

CACAO training part I

Jim Hu and Suzi Aleksander

For UW Parkside

Fall 2014

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
 - Environmental responses
 - 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. Hofl[¶], 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

Contributed by Morris Goodman, December 30, 2003

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 representing transcripts of all or most h detected in human ACC samp 15,000, in gorilla and chimpan obtained from gene express expectation that the non-hum gorilla) should be more like humans. Instead, the chimpan human than like the gorilla; panzees are the sister group d biguous expression changes cesses and molecular functio 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 increased 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

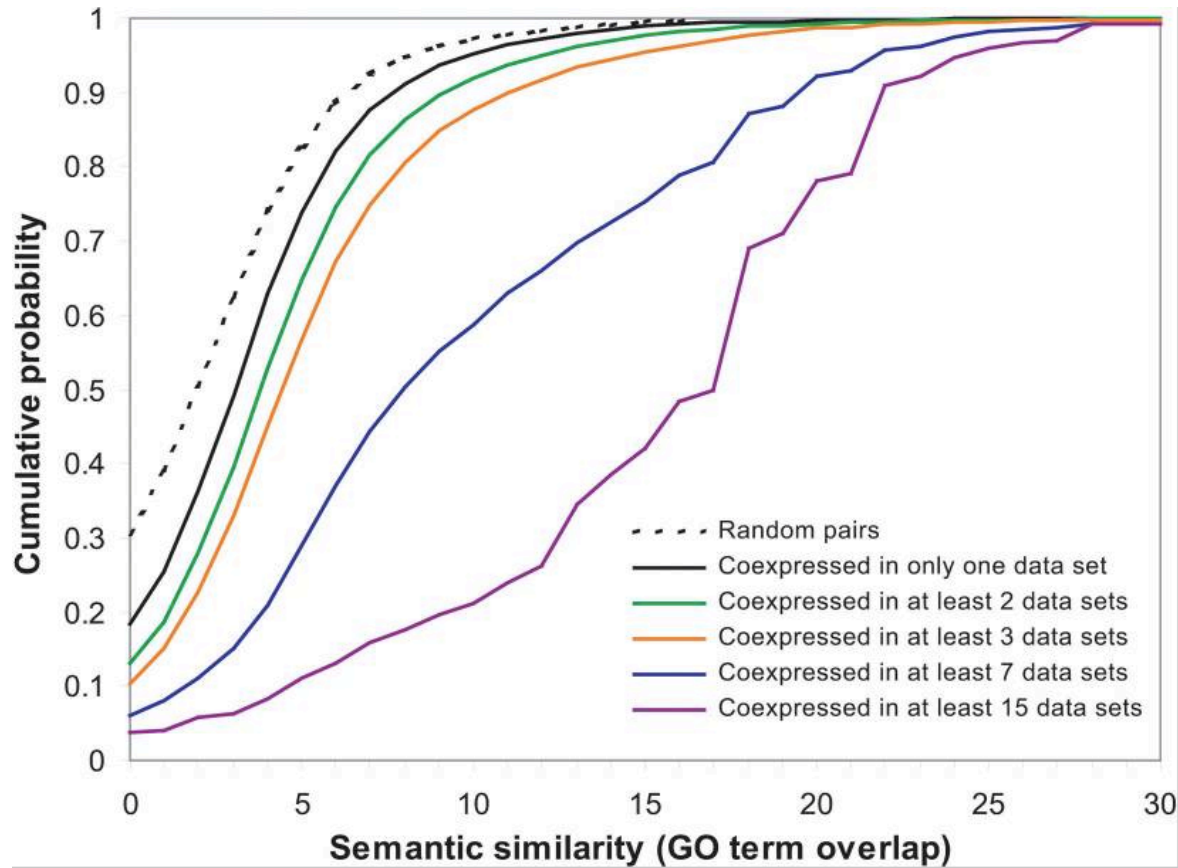
Among important changes in the ancestry of both humans and chimpanzees, but to a greater extent in humans, are the up-regulated expression profiles of aerobic energy metabolism genes and neuronal function-related genes, suggesting that increased neuronal activity required increased supplies of energy.

structured 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

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

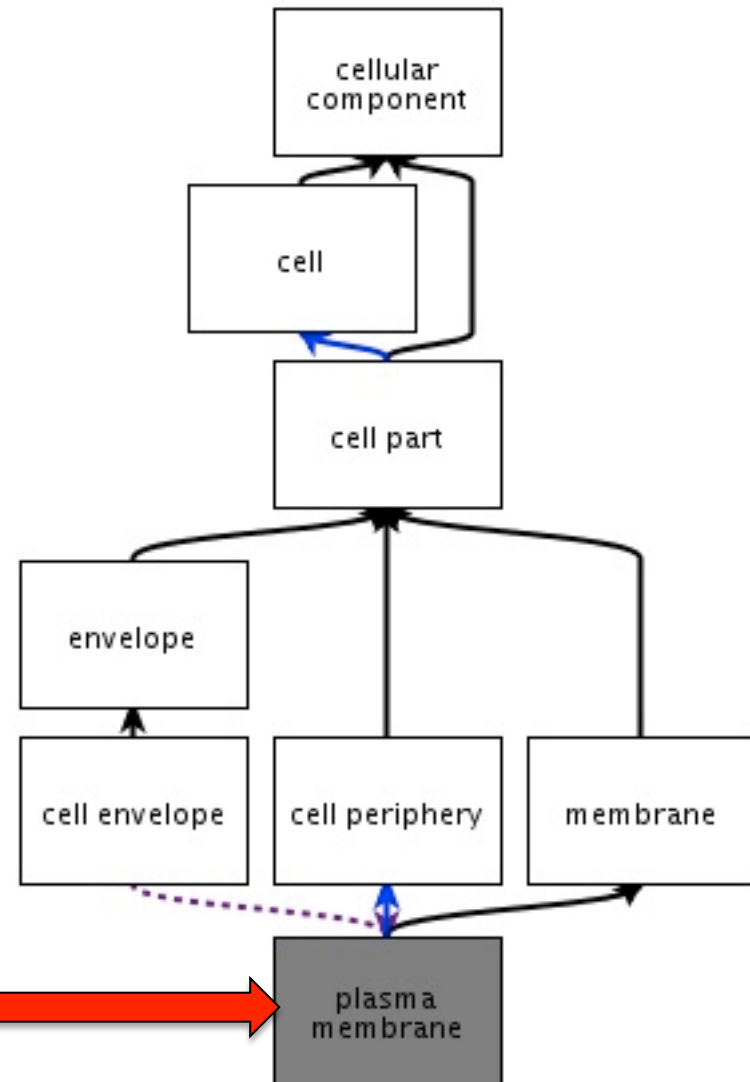
What do we mean by function?

- 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)

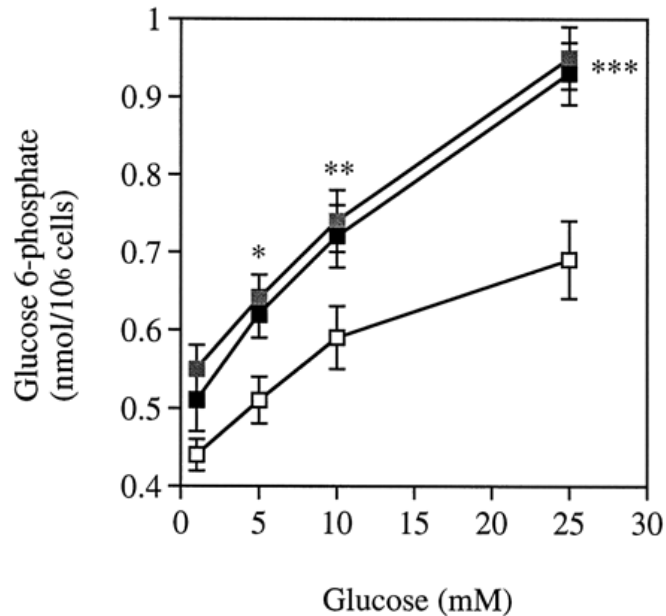
- 3 aspects (ontologies) :
 - Molecular Function
 - Biological Process
 - Cellular Component
- Controlled vocabulary
 - ID number for computers
 - Name and definition for humans
- Relationships



