# CACAO Training
# Part 1

Jim Hu and Suzi Aleksander

Fall 2015

# Outline

- Part 1: Gene Ontology and functional annotation
  - How known functions are used to reveal new knowledge
  - Gene Ontology
  - What is an annotation?
  - CACAO
- Part 2: Making annotations and challenges

# Leveraging what we know about function

# Leveraging What We Know About Function

- **Functional profiling**: For a list of genes, what functions are important?
  - Genes turned up or down together
    - Disease states
    - Environmentalresponses
    - Genotypes
    - …
  - Genes encoding proteins that physically interact
  - Genes conserved in specific taxa
  - Genes found in specific microbial communities

# Functional Profiling Example

## Sister grouping of chimpanzees and humans as revealed by genome-wide phylogenetic analysis of brain gene expression profiles

Monica Uddin[†‡], Derek E. Wildman[†‡], Guozhen Liu[†§], Wenbo Xu[§], Robert M. Johnson[¶], Patrick R. Hof[∥], Gregory Kapatos[†,††], Lawrence I. Grossman[†], and Morris Goodman[†‡,‡‡]

[†]Center for Molecular Medicine and Genetics, Departments of [‡]Anatomy and Cell Biology, [¶]Biochemistry and Molecular Biology, and [††]Psychiatry and Behavioral Neurosciences, Wayne State University School of Medicine, 540 East Canfield Avenue, Detroit, MI 48201; [§]Bioinformatics Facility, 5107 Biological Science Building, 5047 Gullen Mall, Detroit, MI 48202; and [∥]Department of Neurobiology, Mount Sinai School of Medicine, One Gustave L. Levy Place, New York, NY 10029

Gene expression profiles from the anterior cingulate cortex (ACC) of human, chimpanzee, gorilla, and macaque samples provide clues about genetic regulatory changes in human and other catarrhine primate brains. The ACC, a cerebral neocortical region, has human-specific histological features. Physiologically, an individual's ACC displays increased activity during that individual's performance of cognitive tasks. Of ≈45,000 probe sets on microarray chips represent-ing transcripts of all or most h... detected in human ACC samp... 15,000, in gorilla and chimpan... obtained from gene expressi... expectation that the non-hu... gorilla) should be more like... humans. Instead, the chimpa... human than like the gorilla;... panzees are the sister group... biguous expression changes... cesses and molecular functi... represented in the data, the ch... apparent regulatory evolutio... important changes in the ance... but to a greater extent in hu... profiles of aerobic energy metabolism genes and neuronal function-related genes, suggesting that increased neuronal activity required increased supplies of energy.

more vulnerable to Alzheimer's disease than are other pyramidal neurons (17, 18). Physiologically, brain imaging results show in-creased activity in an individual's ACC when that individual is engaged in cognitive tasks (19–21). The ACC participates in decision making when interfering choices are present, a cognitive role involved in executive function (22). In view of these histological and physiological findings, it seemed likely to us that comparative...

...structed the phylogenetic history of the ACC gene expression profiles by treating each probe set as a single character, e.g., analogous to a single genomic locus or a single position in a
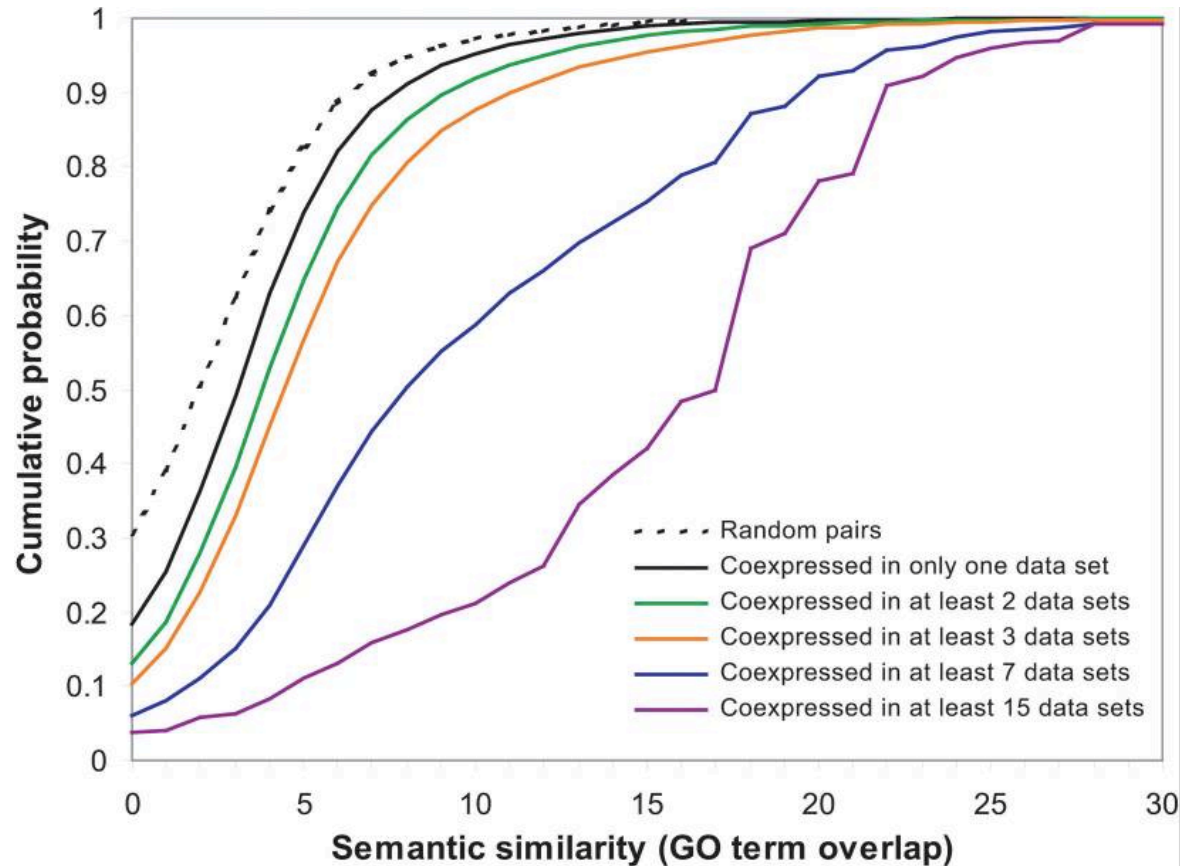
> Among important changes in the ancestry of both humans and chimpanzees, but to a greater extent in humans, are the up-regulated expression profiles of aerobic energy metabolism genes and neuronal function-related genes, suggesting that increased neuronal activity required increased supplies of energy.

