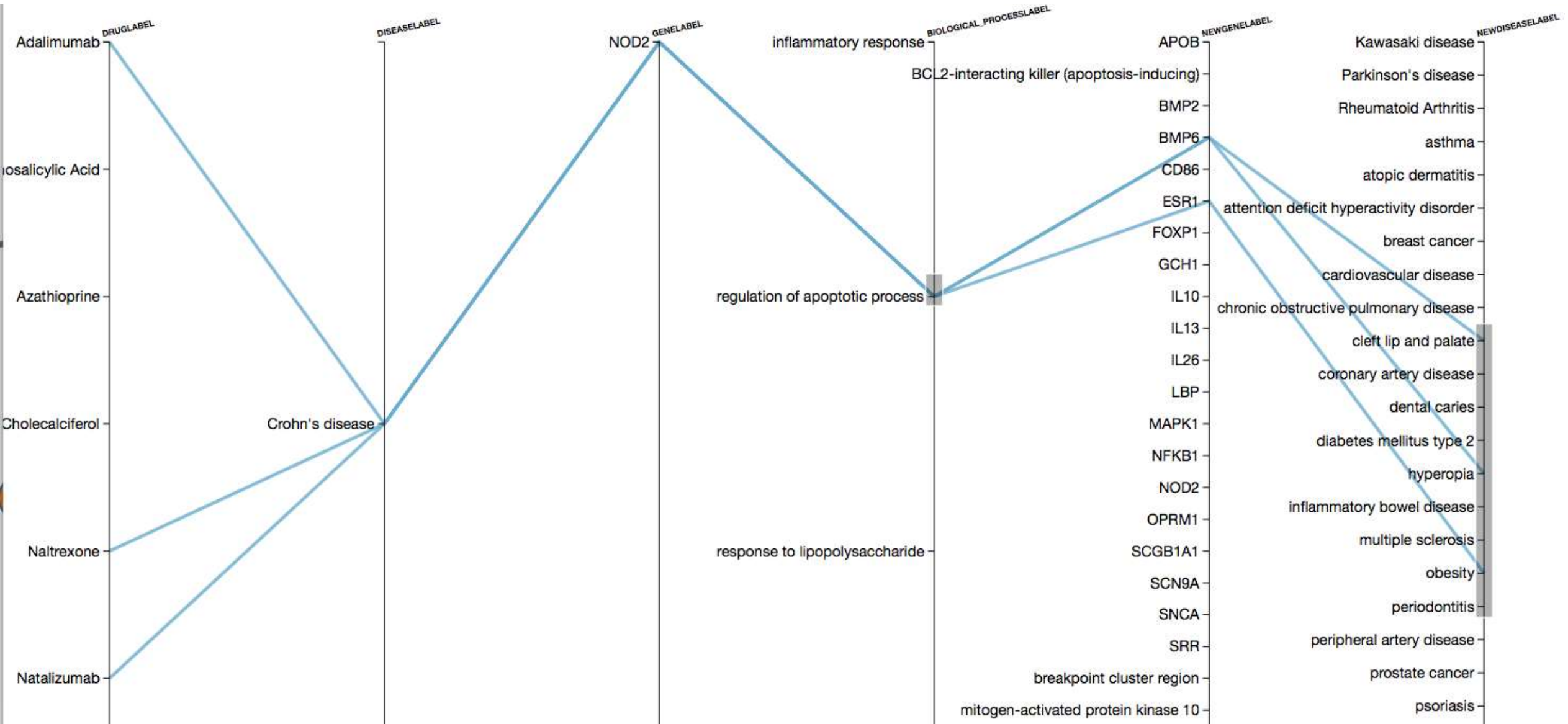
Is there a connecting path in the knowledge graph?
Is it meaningful?