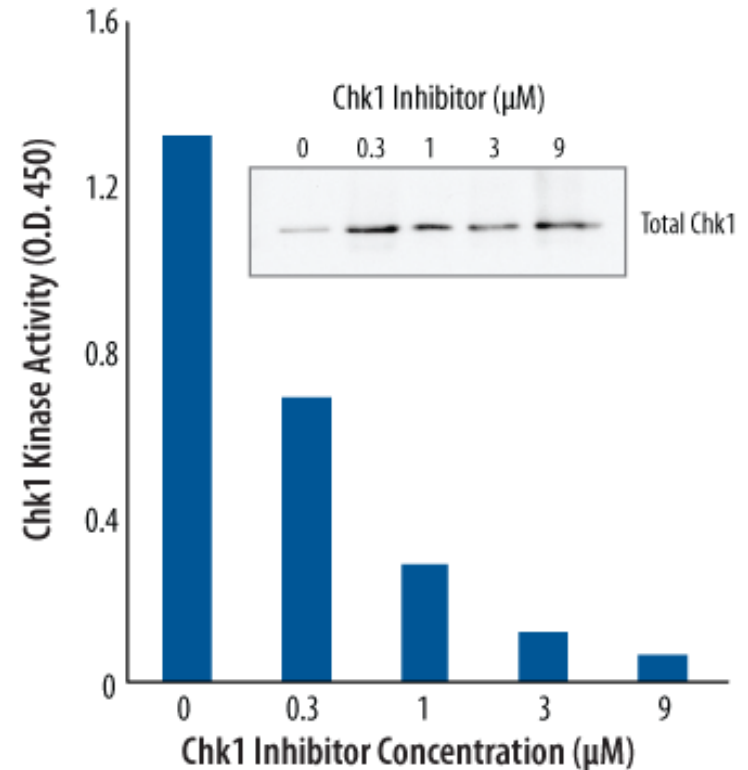
GO:0005886

Molecular Function

- activities = what a protein can do by itself



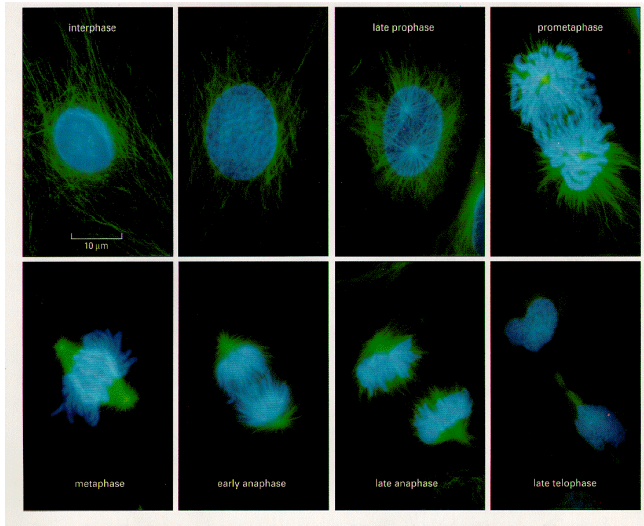
GO:0004347 hexokinase activity



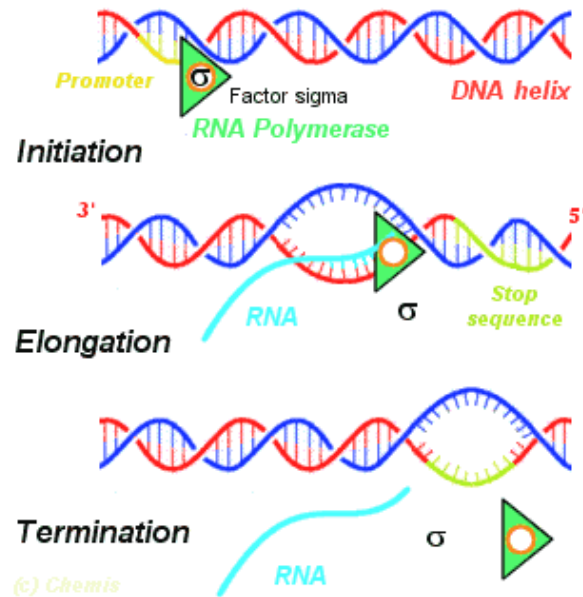
GO:0016301 Kinase activity

Biological Process

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



GO:0051301
cell division



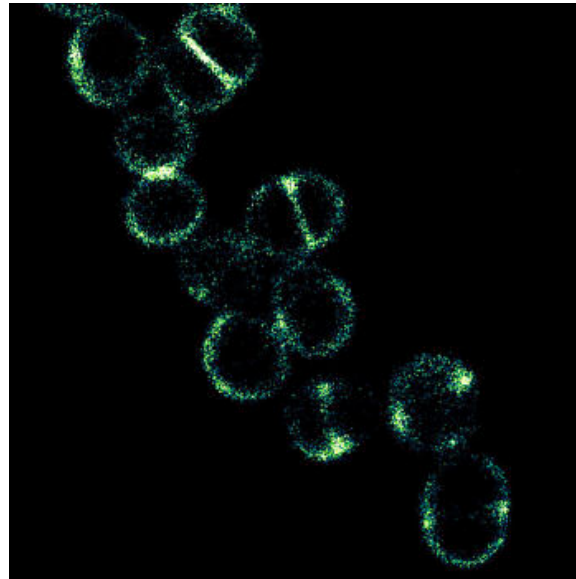
GO:0006351
transcription, DNA dependent

Cellular Component

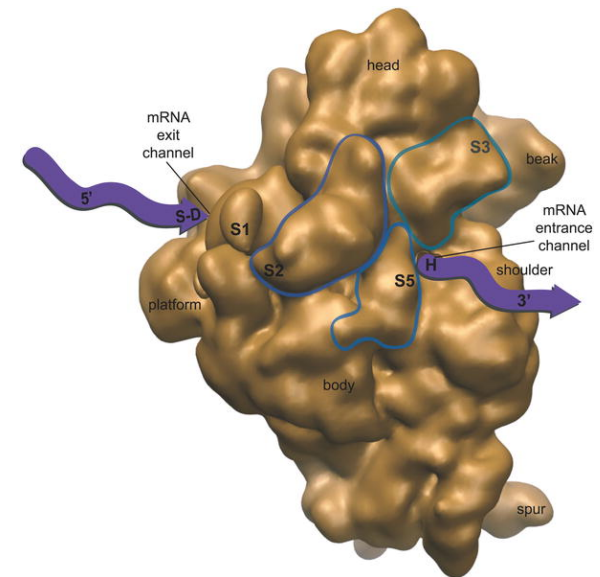
- where a gene product acts
 - Subcellular location
 - Multicomponent complex



GO:0005739
mitochondrion



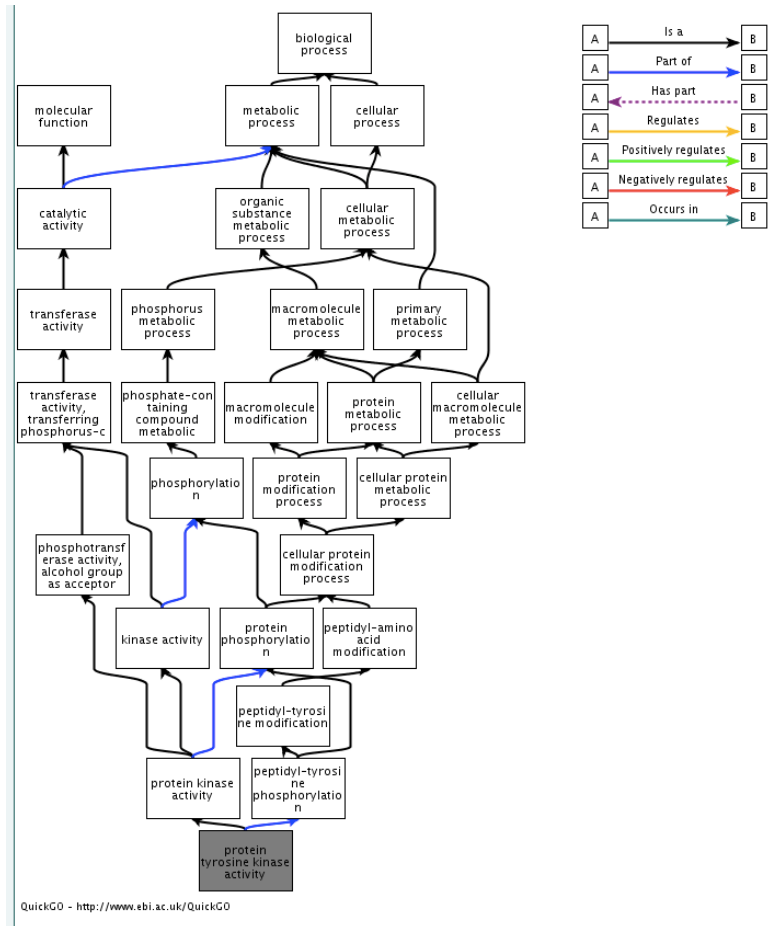
GO:0009274
peptidoglycan-based cell wall