Uddin et al. (2004) PNAS 101:2957-2962

# Leveraging What We Know About Function

- **Guilt by association**: For a gene of unknown function, can we infer its function from genes of known function:
    - that are coexpressed across many conditions
    - that are homologs
    - that are coinherited across evolution
    - that physically interact in a multiprotein complex

# Coexpression Correlates with Common Function



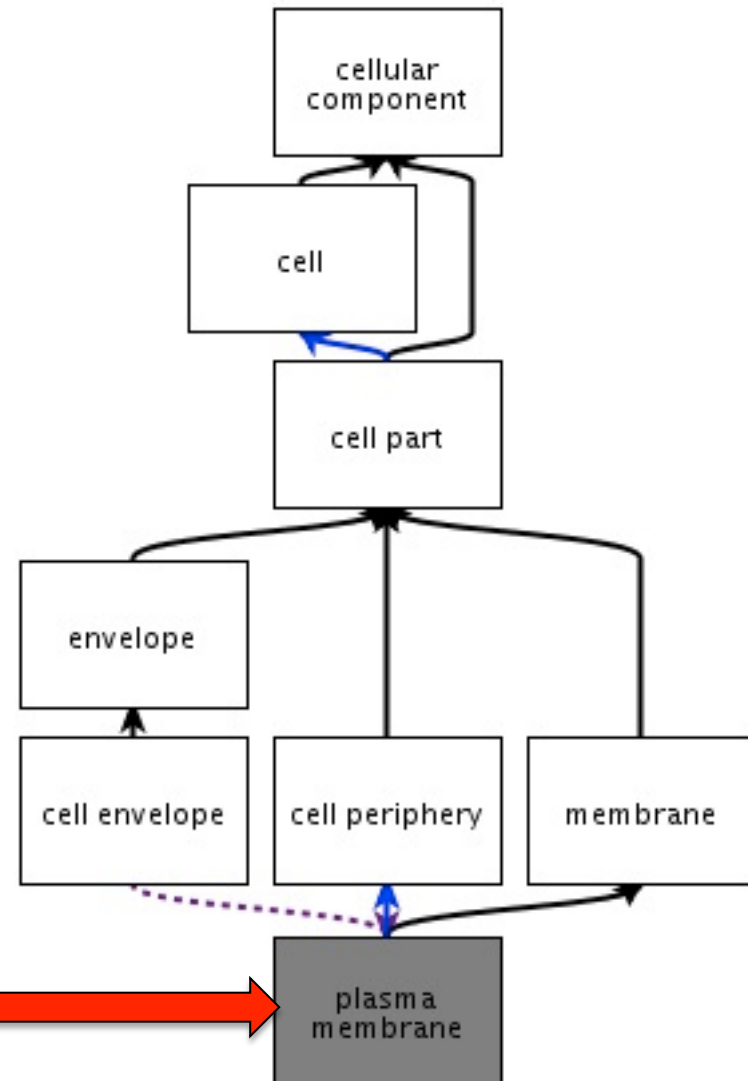Lee et al. (2004) Genome Research 14:1085-1094

# What Do We Mean by Function?

- Massive body of published knowledge
  - Almost useless by itself!!
- We need
  - Knowledge that computers can analyze
  - Common vocabulary across different organisms
  - Disambiguation of synonyms
  - Connection of related ideas that are more or less specific
    - Examples:
      - polygon – quadrilateral – rectangle – square
      - Enzyme – kinase – protein kinase – protein tyrosine kinase

# GENE ONTOLOGY

# Gene Ontology (GO)

- 3 aspects (ontologies) :
  - Molecular Function
  - Biological Process
  - Cellular Component
- Controlled vocabulary
  - ID number for computers
  - Name and definition for humans
- Relationships



GO:0005886

# Molecular Function

- activities = what a protein can do by itself



GO:0004347    hexokinase activity



GO:0016301   Kinase activity

# Biological Process

- a commonly recognized series of events
  - Including, but not just biochemical pathways



