



Protein Ontology (PRO)

for integration of knowledge on proteins, complexes and PTMs



Using Ontologies for Organizing Plant and Animal Genomics Data
January 14, 2012

Cathy H. Wu, Ph.D.

Director, Protein Information Resource (PIR)
Edward G. Jefferson Chair and Director

Center for Bioinformatics & Computational Biology, University of Delaware
Professor of Biochemistry & Molecular Biology, Georgetown University

PRO in OBO Foundry

OBO Foundry

- Establishing a set of principles to create a suite of orthogonal interoperable reference ontologies

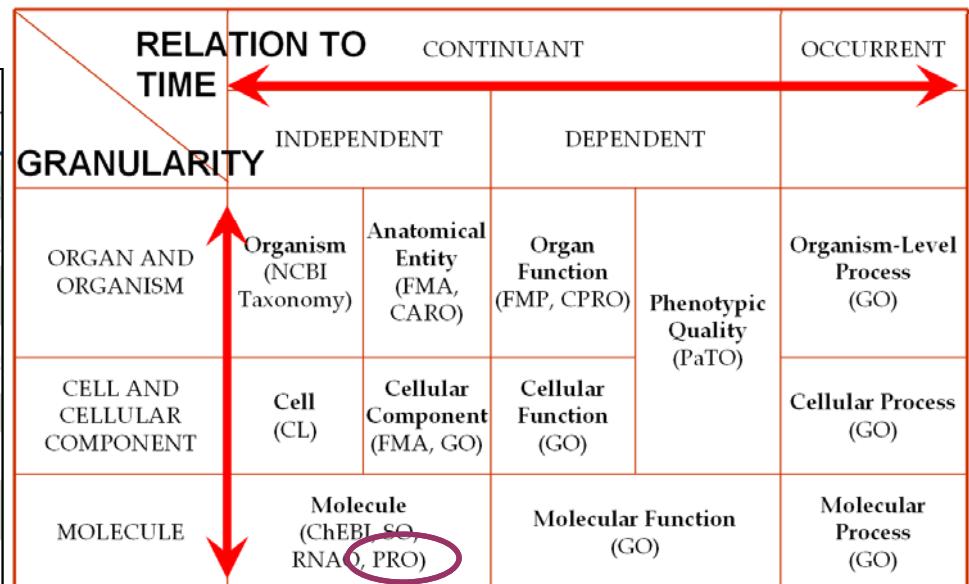


The Open Biological and Biomedical Ontologies

Protein Ontology (PRO)

- One of the first set of six OBO Foundry ontologies
- Reference Ontology for Proteins

OBO Foundry ontologies			
Title	Domain	Prefix	File
Biological process	biological process	GO	gene ontology edit.obo
Cellular component	anatomy	GO	gene ontology edit.obo
Chemical entities of biological interest	biochemistry	CHEBI	chebi.obo
Molecular function	biological function	GO	gene ontology edit.obo
Phenotypic quality	phenotype	PATO	quality.obo
PRotein Ontology (PRO)	proteins	PR	pro.obo
Xenopus anatomy and development	anatomy	XAO	xenopus_anatomy.obo
Zebrafish anatomy and development	anatomy	ZFA	zebrafish_anatomy.obo



The Protein Ontology: a structured representation of protein forms and complexes

Natale DA, Arighi CN, Barker WC, Blake JA, Bult CJ, Caudy M, Drabkin HJ, D'Eustachio P, Esvikov AV, Huang H, Nchoutmboube J, Roberts NV, Smith B, Zhang J, Wu CH. (2011)
Nucleic Acids Res. 39, D539-545 [PMID:20935045]