GO:0005840
ribosome

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

Q annotation

an•no•ta•tion |,anə'tā SH ən|

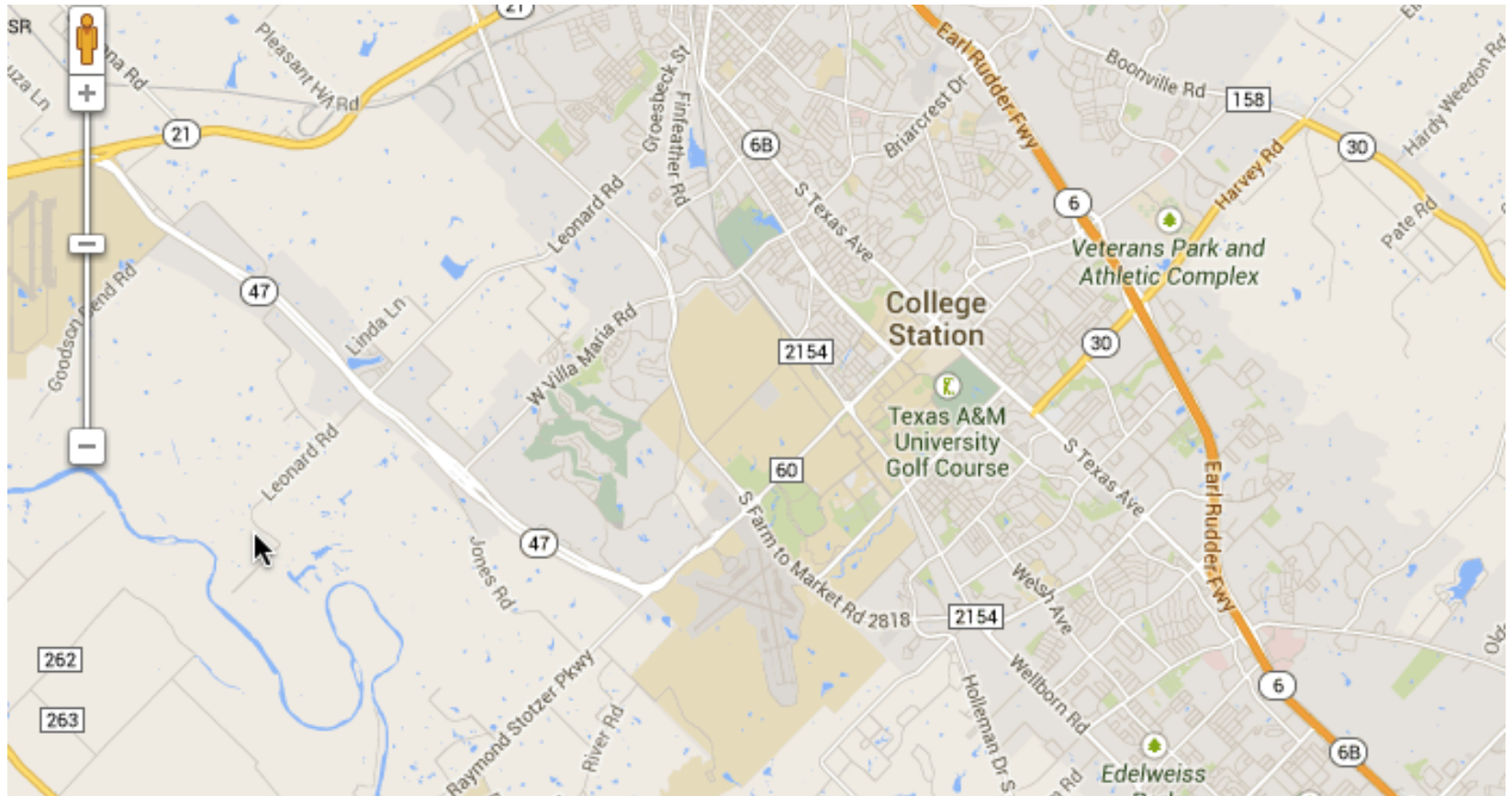
noun

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

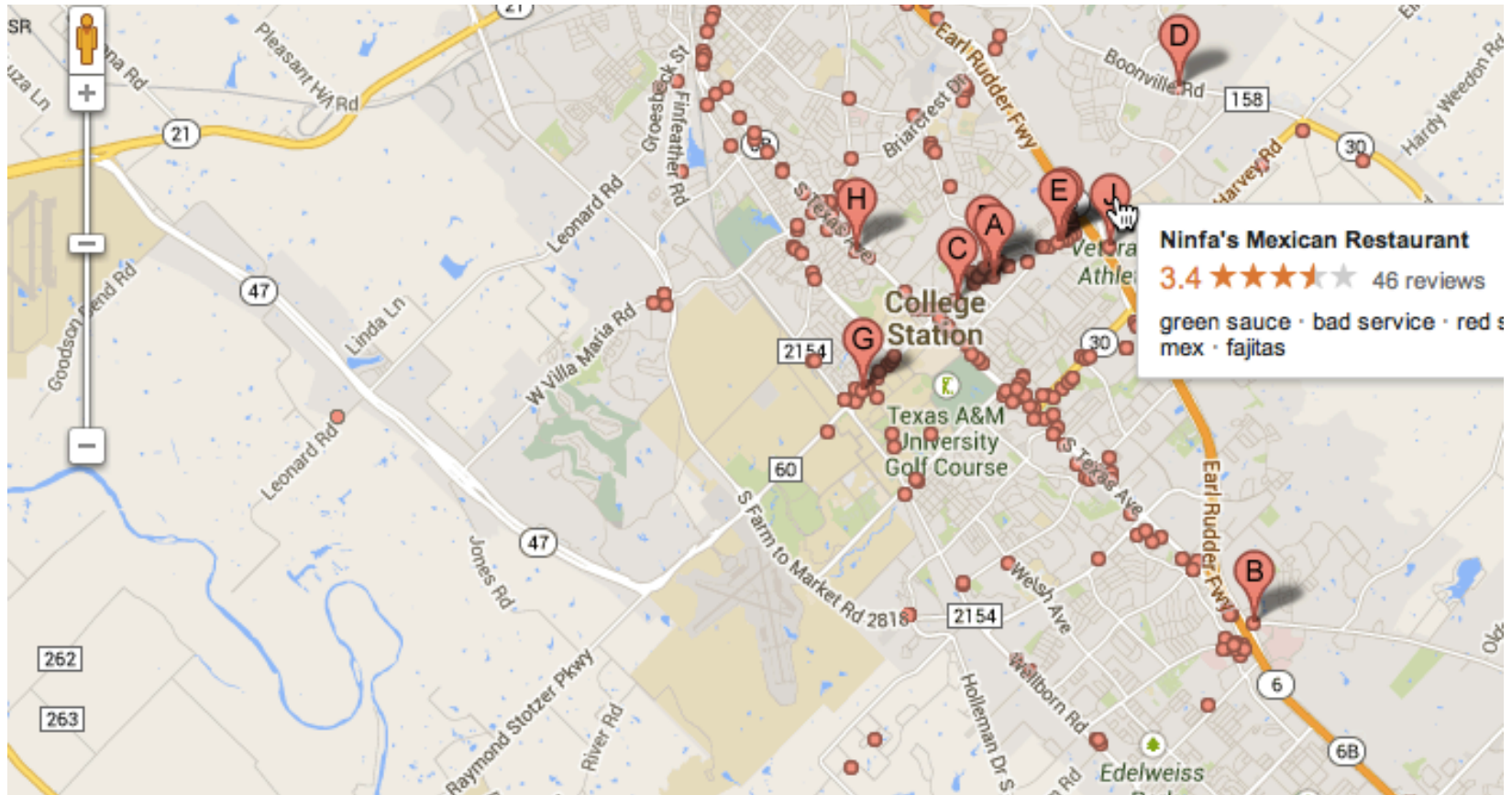
- the action of annotating a text or diagram : *annotation of prescribed texts.*

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?

Enter HAMLET. [Full Summary](#)

HAMLET

56 To be, or not to be: that is the question:
57 Whether 'tis nobler in the mind to suffer
58 The slings and arrows of outrageous fortune,
59 Or to take arms against a sea of troubles,
60 And by opposing end them? To die, to sleep—
61 No more—and by a sleep to say we end
62 The heart-ache and the thousand natural shocks
63 That flesh is heir to, 'tis a consummation
64 Devoutly to be wish'd. To die, to sleep;
65 To sleep: perchance to dream: ay, there's the rub;
66 For in that sleep of death what dreams may come
67 When we have shuffled off this mortal coil,
68 Must give us pause: there's the respect
69 That makes calamity of so long life;
70 For who would bear the whips and scorns of time,
71 The oppressor's wrong, the proud man's contumely,
72 The pangs of despised love, the law's delay,
73 The insolence of office and the spurns
74 That patient merit of the unworthy takes,
75 When he himself might his quietus make
76 With a bare bodkin? Who would fardels bear,
77 To grunt and sweat under a weary life,
78 But that the dread of something after death,
79 The undiscover'd country from whose bourn
80 No traveller returns, puzzles the will
81 And makes us rather bear those ills we have
82 Than fly to others that we know not of?
83 Thus conscience does make cowards of us all;
84 And thus the native hue of resolution

Notes to Hamlet...ct 3, Scene 1

57. **suffer:** endure patiently.

58. **slings:** *i.e.*, projectiles launched from slings.