GO:0051301
cell division



GO:0006351
transcription, DNA dependent

# Cellular Component

- where a gene product acts
  - Subcellular location
  - Multicomponent complex



GO:0005739
mitochondrion

GO:0009274
peptidoglycan-based cell wall

GO:0005840
ribosome

From visualphotos.com, epmm.group.shef.ac.uk, wikimedia commons

# GO terms

- ID numbers
- Definitions
- Relationships
  - Directed Acyclic Graph
- GO terms provide a way to describe functions, now we have to associate them with genes!
  - AKA GO annotation

# GO Annotation

# What is Annotation?

**Dictionary** Thesaurus    🔍 annotation

## an•no•ta•tion |ˌanəˈtā sH ən|

noun

a note of explanation or comment added to a text or diagram : *marginal annotations.*

• the action of annotating a text or diagram : *annotation of prescribed texts.*

ORIGIN late Middle English : from French, or from Latin ***annotatio(n-)****,* from the verb ***annotare*** (see ANNOTATE ).

# What is Annotation?

# What is Annotation?

# What is Annotation?



http://www.shakespeare-navigators.com/hamlet/H31.html

# Levels of Annotation for Genomes

- Metadata
  - What is this genome?
- Features
  - Where are things in the sequence?
- Products
  - What do we know about the features?
- Systems
  - What do the products do?
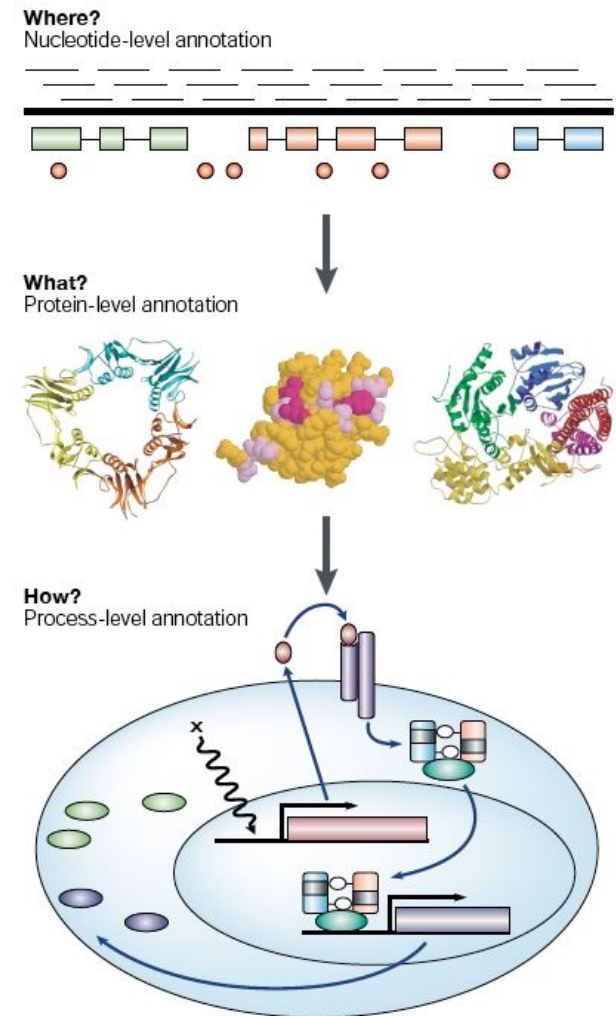    - individually?
    - Working together?



Where?
Nucleotide-level annotation

What?
Protein-level annotation

How?
Process-level annotation

Figure 1 | The three layers of genome annotation: where, what and how?

# Functional Annotation w/GO

- Annotation: a note that is made while reading any form of text

- GO Annotation: a database entry in a **specific format** that associates a **GO term** with a **gene product** made based on **evidence** in a peer-reviewed **paper**
  - Specific format makes the annotations readable by both computers and humans
  - GO annotations capture the chain of evidence for how functions were inferred from experiments
  - More when we talk about CACAO

# Where Do Annotations Come From?



Literature

Datasets

Biocurators
(rate limiting)

Database

# Databases Need Help!

- >24 million peer-reviewed articles in PubMed
- Many millions of proteins recorded in UniProt



http://www.uniprot.org

CACAO

# What is CACAO?

- **C**ommunity **A**ssessment of **C**ommunity **A**nnotation with **O**ntologies  (CACAO)
  - Annotation of gene function
  - Competition
    - Within a class
    - Between teams at different schools
    - More details next week

# How Does CACAO Work?

- Working in teams we will use the GONUTS website:
  - http://gowiki.tamu.edu
- Multiple innings: each is two weeks
  - Annotation week: you make annotations on the website to get points
  - Challenge week: you challenge annotations made by other teams to steal their points
- Open week
  - You can challenge AND annotate in the same week
  - RARE but help balance out holiday breaks, etc. between schools
- You can make as many annotations as you want.
  - You pick the topic
  - You have to convince us that they are correct.
    - The default is that they are wrong!!
- Your annotations could end up in databases used by researchers all over the world

# How Does CACAO Work?

- Getting help is not cheating!
  - Talk to your teammates
  - Ask us questions
  - Talk to other professors
  - Email authors of papers

# What To Annotate

- You can start with a paper
  - Find the proteins discussed, potential GO terms
- You can start with a GO term
  - Modify PubMed search with keywords or organism
- You can start with a protein
  - Find papers about the protein
- Don't get stuck on what you started with
  - Your first paper may not have **experiments** about function (Reviews)
  - Reading about your initial protein may lead you to better information about other proteins