<http://tinyurl.com/gpfr9kj>

Beta result viewer,
<http://jonaskress.github.io/>

<http://tinyurl.com/jmoczaq>

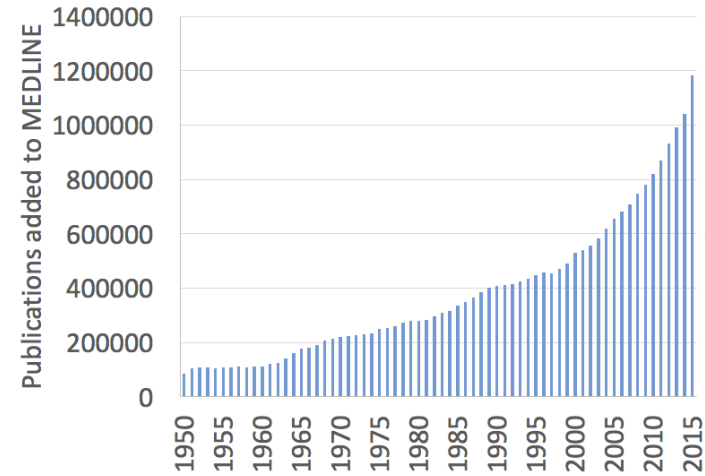


SPARQL endpoints of interest

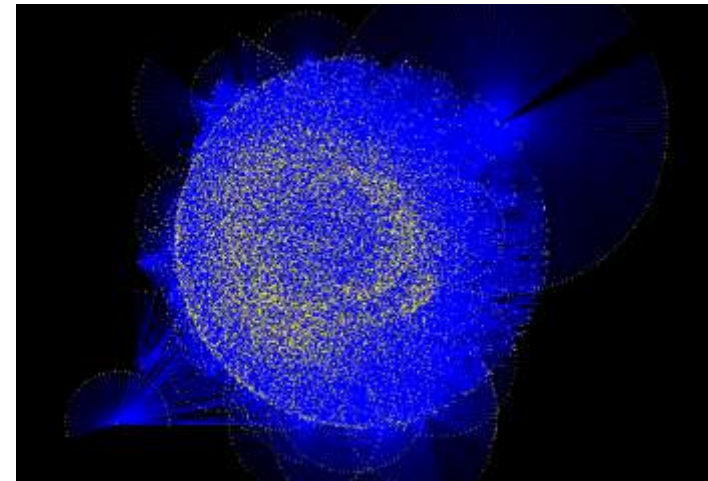
- Wikidata <http://query.wikidata.org>
- UniProt <http://sparql.uniprot.org>
- MeSH <https://id.nlm.nih.gov/mesh/query>
- EBI <https://www.ebi.ac.uk/rdf/documentation/sparql-endpoints>
- Bio2RDF <https://github.com/bio2rdf/bio2rdf-scripts/wiki/Query-repository>

2 problems with knowledge graphs

Not enough knowledge in the graph
text and data mining
crowdsourcing ?



Too much knowledge in the graph
sorting algorithms
visualizations



<http://i9606.blogspot.com/2010/05/gene-wiki-hairball-1.html>

Plan for Thursday / Homework

- Implement and apply an ABC Model style hypothesis generating program
- Assignment: write the program, explain its logic, explain how you used it to generate a hypothesis, explain the hypothesis
- A Jupyter notebook with Python code will be provided to get you started
- If you do not want to program, there will be another option using online tools.

Suggested Reading

- Ontology
 - Biomedical Ontologies in Action: Role in Knowledge Management, Data Integration and Decision Support
 - <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2592252/>
 - Gene Ontology: tool for the unification of biology
 - <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3037419/>
- Knowledge-based hypothesis generation
 - Fish oil, Raynaud's syndrome and undiscovered public knowledge
 - <http://muse.jhu.edu/article/403510/pdf>
 - Knowledge discovery by automated identification and ranking of implicit relationships
 - <http://bioinformatics.oxfordjournals.org/content/20/3/389.full.pdf>
- Text mining
 - Literature mining for the biologist: from information retrieval to biological discovery
 - <http://www.nature.com/nrg/journal/v7/n2/full/nrg1768.html>