60. **To die, to sleep— / No more—:** This sequence puzzles me. "To sleep" seems to be a comforting way of describing what it means "to die," but "No more" could mean "to dream no more"; remember that Hamlet said to Rosencrantz and Guildenstern, "[I could be bounded in a nutshell and count / myself a king of infinite space, were it not that I / have bad dreams.](#)" On the other hand, "No more" could be all-encompassing: no more "slings and arrows"; no more "sea of troubles"; no more questions about what would be "nobler in the mind."

63. **consummation:** completion, end.

65. **rub:** *i.e.*, obstacle, catch. The term comes from the game Americans know as lawn bowling, in which "A rub is some fault in the surface of the green that stops a bowl or diverts it from its intended direction" ([World Wide Words: Michael Quinton writes on International English from a British Viewpoint](#)).

67. **shuffled off:** sloughed, cast off. **this mortal coil:** the turmoil of this mortal life.

68. **respect:** consideration.

69. **of so long life:** so long-lived.

70. **bear the whips and scorns of time:** *i.e.*, endure the punishments and insults that always come with the passage of time.

<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?

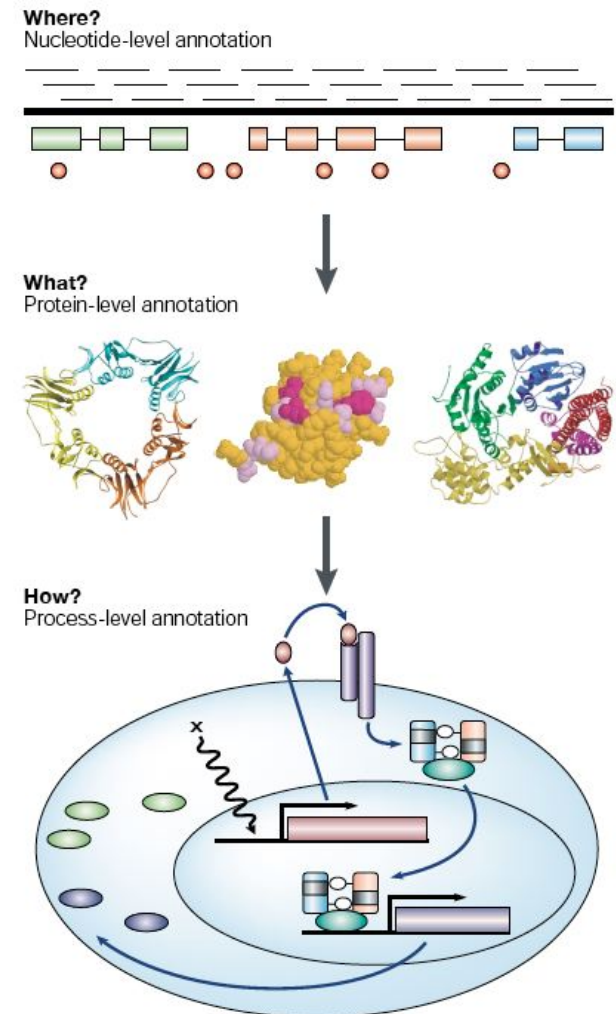


Figure 1 | **The three layers of genome annotation: where, what and how?**

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?

Journal home > Archive > Letters to Nature > Abstract

Letters to Nature

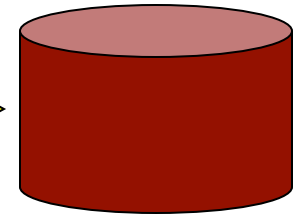
Nature 425, 628-633 (9 October 2003) | doi: 10.1038/nature02030

Basal body dysfunction is a likely cause of pleiotropic Bardet-Biedl syndrome

Stephen J. Ansley^{1,2,3}, Jose L. Badano^{1,2}, Oliver E. Blacque^{3,4,5}, Josephine Hill³, Bethan E. Hoskins^{1,2}, Carmen C. Leitch¹, Jun Chul Kim³, Alison J. Ross³, Erica R. Eichers⁵, Tanya M. Teslovich¹, Allan K. Mah³, Robert C. Johnson³, John C. Cavender², Richard Alan Lewis^{3,6}, Michel R. Leroux³, Philip L. Beales³ and Nicholas Katsanis^{1,2}

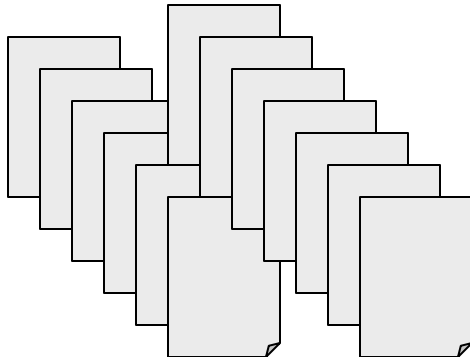
Bardet-Biedl syndrome (BBS) is a genetically heterogeneous disorder characterized primarily by retinal dystrophy, obesity, polydactyly, renal malformations and learning disabilities. Although five BBS genes have been cloned^{1, 2, 3, 4, 5, 6}, the molecular basis of this syndrome remains elusive. Here we show that BBS is probably caused by a defect at the basal body of ciliated cells. We have cloned a new BBS gene, *BBS8*, which encodes a protein with a prokaryotic domain, *pilF*, involved in pilus formation and twitching motility. In one family, a homozygous null *BBS8* mutation leads to BBS with randomization of left-right body axis symmetry, a known defect of the nodal cilium. We have also found that *BBS8* localizes specifically to ciliated structures, such as the connecting cilium of the retina and columnar epithelial cells in the lung. In cells, *BBS8* localizes to centrosomes and basal bodies and interacts with PCMT1, a protein probably involved in cillogenesis. Finally, we demonstrate that all available *Caenorhabditis elegans* BBS homologues are expressed exclusively in ciliated neurons, and contain regulatory elements for RFX, a transcription factor that modulates the expression of genes associated with cillogenesis and intraflagellar transport.

• Top



Database

Literature



Datasets

Biocurators
(rate limiting)

Databases need help!

- >21 million peer-reviewed articles in PubMed
- Many millions of proteins recorded in UniProt

The screenshot shows the UniProtKB search interface. At the top, the UniProt logo and 'UniProtKB' are visible. Below this are navigation tabs for 'Search', 'Blast', 'Align', 'Retrieve', and 'ID Mapping *'. The 'Search' tab is active. Under 'Search in', a dropdown menu is set to 'Protein Knowledgebase (UniProtKB)'. The 'Query' field contains the text 'human'. To the right of the query field are buttons for 'Search', 'Advanced Search »', and 'Clear'. Below the search bar, the results summary reads: '1 - 25 of 1,093,299 results for human in UniProtKB sorted by score descending'. There are also links to 'Browse by taxonomy, keyword, gene ontology, enzyme class or pathway' and a filter for 'Reduce sequence redundancy to 100%, 90% or 50%'. At the bottom, the 'Results' section has a 'Customize' button. Below this, there are two filter options: 'Show only reviewed (45,159) ★ UniProtKB/Swiss-Prot) or unreviewed (1,048,140) ★ (UniProtKB/TrEMBL) entries'. The numbers '45,159' and '1,048,140' are highlighted with red boxes.

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
- 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
 - Start with a GO term
- You can start with a protein
 - Find papers about the protein
- Either way, don't get stuck on what you started with
 - Your first paper may not have **experiments** about function
 - 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 **paper**

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 should 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

NCBI Resources ▾ How To ▾

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

PubMed Hu AND McIntosh

 RSS [Save search](#) [Limits](#) [Advanced](#)

[Display Settings:](#) Summary, 20 per page, Sorted by Recently Added

[Send to:](#)

Results: 10

[GONUTS: the Gene Ontology Normal Usage Tracking System.](#)

1. Renfro DP, **McIntosh** BK, Venkatraman A, Siegele DA, **Hu** JC.

Nucleic Acids Res. 2012 Jan;40(1):D1262-9. Epub 2011 Nov 22.

PMID: 22110029

[Related citations](#)

22110029

Pubmed record

Display Settings: Abstract

Send to:

Mol Syst Biol. 2012 May 8;8:581. doi: 10.1038/msb.2012.13.

Prediction and identification of sequences coding for orphan enzymes using genomic and metagenomic neighbours.

Yamada T¹, Waller AS, Raes J, Zelezniak A, Perchat N, Perret A, Salanoubat M, Patil KR, Weissenbach J, Bork P.

Author information

Abstract

Despite the current wealth of sequencing data, one-third of all biochemically characterized metabolic enzymes lack a corresponding gene or protein sequence, and as such can be considered orphan enzymes. They represent a major gap between our molecular and biochemical knowledge, and consequently are not amenable to modern systemic analyses. As 555 of these orphan enzymes have metabolic pathway neighbours, we developed a global framework that utilizes the pathway and (meta)genomic neighbour information to assign candidate sequences to orphan enzymes. For 131 orphan enzymes (37% of those for which (meta)genomic neighbours are available), we associate sequences to them using scoring parameters with an estimated accuracy of 70%, implying functional annotation of 16,345 gene sequences in numerous (meta)genomes. As a case in point, two of these candidate sequences were experimentally validated to encode the predicted activity. In addition, we augmented the currently available genome-scale metabolic models with these new sequence-function associations and were able to expand the models by on average 8%, with a considerable change in the flux connectivity patterns and improved essentiality prediction.