# Functional Annotation w/GO

- Annotation:  a note that is made while reading any form of text


- GO Annotation:  a database entry in a **specific format** that associates a **GO term** with a **gene product** made based on **evidence** in a peer-reviewed <span style="color:red">**paper**</span>

# Starting with a paper

- Need a scientific paper with experimental data
  - Use PubMed: http://www.ncbi.nlm.nih.gov/pubmed/
    - Or use an alias like http://pubmed.com
  - No review articles, no books, no textbooks, no Wikipedia articles, no class notes…
  - BUT you could start with those!
  - DON'T start with the first paper you see from a random PubMed search

# Starting with a paper

- Need a scientific paper with experimental data
  - PubMed review?
  - We refer to the paper through the PMID number
    - Not the full citation

# PubMed Record

# Getting the Full Text



- The abstract is not enough for an annotation
  - But, may be enough to reject a paper!!!

# Getting the Full Text

- Some papers are open access
  - Pubmed Central
  - Journal sites

- Others are pay only
  - Don't pay real $$!
    - Your library may have subscriptions
    - Pick a different paper
    - Email the author and ask for a pdf
      - Send us a copy

# Alternative Path: Start w/Full Text

# Alternative Path: Start w/Full Text

# Beware!

- Good science ≠ good for annotation

## Second Extracellular Loop of Human Glucagon-like Peptide-1 Receptor (GLP-1R) Differentially Regulates Orthosteric but Not Allosteric Agonist Binding and Function*[S]

Cassandra Koole[‡], Denise Wootten[‡], John Simms[‡], Emilia E. Savage[‡], Laurence J. Miller[§], Arthur Christopoulos[‡1], and Patrick M. Sexton[‡2]

From the [‡]Drug Discovery Biology, Monash Institute of Pharmaceutical Sciences and Department of Pharmacology, Monash University, Parkville, Victoria 3052, Australia and the [§]Department of Molecular Pharmacology and Experimental Therapeutics, Mayo Clinic, Scottsdale, Arizona 85259

**Background:** The ECL2 of the GLP-1R is critical for GLP-1 peptide-mediated selective signaling.
**Results:** Mutation of most ECL2 residues to alanine results in changes in binding and/or efficacy of oxyntomodulin and exendin-4 but not allosteric agonists.
**Conclusion:** ECL2 of the GLP-1R has ligand-specific as well as general effects on peptide agonist-mediated receptor activation.
**Significance:** This work provides insight into control of family B GPCR activation transition.

# Beware!

- Good science ≠ good for annotation

## Robust design and optimization of retroaldol enzymes

Eric A. Althoff,[1,2] Ling Wang,[1] Lin Jiang,[1,3] Lars Giger,[4] Jonathan K. Lassila,[5] Zhizhi Wang,[1] Matthew Smith,[1] Sanjay Hari,[1] Peter Kast,[4] Daniel Herschlag,[5] Donald Hilvert,[4] and David Baker[1]*

[1]Department of Biochemistry, University of Washington and HHMI, Seattle, Washington 98195
[2]Arzeda Corp., Seattle, Washington 98102
[3]Department of Biological Chemistry, UCLA, Los Angeles, California 90095
[4]Laboratory of Organic Chemistry, ETH Zurich, 8093 Zurich, Switzerland
[5]Department of Biochemistry, Stanford University, Stanford, California 94305

# Beware!

- Good science $\neq$ good for annotation

**Short Article**

Cell PRESS

# Vitamin C Enhances the Generation of Mouse and Human Induced Pluripotent Stem Cells

Miguel Angel Esteban,[1,6] Tao Wang,[1,6] Baoming Qin,[1,6] Jiayin Yang,[1] Dajiang Qin,[1] Jinglei Cai,[1] Wen Li,[1] Zhihui Weng,[1] Jiekai Chen,[1] Su Ni,[1] Keshi Chen,[1] Yuan Li,[1] Xiaopeng Liu,[1] Jianyong Xu,[1] Shiqiang Zhang,[1] Feng Li,[1] Wenzhi He,[1] Krystyna Labuda,[2] Yancheng Song,[3] Anja Peterbauer,[4] Susanne Wolbank,[2] Heinz Redl,[2] Mei Zhong,[5] Daozhang Cai,[3] Lingwen Zeng,[1] and Duanqing Pei[1,*]
[1]Stem Cell and Cancer Biology Group, Key Laboratory of Regenerative Biology, South China Institute for Stem Cell Biology and Regenerative Medicine, Guangzhou Institutes of Biomedicine and Health, Chinese Academy of Sciences, Guangzhou 510663, China
[2]Ludwig Boltzmann Institute for Clinical and Experimental Traumatology, Austrian Cluster for Tissue Regeneration, Vienna 1200, Austria

# Beware!

- Good science ≠ good for annotation

Neurobiology of Disease

# Excess Phosphoinositide 3-Kinase Subunit Synthesis and Activity as a Novel Therapeutic Target in Fragile X Syndrome

Christina Gross,[1] Mika Nakamoto,[2]* Xiaodi Yao,[1]* Chi-Bun Chan,[3] So Y. Yim,[1] Keqiang Ye,[3] Stephen T. Warren,[2,4,5] and Gary J. Bassell[1,6]

Departments of [1]Cell Biology, [2]Human Genetics, [3]Pathology and Laboratory Medicine, [4]Biochemistry, [5]Pediatrics, and [6]Neurology, Emory University School of Medicine, Atlanta, Georgia 30322

# Functional Annotation w/GO

- Annotation:  a note that is made while reading any form of text

- GO Annotation:  a database entry in a **specific format** that associates a **GO term** with a **gene product** made based on **evidence** in a peer-reviewed **paper**

# Finding Proteins

- Search UniProt for something interesting
- Look in UniProt for the protein(s) in the paper you are reading.