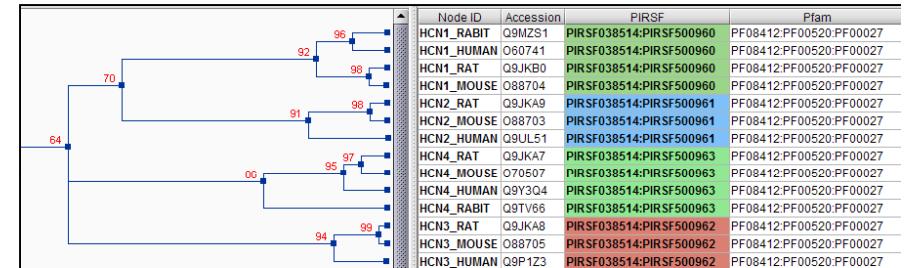


PRO Overview

PRO in OBO Foundry to represent protein entities

Three sub-ontologies to connect protein types necessary to model biology

- **Ontology for Protein Evolution (ProEvo):** Captures protein classes reflecting evolutionary relatedness of whole proteins



<input type="checkbox"/> PR:000000676	<i>potassium/sodium hyperpolarization-activated cyclic nucleotide-gated channel protein</i>	family
<input checked="" type="checkbox"/> PR:000000705	<i>potassium/sodium hyperpolarization-activated cyclic nucleotide-gated channel 1</i>	gene
<input checked="" type="checkbox"/> PR:000000706	<i>potassium/sodium hyperpolarization-activated cyclic nucleotide-gated channel 2</i>	gene
<input checked="" type="checkbox"/> PR:000000707	<i>potassium/sodium hyperpolarization-activated cyclic nucleotide-gated channel 3</i>	gene
<input checked="" type="checkbox"/> PR:000000708	<i>potassium/sodium hyperpolarization-activated cyclic nucleotide-gated channel 4</i>	gene

- **Ontology for Protein Forms (ProForm):** Captures different protein forms of a given gene locus from genetic variations, alternative splicing, proteolytic cleavage, PTMs

<input type="checkbox"/> PR:000002184	<i>Bcl2 antagonist of cell death</i>	gene	
<input checked="" type="checkbox"/> PR:000002280	<i>Bcl2 antagonist of cell death isoform 1</i>	sequence	
<input checked="" type="checkbox"/> PR:000003084	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated form</i>	modification	Q35147; Q61337; Q92934
<input checked="" type="checkbox"/> PR:000003085	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated 1</i>	modification	Q35147-1; Q61337-1
<input checked="" type="checkbox"/> PR:000003086	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated 2</i>	modification	
<input checked="" type="checkbox"/> PR:000003087	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated 3</i>	modification	Q61337-1:pS112/pS136
<input checked="" type="checkbox"/> PR:000003233	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated 4</i>	modification	Q61337-1:pS112/pS136/pS155
<input checked="" type="checkbox"/> PR:000003238	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated 5</i>	modification	Q35147-1:pS112
<input checked="" type="checkbox"/> PR:000003269	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated 6</i>	modification	Q61337-1:pT201
<input checked="" type="checkbox"/> PR:000025849	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated 7</i>	modification	Q61337-1:pS136
<input checked="" type="checkbox"/> PR:000025850	<i>Bcl2 antagonist of cell death isoform 1 phosphorylated 8</i>	modification	Q61337-1:pS128/pS136