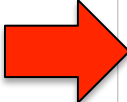
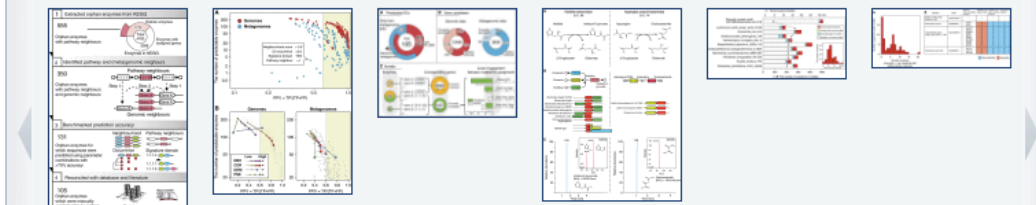
PMID: 22569339 [PubMed - indexed for MEDLINE] PMID: PMC3377989 [Free PMC Article](#)



Full text links



Images from this publication. [See all images \(6\)](#) [Free text](#)



Getting the full text

US National Library of Medicine
National Institutes of Health

Advanced Help

Display Settings: Abstract Send to:

Mol Syst Biol. 2012 May 8;8:581. doi: 10.1038/msb.2012.13.

Prediction and identification of sequences coding for orphan enzymes using genomic and metagenomic neighbours.

Yamada T¹, Waller AS, Raes J, Zeleznik A, Perchat N, Perret A, Salanoubat M, Patil KR, Weissenbach J, Bork P.

[+](#) Author information

Abstract

Despite the current wealth of sequencing data, one-third of all biochemically characterized metabolic enzymes lack a corresponding gene or protein sequence, and as such can be considered orphan enzymes. They represent a major gap between our molecular and biochemical knowledge, and consequently are not amenable to modern systemic analyses. As 555 of these orphan enzymes have metabolic pathway neighbours, we developed a global framework that utilizes the pathway and (meta)genomic neighbour information to assign candidate sequences to orphan enzymes. For 131 orphan enzymes (37% of those for which (meta)genomic neighbours are available), we associate sequences to them using scoring parameters with an estimated accuracy of 70%, implying functional annotation of 16,345 gene sequences in numerous (meta)genomes. As a case in point, two of these candidate sequences were experimentally validated to encode the predicted activity. In addition, we augmented the currently available genome-scale metabolic models with these new sequence-function associations and were able to expand the models by on average 8%, with a considerable change in the flux connectivity patterns and improved essentiality prediction.

PMID: 22569339 [PubMed - indexed for MEDLINE] PMID: PMC3377989 [Free PMC Article](#)

[f](#) [t](#) [+](#)

Images from this publication. [See all images \(6\)](#) [Free text](#)



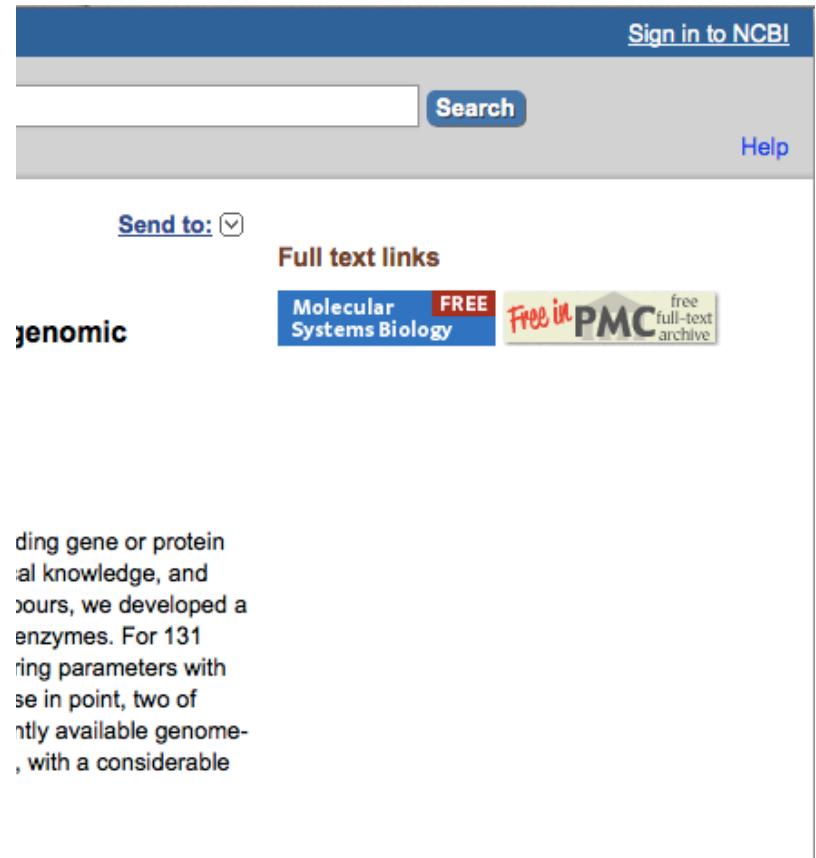
Full text links



- The abstract is not enough
 - 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



The screenshot shows the top navigation bar of a PubMed page. It includes a blue header with the text "Sign in to NCBI" on the right. Below the header is a search bar with a "Search" button. A "Help" link is visible in the bottom right corner of the header area. Below the search bar, there is a "Send to:" dropdown menu and a "Full text links" section. The "Full text links" section contains two buttons: "Molecular Systems Biology" with a "FREE" badge, and "Free in PMC" with a "free full-text archive" badge. Below these buttons, there is a snippet of text from a paper, which is partially cut off and reads: "ding gene or protein al knowledge, and ours, we developed a enzymes. For 131 ring parameters with se in point, two of ntly available genome-, with a considerable".

Alternative path: Start w/Full Text

NCBI Resources How To

PMC
US National Library of Medicine
National Institutes of Health

PMC gene function AND enzyme
Save search Journal List Limits Advanced

Display Settings: Summary, 20 per page, Sorted by Default order Send to:

Results: 1 to 20 of 596771 << First < Prev Page 1 of 29839 Next > Last >>

[Recent advances in 2D and 3D in vitro systems using primary hepatocytes, alternative hepatocyte sources and non-parenchymal liver cells and their use in investigating mechanisms of hepatotoxicity, cell signaling and ADME](#)

1. **Patricio Godoy, Nicola J. Hewitt, Ute Albrecht, Melvin E. Andersen, Nariman Ansari, Sudin Bhattacharya, Johannes Georg Bode, Jennifer Bolleyn, Christoph Borner, Jan Böttger, Albert Braeuning, Robert A. Budinsky, Britta Burkhardt, Neil R. Cameron, Giovanni Camussi, Chong-Su Cho, Yun-Jaie Choi, J. Craig Rowlands, Uta Dahmen, Georg Damm, Olaf Dirsch, María Teresa Donato, Jian Dong, Steven Dooley, Dirk Drasdo, Rowena Eakins, Karine Sá Ferreira, Valentina Fonsato, Joanna Fraczek, Rolf Gebhardt, Andrew Gibson, Matthias Glanemann, Chris E. P. Goldring, María José Gómez-Lechón, Geny M. M. Groothuis, Lena Gustavsson, Christelle Guyot, David Hallifax, Seddik Hammad, Adam Hayward, Dieter Häussinger, Claus Hellerbrand, Philip Hewitt, Stefan Hoehme, Hermann-Georg Holzhütter, J. Brian Houston, Jens Hrach, Kiyomi Ito, Hartmut Jaeschke, Verena Keitel, Jens M. Kelm, B. Kevin Park, Claus Kordes, Gerd A. Kullak-Ublick, Edward L. LeCluyse, Peng Lu, Jennifer Luebke-Wheeler, Anna Lutz, Daniel J. Maltman, Madlen Matz-Soja, Patrick McMullen, Irmgard Merfort, Simon Messner, Christoph Meyer, Jessica Mwinyi, Dean J. Naisbitt, Andreas K. Nussler, Peter Olinga, Francesco Pampaloni, Jingbo Pi, Linda Pluta, Stefan A. Przyborski, Anup Ramachandran, Vera Rogiers, Cliff Rowe, Celine Schelcher, Kathrin Schmich, Michael Schwarz, Bijay Singh, Ernst H. K. Stelzer, Bruno Stieger, Regina Stöber, Yuichi Sugiyama, Ciro Tetta, Wolfgang E. Thasler, Tamara Vanhaecke, Mathieu Vinken, Thomas S. Weiss, Agata Widera, Courtney G. Woods, Jinghai James Xu, Kathy M. Yarborough, Jan G. Hengstler**
Arch Toxicol. 2013; 87(8): 1315–1530. Published online 2013 August 23. doi: 10.1007/s00204-013-1078-5
PMCID: PMC3753504
[Article](#) [PubReader](#) [PDF-9.1M](#) [Supplementary Material](#)