**No matter what, you will need to find the protein's accession on UniProt (http://uniprot.org)**

⬇

**Use that accession to make a page for that protein on GONUTS (http://gowiki.tamu.edu)**

⬇

**Add your GO annotations to the protein's page on GONUTS**

# UniProt (http://www.uniprot.org)

- If you have a paper, look for an accession
  - UniProt accession
  - NCBI Gene ID
- If you don't have an accession, search by name/keyword

# UniProt Search Results

- Multiple entries
  - Find the right one
  - Icons
    - Gold = Swissprot = reviewed
    - Plain = TrEMBL = automated

# UniProt Records

- Lots of information to help you
  - Summary of existing GO annotations
    - Link to QuickGO for complete set of existing annotations
  - Information about the protein

# Make Sure You Have the Right Protein

- Right species/strain
- Not a fragment
- Sometimes UniProt has multiple entries for the same protein
  - Gold star = SwissProt = reviewed
  - Blank star = TrEMBL = computational entry
- Sometimes the protein you want is not in UniProt
  - May want to find another paper/protein
- Ask for help
  - OK to email the UniProt help desk
  - check your reasoning with us!

# Create a Protein Page in GONUTS

# Entering/Editing Annotations

# Functional Annotation w/GO

- Annotation: a note that is made while reading any form of text

- GO Annotation: a database entry in a **specific format** that associates a **GO term** with a **gene product** made based on **evidence** in a peer-reviewed **paper**

# Finding GO Terms

- GONUTS: http://gowiki.tamu.edu
- QuickGO: http://www.ebi.ac.uk/QuickGO
- AmiGO2: http://amigo.geneontology.org/amigo

Bmcintosh  my talk  my preferences  my watchlist  my contributions  log out

go term | discussion | edit | history | delete | protect | watch | purge

GONUTS is undergoing some *major* debugging for Pecan.
Please expect blank pages and some delays in updating.
[ Email comments to Daniel. ]

# GO:0004713 ! protein tyrosine kinase activity

**id:** GO:0004713

**name:** protein tyrosine kinase activity
**namespace:** molecular_function
**alt_id:** GO:0004718
**def:** "Catalysis of the reaction: ATP + a protein tyrosine = ADP + protein tyrosine phosphate." [EC:2.7.10]
**subset:** gosubset_prok
**synonym:** "JAK" NARROW []
**synonym:** "Janus kinase activity" NARROW []
**synonym:** "protein-tyrosine kinase activity" EXACT []
**xref:** EC:2.7.10
**xref:** MetaCyc:EC-2.7.10
**xref:** Reactome:11065 "protein tyrosine kinase activity"
**is_a:** GO:0004672 ! protein kinase activity

AmiGO

**Last version checked**
date: 14:01:2011 17:26
saved-by: rfoulger
auto-generated-by: OBO-Edit 2.0

**Last updated**
date: 08:10:2010 13:21
saved-by: dph
auto-generated-by: OBO-Edit 2.0

Gene Ontology Home
The contents of this box are automatically generated. You can help by adding information to the "Notes"

## Usage Notes                                                [edit]

## References                                                 [edit]

See Help:References for how to manage references in GONUTS.

## Child Terms

This term has the following 4 child terms.

- [+] GO:0004714 - transmembrane receptor protein tyrosine kinase activity (13)
- [ ] GO:0004715 - non-membrane spanning protein tyrosine kinase activity (1)
- [+] GO:0004716 - receptor signaling protein tyrosine kinase activity (1)
- [+] GO:0035400 - histone tyrosine kinase activity (1)

## Pages in category "GO:0004713 ! protein tyrosine kinase activity"

The following 200 pages are in this category, out of 732 total.

Show articles starting with: --All results-- Go

(previous 200) (next 200)

| C | C cont. | F cont. |
|---|---------|---------|
| CHICK:A0M8T9 | CHICK:Q90960 | FB:Tk4 |
| CHICK:A0SVH2 | CHICK:Q90961 | FB:Tk6 |
| CHICK:BTK | CHICK:Q90962 | FB:tor |

# GO:0004713 ! protein tyrosine kinase activity

**id:** GO:0004713

**name:** protein tyrosine kinase activity

**namespace:** molecular_function

**alt_id:**GO:0004718

**def:** "Catalysis of the reaction: ATP + a protein tyrosine = ADP + protein tyrosine phosphate." [EC:2.7.10]

**subset:** gosubset_prok

**synonym:** "JAK" NARROW []

**synonym:** "Janus kinase activity" NARROW []

**synonym:** "protein-tyrosine kinase activity" EXACT []

**xref:** EC:2.7.10

**xref:** MetaCyc:EC-2.7.10

**xref:** Reactome:11065 "protein tyrosine kinase activity"

**is_a:** GO:0004672 ! protein kinase activity

AmiGO



| Last version checked | Last updated |
|---|---|
| date: 14:01:2011 17:26 | date: 08:10:2010 13:21 |
| saved-by: rfoulger | saved-by: dph |
| auto-generated-by: OBO-Edit 2.0 | auto-generated-by: OBO-Edit 2.0 |