Need for Representing Proteins Forms

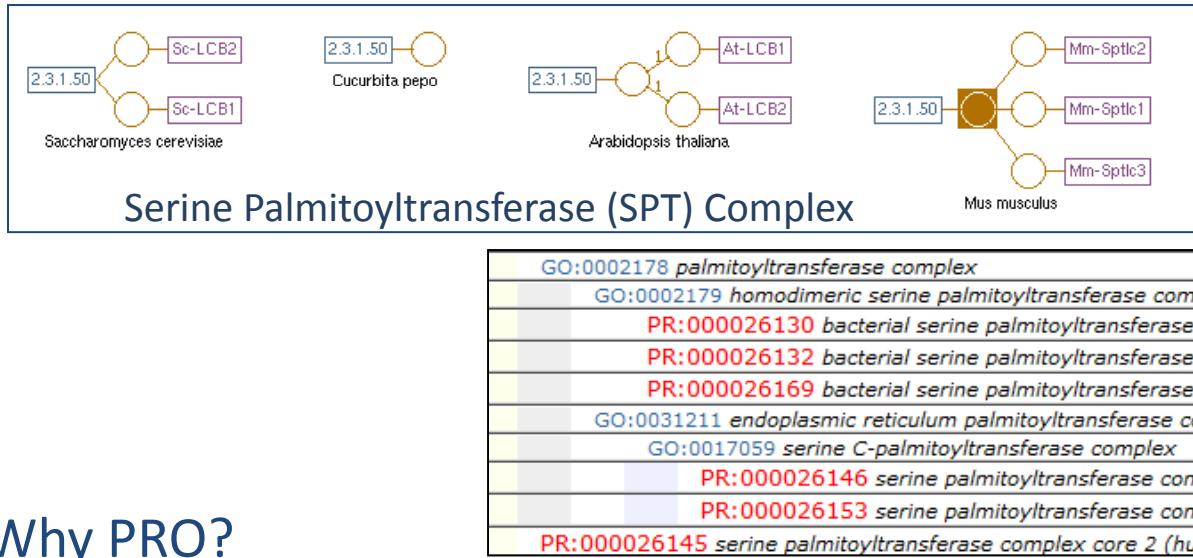
Alternative Splicing/Genetic Variation/PTM/Cleavage

Gene	Protein Form	Distinctive Attributes
SMAD2	Long isoform phosphorylated (PR:000000468)	<i>NOT has_function</i> GO:00036677 DNA binding
	Short isoform phosphorylated (PR:000000469)	<i>has_function</i> GO:00036677 DNA binding
CUL1	Unmodified form (PR:000002507)	<i>NOT part_of</i> GO:0019005 SCF ubiquitin ligase complex
	Neetylated form (PR:000000542)	<i>part_of</i> GO:0019005 SCF ubiquitin ligase complex
CD14	Membrane form (PR:000002149)	<i>located_in</i> GO:0005886 plasma membrane
	Soluble form (PR:000002147)	<i>located_in</i> GO:0005615 extracellular space
ROCK1	Full length (PR:000002529)	<i>has_function</i> GO:0004674 protein serine/threonine kinase activity
	Cleaved form (PR:000000563)	<i>Increased has_function</i> GO:0004674 protein serine/threonine kinase activity
CREBBP	Variant R → P(1378) (PR:000000266)	<i>agent_in MIM:180849, RUBINSTEIN-TAYBI SYNDROME</i> SO:1000118, <i>loss_of_function_of_polypeptide</i>

The diagram illustrates how specific protein forms are annotated with distinct attributes. Arrows point from the 'Distinctive Attributes' column to five categories: Function, Association, Localization, Modification, and Disease. The 'Function' arrow points to the 'NOT has_function' and 'has_function' annotations for SMAD2. The 'Association' arrow points to the 'NOT part_of' and 'part_of' annotations for CUL1. The 'Localization' arrow points to the 'located_in' annotations for CD14. The 'Modification' arrow points to the 'Increased has_function' annotation for ROCK1. The 'Disease' arrow points to the 'agent_in' and 'SO:1000118' annotations for CREBBP.