[Phosphoinositides: Tiny Lipids With Giant Impact on Cell Regulation](#)

2. **Tamas Balla**
Physiol Rev. 2013 July; 93(3): 1019–1137. doi: 10.1152/physrev.00028.2012
PMCID: PMC3962547
[Article](#) [PubReader](#)

Alternative path: Start w/Full Text

NCBI Resources How To

PMC US National Library of Medicine National Institutes of Health

Search

Limits Advanced Journal list

Journal List > Mol Syst Biol > v.8; 2012 > PMC3377989

PubReader for full text click here to

molecular systems biology Setting standards in Systems Biology

Mol Syst Biol. 2012; 8: 581. PMID: PMC3377989
Published online May 8, 2012. doi: [10.1038/msb.2012.13](https://doi.org/10.1038/msb.2012.13)

Prediction and identification of sequences coding for orphan enzymes using genomic and metagenomic neighbours

[Takuji Yamada](#)¹, [Alison S Waller](#)¹, [Jeroen Raes](#)^{2,3}, [Aleksandra Salanoubat](#)^{5,6,7}, [Kiran R Patil](#)¹, [Jean Weissenbach](#)^{5,6,7} and

[Author information](#) ▶ [Article notes](#) ▶ [Copyright and License information](#)

This article has been [cited by](#) other articles in PMC.

Abstract

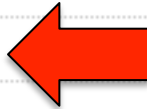
Despite the current wealth of sequencing data, one-third of the genes in prokaryotic genomes lack a corresponding gene or protein sequence.

Formats:
Article | [PubReader](#) | [ePub \(beta\)](#) |

Related citations in PubMed
The CanOE strategy: integrating genomic data across multiple prokaryote genomes to

Links

- MedGen
- Protein
- PubMed
- Taxonomy
- Taxonomy Tree



Beware!

- Good science \neq 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}

Received for publication, September 30, 2011, and in revised form, November 29, 2011 Published, JBC Papers in Press, December 6, 2011, DOI 10.1074/jbc.M111.309369

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 \neq 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*}

¹Department of Biochemistry, University of Washington and HHMI, Seattle, Washington 98195

²Arzeda Corp., Seattle, Washington 98102

³Department of Biological Chemistry, UCLA, Los Angeles, California 90095

⁴Laboratory of Organic Chemistry, ETH Zurich, 8093 Zurich, Switzerland

⁵Department of Biochemistry, Stanford University, Stanford, California 94305

Beware!

- Good science \neq good for annotation

Cell Stem Cell

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,*}

¹Stem 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

²Ludwig Boltzmann Institute for Clinical and Experimental Traumatology, Austrian Cluster for Tissue Regeneration, Vienna 1200, Austria

Beware!

- Good science \neq good for annotation

10624 • The Journal of Neuroscience, August 11, 2010 • 30(32):10624–10638

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 ¹Cell Biology, ²Human Genetics, ³Pathology and Laboratory Medicine, ⁴Biochemistry, ⁵Pediatrics, and ⁶Neurology, 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

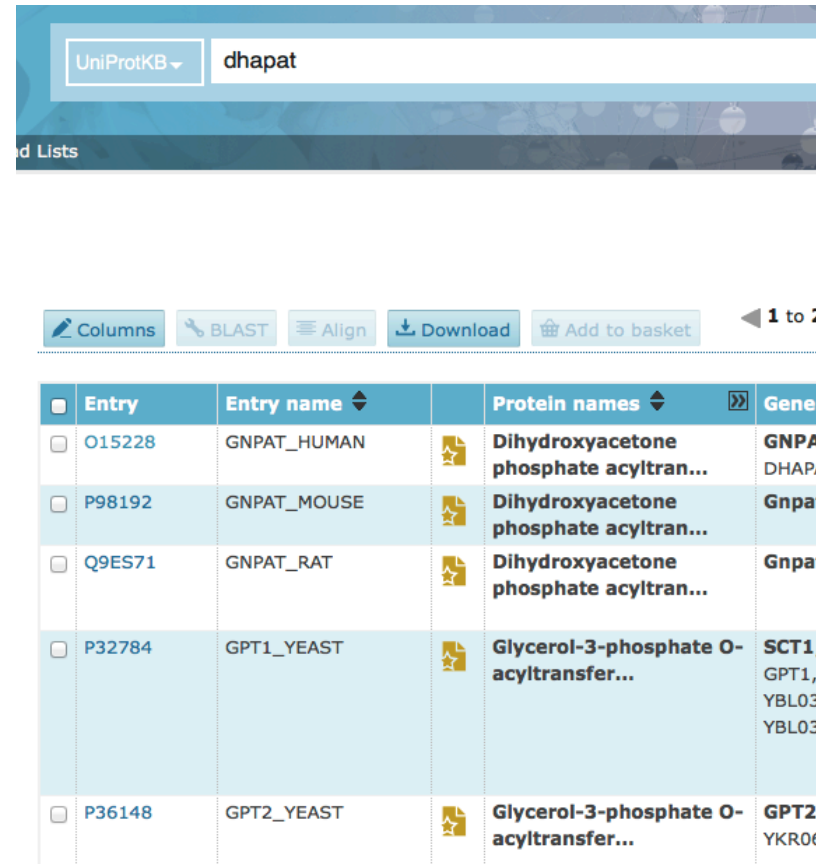
UniProt (<http://www.uniprot.org>)

- 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






The screenshot shows the UniProt website homepage. At the top, there is a search bar with 'UniProtKB' selected and 'DHAPAT' entered. Below the search bar, there are links for 'BLAST', 'Align', 'Upload Lists', and 'Help'. A yellow banner welcomes users to the new website. Below the banner, there is a navigation menu with 'UniProtKB', 'UniRef', 'UniParc', and 'Proteomes'. The 'UniProtKB' section shows 'Swiss-Prot (546,238)' and 'TrEMBL (82,126,897)'. The 'Supporting data' section includes 'Literature citations', 'Taxonomy', 'Subcellular locations', 'Cross-ref. databases', 'Diseases', and 'Keywords'. A 'News' section on the right features a tweet about Ubiquitin and a link to 'News archive'. At the bottom, there are links for 'Getting started', 'UniProt data', and 'Protein spotlight'.

UniProt search results

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



The screenshot shows the UniProt search interface. At the top, there is a search bar with 'UniProtKB' selected and the search term 'dhapat'. Below the search bar, there are several action buttons: 'Columns', 'BLAST', 'Align', 'Download', and 'Add to basket'. A pagination indicator shows '1 to 2'. The main content is a table with the following columns: 'Entry', 'Entry name', 'Protein names', and 'Gene'. The table contains five rows of search results.

<input type="checkbox"/>	Entry	Entry name		Protein names	Gene
<input type="checkbox"/>	O15228	GNPAT_HUMAN		Dihydroxyacetone phosphate acyltran...	GNPA DHAP
<input type="checkbox"/>	P98192	GNPAT_MOUSE		Dihydroxyacetone phosphate acyltran...	Gnpa
<input type="checkbox"/>	Q9ES71	GNPAT_RAT		Dihydroxyacetone phosphate acyltran...	Gnpa
<input type="checkbox"/>	P32784	GPT1_YEAST		Glycerol-3-phosphate O-acyltransfer...	SCT1 GPT1, YBL03 YBL03
<input type="checkbox"/>	P36148	GPT2_YEAST		Glycerol-3-phosphate O-acyltransfer...	GPT2 YKR06

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



The screenshot shows the UniProt record for GNPAT_HUMAN. At the top, there is a search bar with 'UniProtKB' selected. Below the search bar, the protein name 'GNPAT_HUMAN' is displayed in orange. The main title is 'Dihydroxyacetone phosphate acyltransferase'. Below this, the protein is identified as 'GNPAT, DAPAT, DHAPAT' and 'Homo sapiens (Human)'. A star icon indicates it is a 'Reviewed' entry with a score of 5.0, and a note states 'Experimental evidence at protein level'. A navigation bar contains buttons for 'BLAST', 'Align', 'Format', 'Add to basket', 'History', and 'Comments'. The 'Function' section is highlighted with an orange underline and contains three sub-sections: 'Catalytic activity' (Acyl-CoA + glycerone phosphate = CoA + acylglycerone phosphate), 'Pathway' (Membrane lipid metabolism; glycerophospholipid metabolism), and 'GO - Molecular function' (glycerone-phosphate O-acyltransferase activity, palmitoyl-CoA hydrolase activity, receptor binding). The 'GO - Biological process' section includes 'cellular lipid metabolic process' and 'cerebellum morphogene'.

UniProtKB

d Lists

GNPAT_HUMAN

Dihydroxyacetone phosphate acyltransferase

GNPAT, DAPAT, DHAPAT

Homo sapiens (Human)

Reviewed - 5.0 - Experimental evidence at protein levelⁱ

BLAST Align Format Add to basket History Comments

Functionⁱ

Catalytic activityⁱ
Acyl-CoA + glycerone phosphate = CoA + acylglycerone phosphate.