Gene Ontology Home

The contents of this box are automatically generated. You can help by adding information to the "Notes"

Bmcintosh    my talk    my preferences    my watchlist    my contributions    log out

go term | discussion | edit | history | delete | protect | watch | purge

GONUTS is undergoing some *major* debugging for Pecan.
Please expect blank pages and some delays in updating.
[ Email comments to Daniel. ]

# GO:0004713 ! protein tyrosine kinase activity

**navigation**
- Main Page
- Enter GO at the top
- Help
- What's new
- Report Bug
- Update log
- Annotation Jamborees
- Recent changes
- Create New Gene Page
- Login/Create Account

**courses**
- CACAO
- Peer Reviews

**page contributors**
- Wikientrybot

**search**

[ Go ] [ Search ] G

**toolbox**
- What links here
- Related changes
- Upload file
- Special pages
- Printable version
- Permanent link

**id:** GO:0004713

**name:** protein tyrosine kinase activity
**namespace:** molecular_function
**alt_id:** GO:0004718
**def:** "Catalysis of the reaction: ATP + a protein tyrosine = ADP + protein tyrosine phosphate." [EC:2.7.10]
**subset:** gosubset_prok
**synonym:** "JAK" NARROW []
**synonym:** "Janus kinase activity" NARROW []
**synonym:** "protein-tyrosine kinase activity" EXACT []
**xref:** EC:2.7.10
**xref:** MetaCyc:EC-2.7.10
**xref:** Reactome:11065 "protein tyrosine kinase activity"
**is_a:** GO:0004672 ! protein kinase activity

AmiGO ⬀

**Last version checked**
date: 14:01:2011 17:26
saved-by: rfoulger
auto-generated-by: OBO-Edit 2.0

**Last updated**
date: 08:10:2010 13:21
saved-by: dph
auto-generated-by: OBO-Edit 2.0

Gene Ontology Home ⬀
The contents of this box are automatically generated. You can help by adding information to the "Notes" ⬀

## Usage Notes                                    [edit]

## References                                     [edit]

See Help:References for how to manage references in GONUTS.

## Child Terms

This term has the following 4 child terms.

- [+] GO:0004714 - transmembrane receptor protein tyrosine kinase activity (13)
- [ ] GO:0004715 - non-membrane spanning protein tyrosine kinase activity
- [+] GO:0004716 - receptor signaling protein tyrosine kinase activity (1)
- [+] GO:0035400 - histone tyrosine kinase activity (1)

## Pages in category "GO:0004713 ! protein tyrosine kinase activity"

The following 200 pages are in this category, out of 732 total.

Show articles starting with: [ --All results-- ▾ ] [ Go ]

(previous 200) (next 200)

**C**
- CHICK:A0M8T9
- CHICK:A0SVH2
- CHICK:BTK

**C cont.**
- CHICK:Q90960
- CHICK:Q90961
- CHICK:Q90962

**F cont.**
- FB:Tk4
- FB:Tk6
- FB:tor

# Strategies

- Search for a keyword and browse the ontology for the right term
    - In GONUTS only search "Category" namespace if you get too many hits
    - Look at the parents, children, and relatives
    - Use Google, Wikipedia etc. to find alternative search terms
- Look at terms suggested by others for your protein
    - Computational with the IEA evidence code
    - Curators with TAS or IC
- Look at terms used for homologous proteins in model organisms
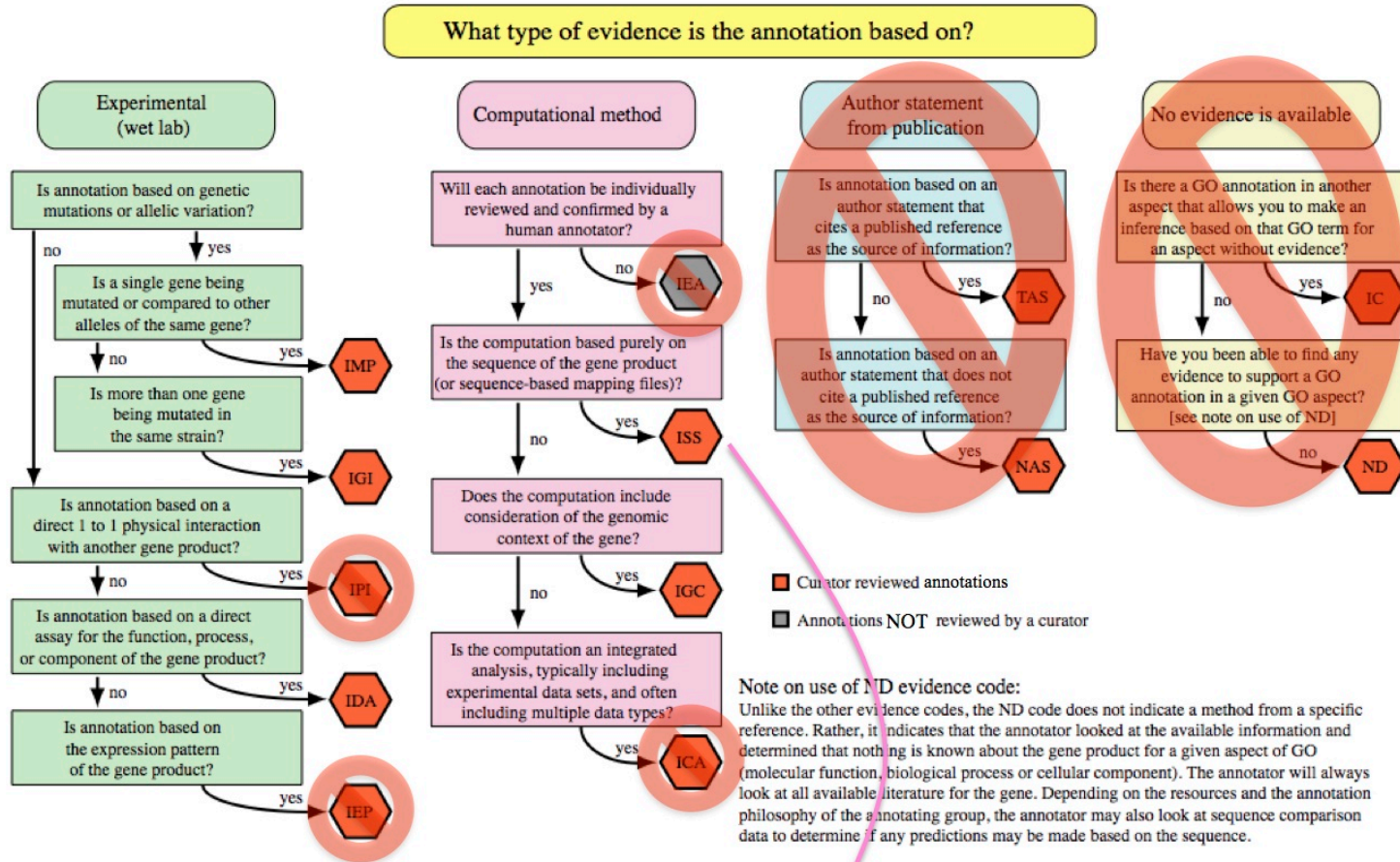
# Functional Annotation w/GO