PRO Overview

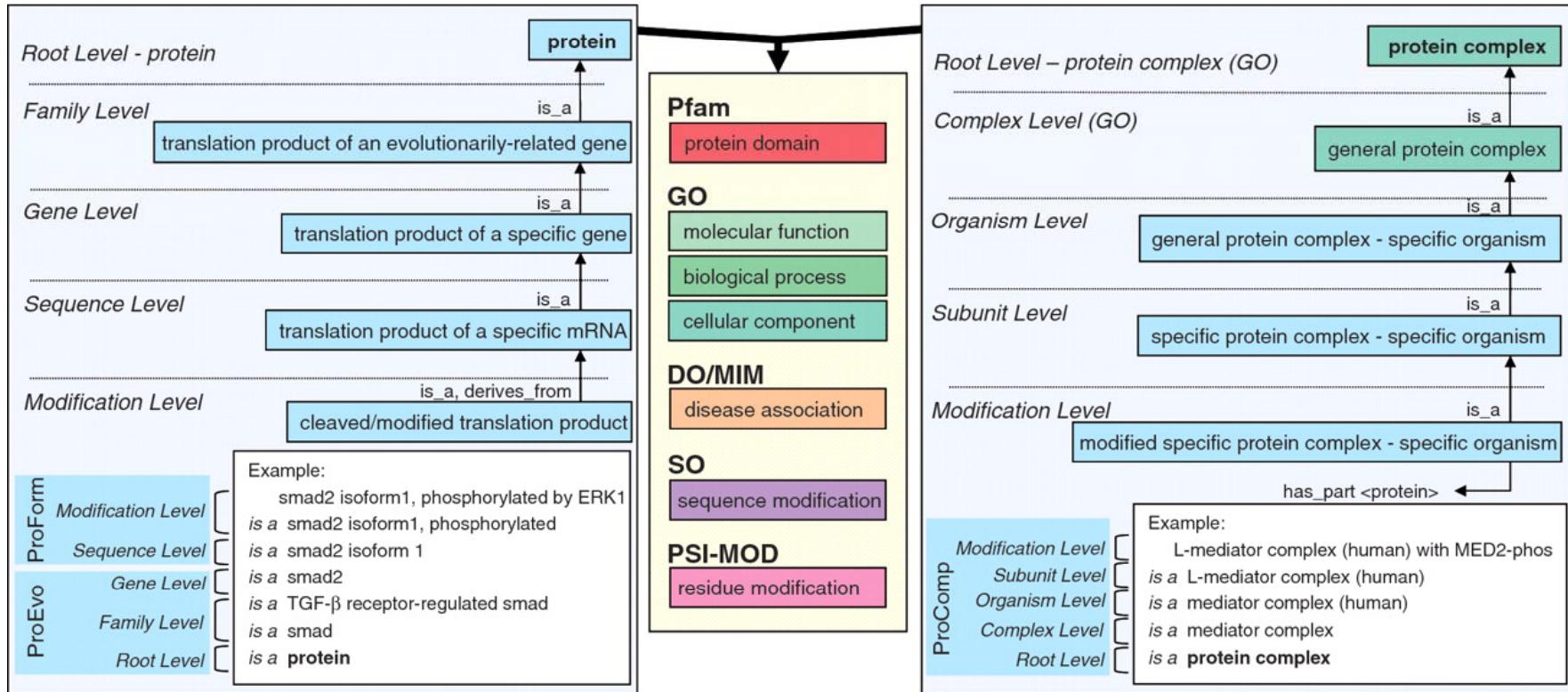
- Ontology for Protein Complexes (ProComp): Captures distinct complexes as they exist in different species and defines complexes through component proteins



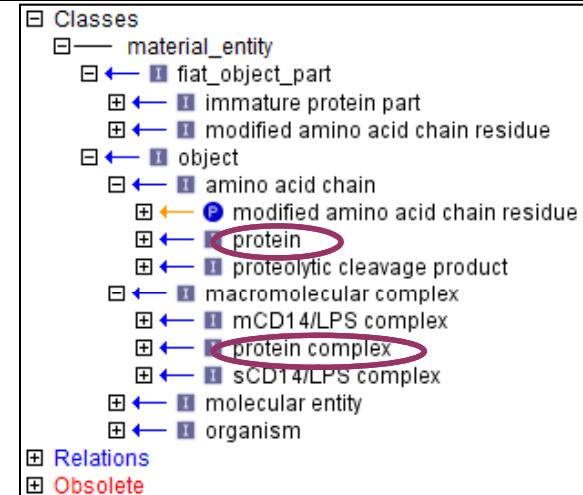
Why PRO?