Pathwayⁱ
Membrane lipid metabolism; glycerophospholipid metabolism.

GO - Molecular functionⁱ

- ▶ glycerone-phosphate O-acyltransferase activity Source: UniProtKB
- ▶ palmitoyl-CoA hydrolase activity Source: UniProtKB
- ▶ receptor binding Source

GO - Biological processⁱ

- ▶ cellular lipid metabolic process Source: Reactome
- ▶ cerebellum morphogene

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

LAMB:VLYS - GONUTS

http://gowiki.tamu.edu/wiki/index.php/LAMB:VLYS

The Spring 2012 season of CACAO has started!

LAMB:VLYS

Species (Taxon ID) *Enterobacteria phage lambda (Bacteriophage lambda)*. ([1] [edit](#))

Gene Name(s) S

Protein Name(s) Holin
gpS protein Lysis protein S Lysis inhibitor

External Links

EMBL	J02459 M14035
PIR	H94164
RefSeq	NP_040644.1 YP_001551775.1
TCDB	1.E.2.1.1
GeneID	2703479 5740919
GenomeReviews	J02459_GR
ProtClusTDB	CLSP2343227
GO	GO:0020002 GO:0016021 GO:0016998 GO:0019835
InterPro	IPR006481
Pfam	PF06708
TIGRFAMs	TIGR01594

Annotations

Qualifier	GO ID	GO term name	Reference	Evidence Code	with/from	Aspect	Notes	Status
	GO:0016020	membrane	GO_REF:0000004	IEA: Inferred from Electronic Annotation	SP_KW:KW-0472	C	Seeded From UniProt	
	GO:0033644	host cell membrane	GO_REF:0000004	IEA: Inferred from Electronic Annotation	SP_KW:KW-1043	C	Seeded From UniProt	

[edit table](#)

Notes

edit table

Annotations


edit table

Contents [hide]
1 Annotations
2 Notes
3 References

[edit]

[edit]

Entering/editing annotations



special page

The Spring 2012 s

TableEdit

LAMBD:VLYS

Qualifier	<input type="text"/>
GO ID	<input type="text"/>
GO term name	
Reference	<input type="text"/>
Evidence Code	<input type="text"/>
with/from	
Aspect	
Notes	<input type="text"/>
Status	Missing: GO ID, evidence, reference
<input type="text"/> Public <input type="button" value="Refresh"/> <input type="button" value="Save Row"/> <input type="button" value="Cancel"/>	

Public rows can be edited or deleted by any user who can edit
Private rows can be edited or deleted by their creator, or by admins

navigation

- Main Page
- Enter GO at the top
- Help
- Report Bug
- Update log
- Annotation Jamborees
- Recent changes
- Create New Gene Page
- Login / Create Account

cacao

- Links about CACAO
- Fall 2011

journal clubs

- Journal Clubs
- Create new literature page

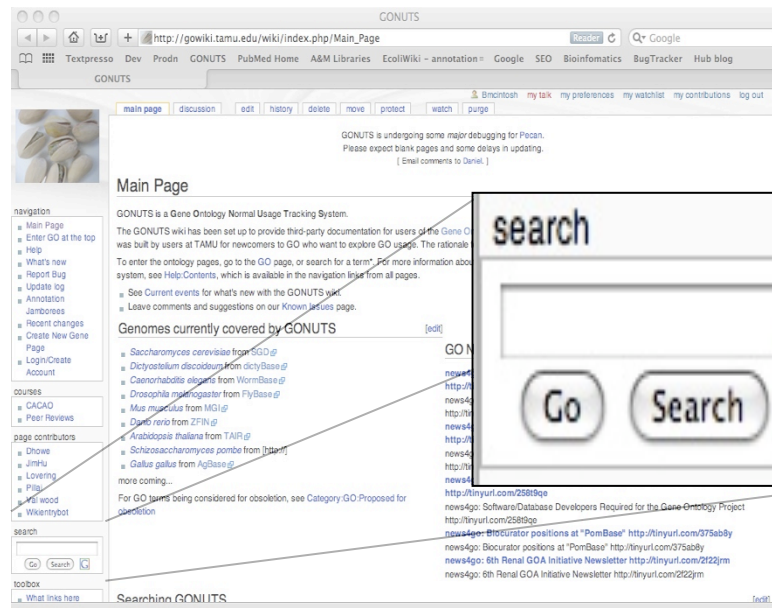
search

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>
- AmiGO: <http://amigo.geneontology.org>



The screenshot shows the GONUTS website interface. At the top, there is a navigation bar with links for 'main page', 'discussion', 'edit', 'history', 'delete', 'move', 'protect', and 'watch'. Below this, a message states: 'GONUTS is undergoing some major debugging for Pecan. Please expect blank pages and some delays in updating. [Email comments to Daniel.]'. The main content area is titled 'Main Page' and describes GONUTS as a 'Gene Ontology Normal Usage Tracking System'. It provides instructions on how to use the system and lists 'Genomes currently covered by GONUTS', including *Saccharomyces cerevisiae*, *Drosophila melanogaster*, and *Mus musculus*. A search box is located in the bottom right corner of the page, with a red arrow pointing to it. The search box contains the text 'search' and has buttons for 'Go', 'Search', and a 'G' icon.

[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.]



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

toolbox

- [What links here](#)
- [Related changes](#)
- [Upload file](#)
- [Special pages](#)
- [Printable version](#)
- [Permanent link](#)

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: rfulger

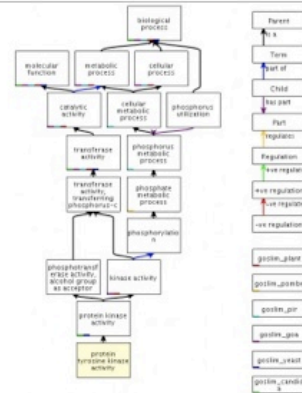
auto-generated-by: OBO-Edit 2.0

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" [↗](#)

GO:0004713 - http://www.geneontology.org/ontology/GO:0004713

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:

(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](#)

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



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

[Gene Ontology Home](#)

The contents of this box are automatically generated. You can help by adding information to the ["Notes"](#)

[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: rfulger

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:

(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


GO:0004713 - http://www.ebi.ac.uk/GO/GO:0004713

Strategies

- Search for a keyword and browse the ontology for the right term
 - In GONUTS only search categories 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