- Annotation:  a note that is made while reading any form of text


- GO Annotation:  a database entry in a **specific format** that associates a **GO term** with a **gene product** made based on **evidence** in a peer-reviewed **paper**

# Evidence Codes for CACAO

- Evidence codes describe the type of work or analysis done by the authors
  - IDA: Inferred from Direct Assay
  - IMP: Inferred from Mutant Phenotype
    - NOT just for mutations! Includes inferred from inhibition in vivo by drugs, RNAi, etc.
  - IGI: Inferred from Genetic Interaction
  - ISO: Inferred from Sequence Orthology
  - ISA: Inferred from Sequence Alignment
  - ISM: Inferred from Sequence Model
  - IGC: Inferred from Genomic Context

- Expert biocurators get to use others, but we restrict them for CACAO. If it's not one of these 7, your annotation is incorrect!!!

- http://gowiki.tamu.edu/wiki/index.php/evidence_codes

# Decision Tree to Choose Evidence

What type of evidence is the annotation based on?



**ALLOWED CODES FOR ALL CACAO STUDENTS:**

1. **IDA: Inferred from Direct Assay**
2. **IMP: Inferred from Mutant Phenotype**
3. **IGI: Inferred from Genetic Interaction** - requires with/from field to be filled in
4. **ISO: Inferred from Sequence Orthology** - requires with/from field to be filled in
5. **ISA: Inferred from Sequence Alignment** - requires with/from field to be filled in
6. **ISM: Inferred from Sequence Model** - requires with/from field to be filled in
7. **IGC: Inferred from Genomic Context**

Use one of these three codes (ISO, ISA, ISM) if the Decision Tree points you to ISS

# Evidence Pull-Down Menu

# Some Evidence Types Require More Information

- With/from

- Evidence from sequence comparison
  - With the protein accession for the protein you are comparing to
    - That protein must have experimental annotation to the same GO term

- Evidence from computational analysis
  - With the reference for the analysis tool

- Evidence from genetic interaction
  - With the other gene(s) your protein is interacting with

# Evidence Codes for CACAO

- Picking the right evidence code is important
- Use the evidence code decision tree
  - http://gowiki.tamu.edu/wiki/images/3/32/CACAO_decisiontree.pdf
- Use the evidence code guidelines at the GO consortium website:
  - http://www.geneontology.org/GO.evidence.shtml
- Discuss!

# Note Required for CACAO

# Example Paper

http://www.ncbi.nlm.nih.gov/pubmed/8227000
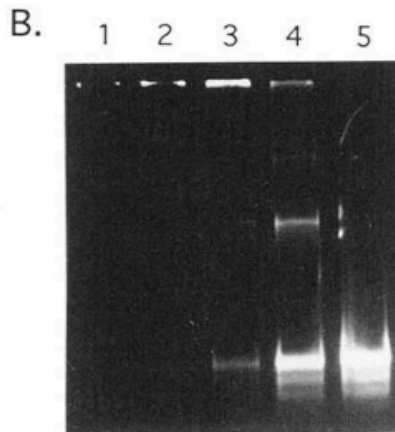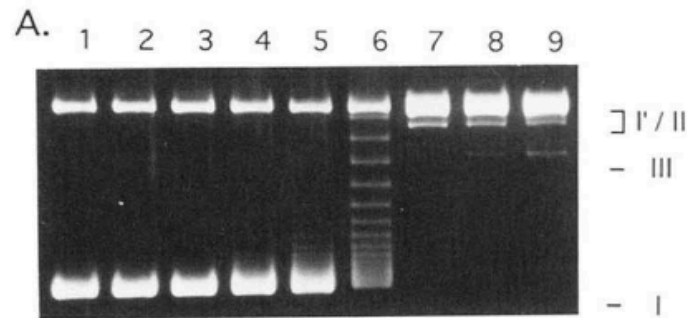=
http://www.jbc.org/content/268/32/24481.full.pdf