- Provides formalization to support precise annotation of specific protein classes/forms/complexes, allowing accurate and consistent data mapping, integration and analysis
- Allows specification of relationships between PRO and other ontologies, such as GO, SO (Sequence Ontology), PSI-MOD, ChEBI, MIM/Disease Ontology, PO (Plant Ontology)
- Provides stable unique identifiers to distinct protein types
- Provides a formal structure to support computer-based reasoning based on homology and shared protein attributes, including “ortho-isoform,” “ortho-modified form”

PRO Framework



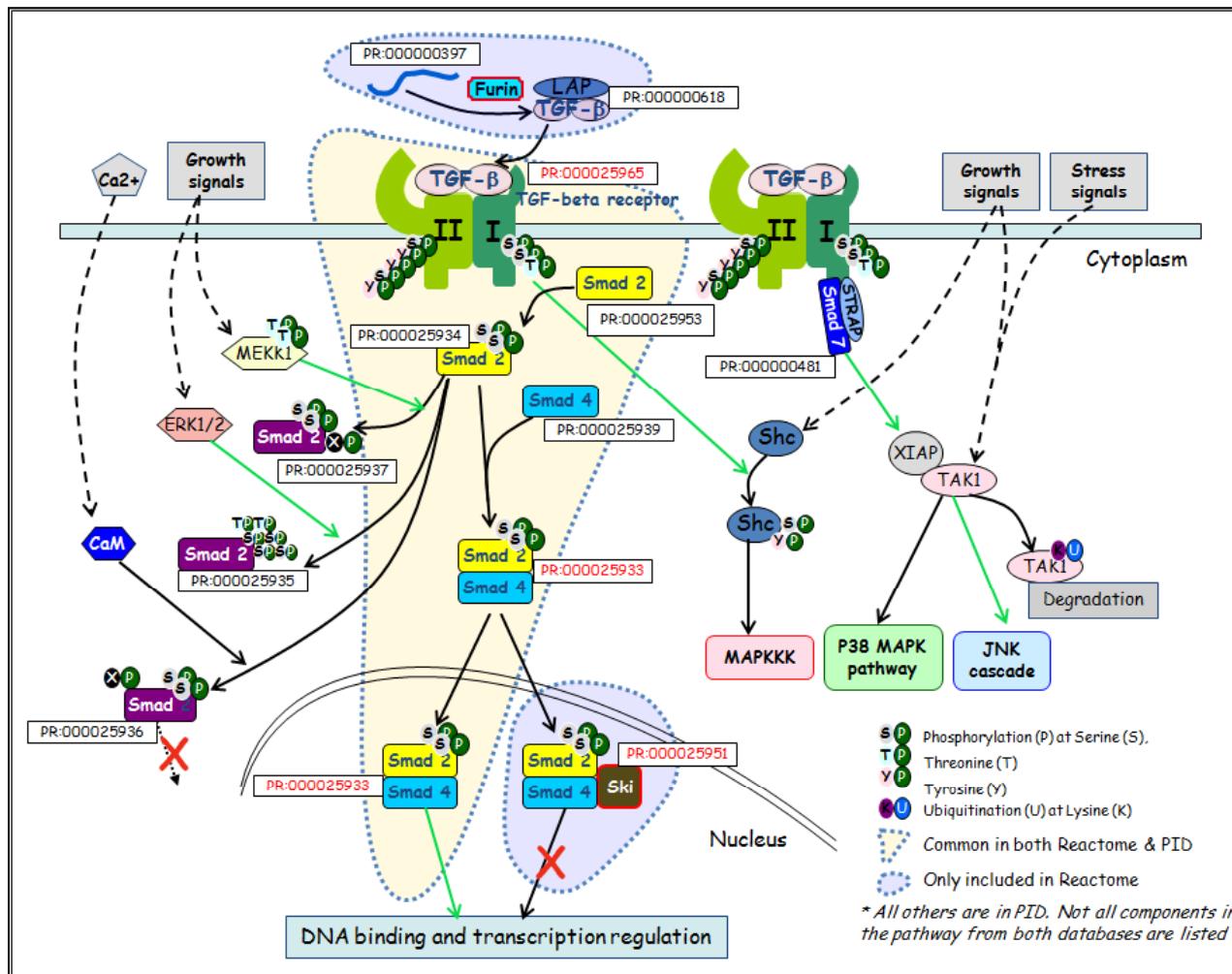
- PRO (ProForm, ProEvo, ProComp) is aligned with other OBO Foundry ontologies under the umbrella of the **Basic Formal Ontology (BFO)**
- PRO terms are defined/annotated using other ontologies and resources via definition of relations or mappings when appropriate