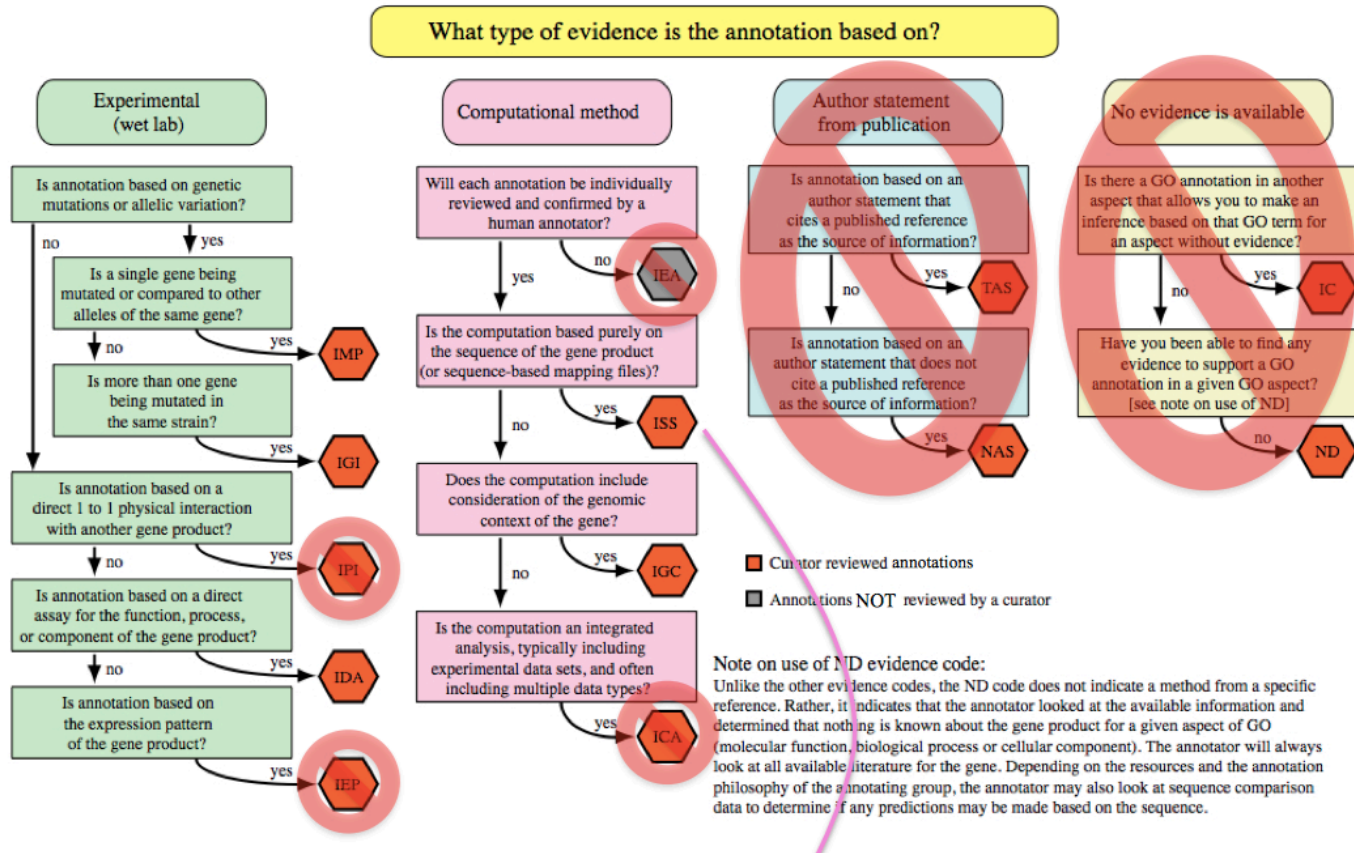
Functional Annotation w/GO

- 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



Evidence pull-down menu

special page

The Spring 2012 s



TableEdit

LAMBD:VLYS

Qualifier	<input type="text"/>
GO ID	<input type="text"/>
GO term name	
Reference	<input type="text"/>
Evidence Code	<input type="text"/>
with/from	
Aspect	
Notes	<input type="text"/>
Status	Missing: GO ID, evidence, reference

Public Refresh Save Row Cancel

Public rows can be edited or deleted by any user who can edit
Private rows can be edited or deleted by their creator, or by admins

navigation

- Main Page
- Enter GO at the top
- Help
- Report Bug
- Update log
- Annotation Jamborees
- Recent changes
- Create New Gene Page
- Login / Create Account

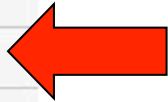
cacao

- Links about CACAO
- Fall 2011

journal clubs

- Journal Clubs
- Create new literature page

search



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

special page

The Spring 2012 s



TableEdit

LAMBD:VLYS

Qualifier	<input type="text"/>
GO ID	<input type="text"/>
GO term name	
Reference	<input type="text"/>
Evidence Code	<input type="text"/>
with/from	
Aspect	
Notes	<input type="text"/>
Status	Missing: GO ID, evidence, reference

Public Refresh Save Row Cancel

Public rows can be edited or deleted by any user who can edit
Private rows can be edited or deleted by their creator, or by admins

navigation

- Main Page
- Enter GO at the top
- Help
- Report Bug
- Update log
- Annotation Jamborees
- Recent changes
- Create New Gene Page
- Login / Create Account

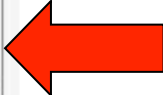
cacao

- Links about CACAO
- Fall 2011

journal clubs

- Journal Clubs
- Create new literature page

search



Example paper

<http://www.ncbi.nlm.nih.gov/pubmed/8227000>

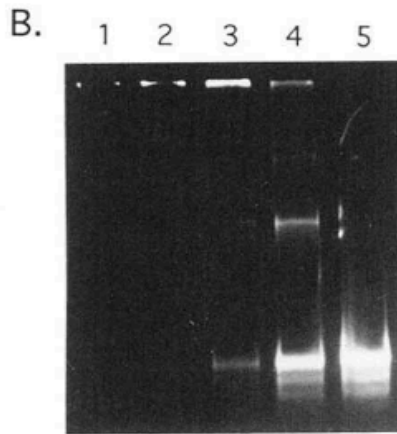
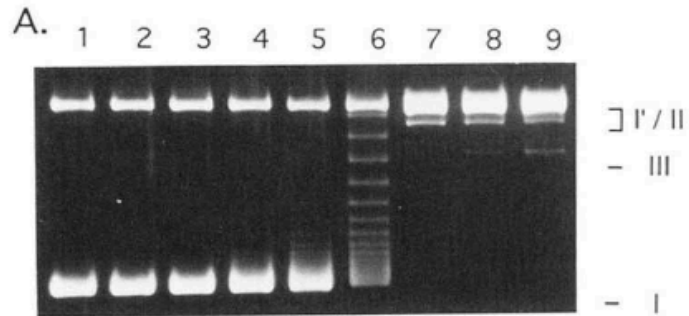
=

<http://www.jbc.org/content/268/32/24481.full.pdf>

What they did

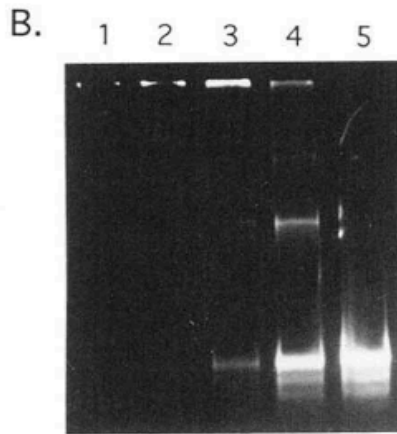
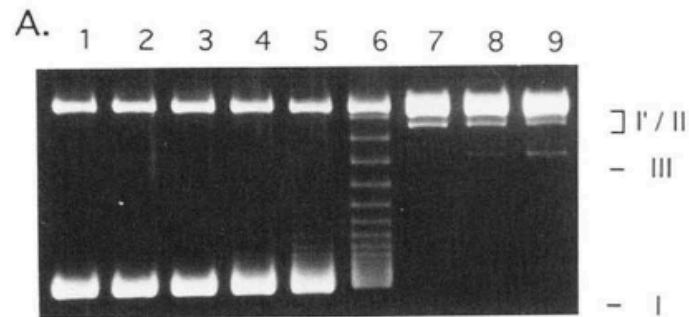
- Finding the proteins
- Do these tell us about the function?
 - Figure 1: sequenced ParC and Part of ParE
 - Figure 2: SDS page of purified proteins
 - Figure 3: Relaxation and decatenation activities of TopoIV
 - ...

Figure 3



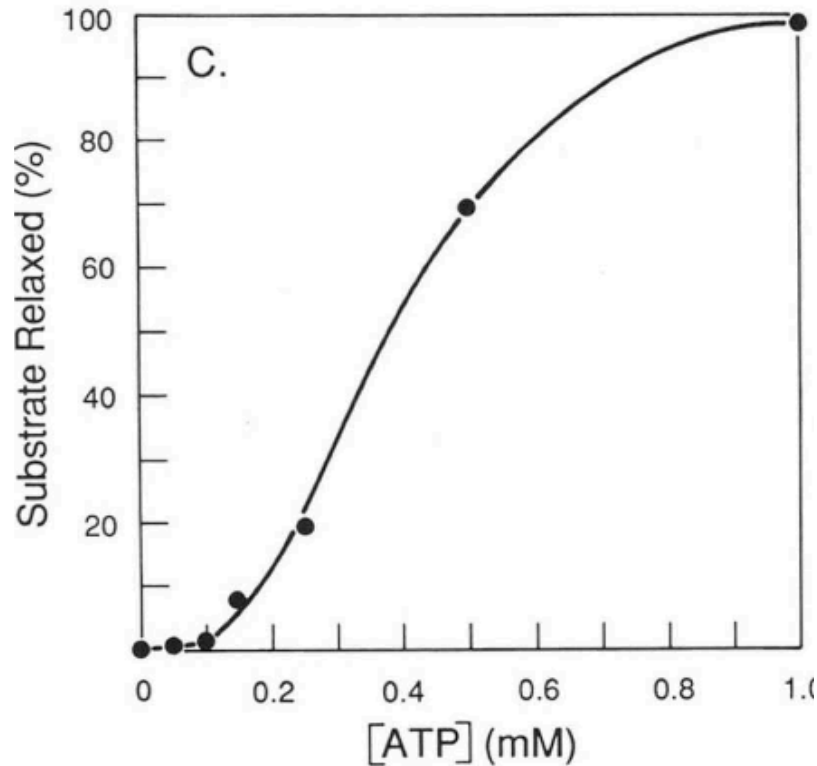
- Panel A: relaxation
- Panel B: decatenation
- What do these mean?

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

Figure 4



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

GO annotation for *E. coli* ParC

TableEdit

ECOLI:PARC

Qualifier	<input type="text" value=""/>
GO ID	<input type="text" value="GO:0003918"/>
GO term name	DNA topoisomerase type II (ATP-hydrolyzing) activity
Reference	PMID: <input type="text" value="8227000"/>
Evidence Code	<input type="text" value="IDA: Inferred from Direct Assay"/>
with/from	
Aspect	F
Notes	<input type="text" value="Topoisomerase assay in Fig 3. ATP dependent decatenation means it is a Type II from Fig 4"/>
Status	complete
<input type="text" value="Public"/> <input type="text" value="Refresh"/> <input type="text" value="Save Row"/> <input type="text" value="Cancel"/>	