# Example Paper: What They Did

- Finding the proteins

- Do these tell us about the function?
  - Figure 1: sequenced ParC and Part of ParE
  - Figure 2: SDS page of purified proteins
  - Figure 3: Relaxation and decatenation activities of TopoIV
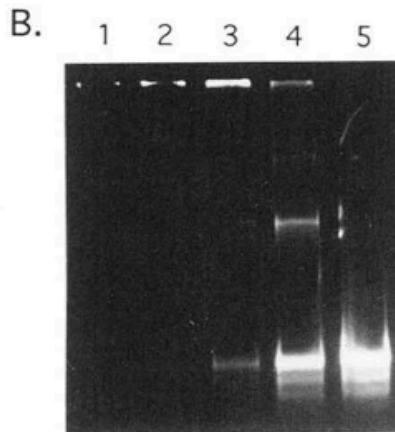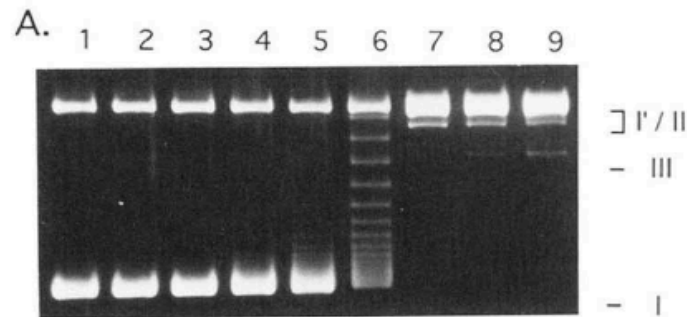  - …

# Example Paper: Figure 3



- Panel A: relaxation
- Panel B: decatenation
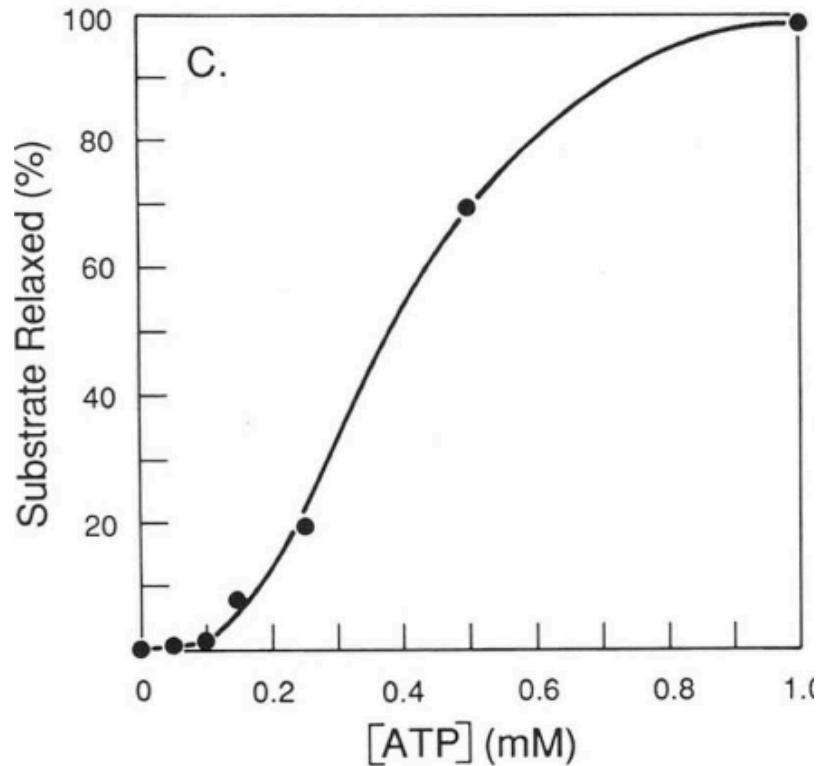- What do these mean?

# Example Paper: Figure 3



- Panel A: relaxation
- Panel B: decatenation
- What do these mean?

- Panel A shows GO: 0003916 ! DNA topoisomerase activity but does not show what kind
- Panel B shows GO: 0061505 ! DNA topoisomerase II activity

# Example Paper: Figure 4



- Shows ATP dependence: GO: 0003918 ! DNA topoisomerase type II (ATP-hydrolyzing) activity

# Example Paper: GO annotation for *E. coli* ParC



## TableEdit

### ECOLI:PARC

| | |
|---|---|
| **Qualifier** | [dropdown] |
| **GO ID** | GO:0003918 |
| **GO term name** | DNA topoisomerase type II (ATP-hydrolyzing) activity |
| **Reference** | PMID: 8227000 |
| **Evidence Code** | IDA: Inferred from Direct Assay |
| **with/from** | |
| **Aspect** | F |
| **Notes** | Topoisomerase assay in Fig 3. ATP dependent decatenation means it is a Type II from Fig 4 |
| **Status** | complete |

Public | Refresh | Save Row | Cancel