PRO in Pathway Context

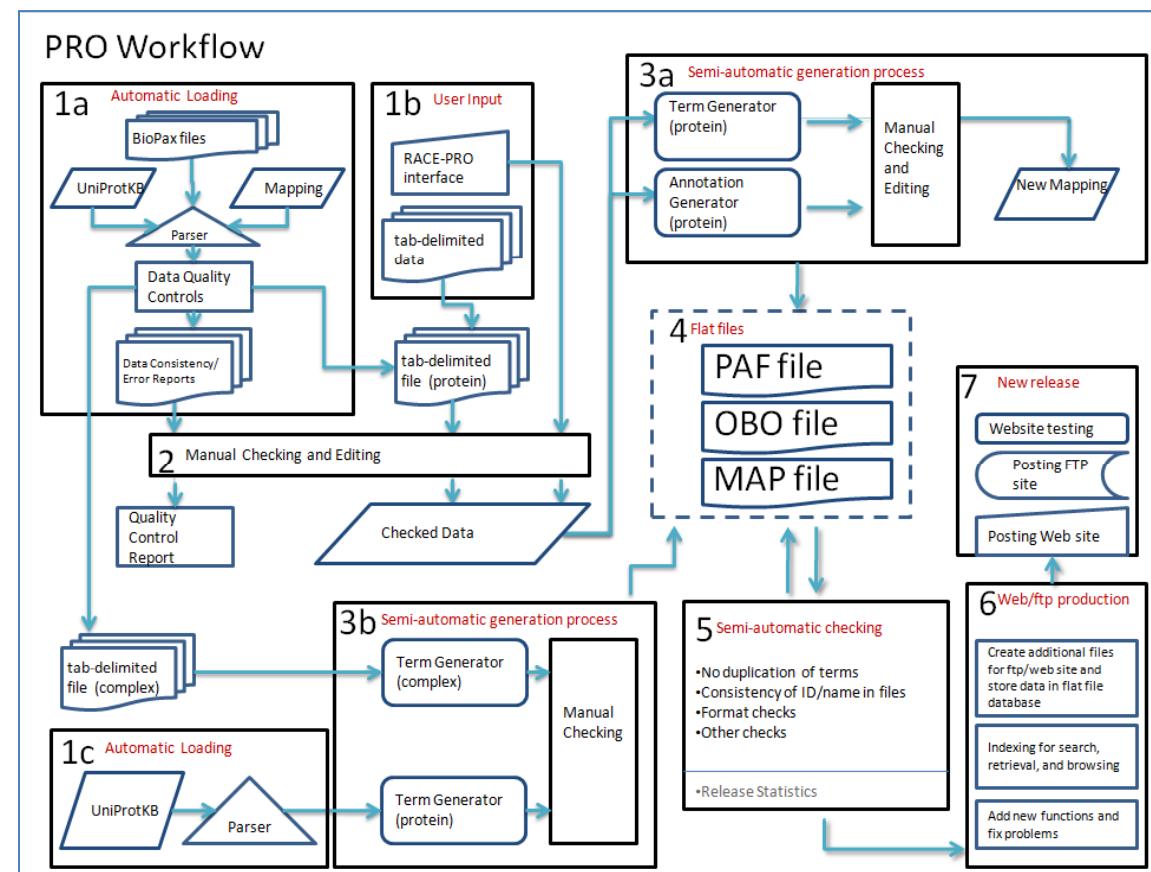
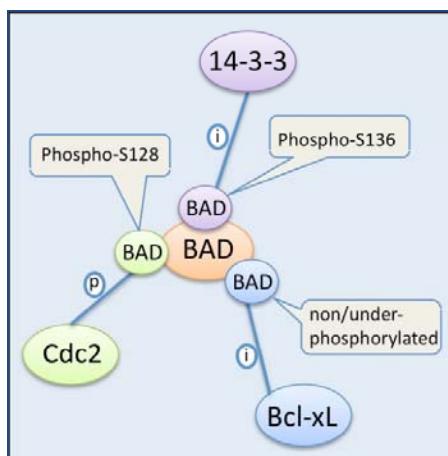
- Representation of protein forms & complexes **in biological/pathway/network context**
- Connecting multiple pathway/network resources

TGF- β Signaling Pathway: *Reactome, PID, coupled with literature curation*



PRO Workflow

- Data Sources
 - Manual annotation (curator, collaborator, user): sourceforge tracker; Race-PRO
 - Semi-automated processing of external databases (e.g., UniProtKB, Reactome, MouseCyc, EcoCyc); coverage of 12 reference genomes forthcoming
- Integration with text mining tool: eFIP (*Functional Impact of Phosphorylation*)
- Distribution Files
 - Ontology (OBO)
 - Annotation (PAF)
 - Mappings (exact; is_a)



Plant PRO Terms

- *Arabidopsis thaliana*: 3842 gene-level terms to be added soon based on the PANTHER mapping of 12 reference genomes
- 59 terms for *Arabidopsis thaliana* and wheat on protein forms and complexes

PRO Home Any field AND Any field + add input box - del input box

Cullin-Related PRO Terms

Display Options □
42 entries | 1 page | 50 / page |

42 selected [\(show\)](#) click to show: selected [Hierarchy](#) selected [OBO / PAF](#) OR related [OBO / PAF](#) [Cytoscape view](#)

<input checked="" type="checkbox"/> PRO ID	PRO Name	PRO Term Definition	Category	Parent
<input checked="" type="checkbox"/> GO:0019005	SCF ubiquitin ligase complex	A ubiquitin ligase complex in which a cullin from the Cul1 subfamily and a RING domain protein form the catalytic core; substrate specificity is conferred by a Skp1 adaptor and an F-box protein. SCF complexes are involved in targeting proteins for degradation by the proteasome. The best characterized complexes are those from yeast and mammals (with core subunits named Cdc53/Cul1, Rbx1/Hrt1/Roc1). [PMID: 15571813 , PMID: 15688063]		GO:0031461
<input checked="" type="checkbox"/> GO:0031461	cullin-RING ubiquitin ligase complex	Any ubiquitin ligase complex in which the catalytic core consists of a member of the cullin family and a RING domain protein; the core is associated with one or more additional proteins that confer substrate specificity. [PMID: 15571813 , PMID: 15688063]		GO:0043234
<input checked="" type="checkbox"/> PR:000000013	cullin	A protein with a core domain composition consisting of a cullin domain that functions as a molecular scaffold responsible for assembling the ROC1/Rbx1 RING-based E3 ubiquitin ligases. [PIRSF:PIRSF017874]	family	PR:000000001
<input checked="" type="checkbox"/> PR:000000043	RING-box protein 1	A RING-box protein 1 related protein that is a translation product of the RBX1 gene or a 1:1 ortholog thereof. [PRO:CNA]	gene	PR:000027961
<input checked="" type="checkbox"/> PR:000000044	RING-box protein 2	A protein that is a translation product of the RNF7 gene or a 1:1 ortholog thereof. [PRO:CNA]	gene	PR:000000004
<input checked="" type="checkbox"/> PR:000000092	fungal/metazoan cullin-1	A cullin that is a translation product of the CUL1 gene or a 1:1 ortholog thereof. [PRO:CNA]	gene	PR:000000013
<input checked="" type="checkbox"/> PR:000000267	cullin-1 isoform 1	A cullin-1 that is a translation product of some mRNA giving rise to a protein with the amino acid sequence represented by UniProtKB: Q13616-1 or a 1:1 ortholog thereof. [PRO:CNA]	sequence	PR:000000092

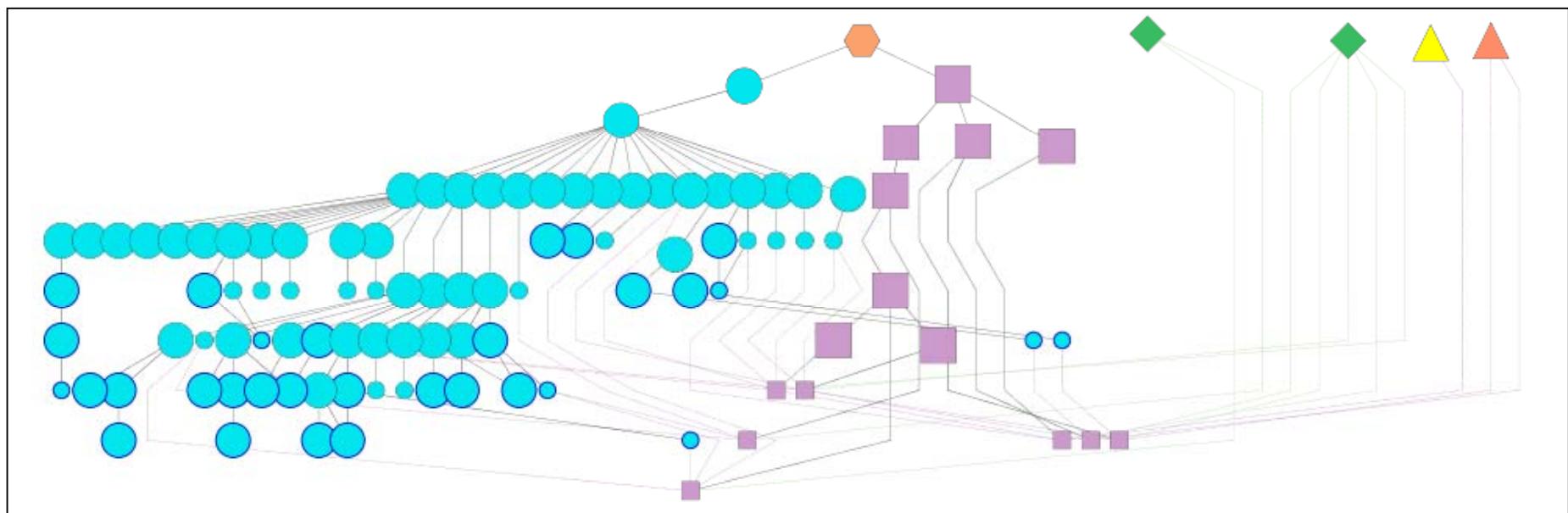
Hierarchical View: Cullin-Related Terms

66 shown of 28321 records		PMID	Taxon	PANTHER	EcoCyc	Definition			
		Synonym	Gene	MGI	HGNC	Pfam	PIRSF	Reactome	UniProtKB
	expand		↓ sort (ID)	↓ sort (STR)		↑ find			Category
	GO:0032991	macromolecular complex							
	GO:0043234	protein complex							
	GO:0031461	cullin-RING ubiquitin ligase complex							
	GO:0019005	SCF ubiquitin ligase complex							
	PR:000026771	SCF beta-TrCP complex (human)							organism-complex
	PR:000028464	SCF(CO11) complex							complex
	PR:000028465	SCF(CO11) complex (<i>Arabidopsis thaliana</i>)							organism-complex
	PR:000028457	SCF(TIR1) complex							complex
	PR:000028458	SCF(TIR1) complex (<i>Arabidopsis thaliana</i>)							organism-complex
	PR:000026772	SCF(Skp2) complex (human)							organism-complex
	PR:000018263	amino acid chain							
	PR:000000001	protein							
	PR:000000013	cullin							family
	PR:000006047	cullin-2							gene
	PR:000006048	cullin-3							gene
	PR:000006049	cullin-4A							gene
	PR:000006050	cullin-4B							gene
	PR:000006051	cullin-5							gene
	PR:000000092	fungal/metazoan cullin-1							gene
	PR:000026757	cullin-1 (human)							organism-gene
	PR:000000267	cullin-1 isoform 1							sequence
		PR:000000427 cullin-1 isoform 1 neddylated form							modification
		PR:000000542 cullin-1 isoform 1 neddylated 1							modification
		PR:000002507 cullin-1 isoform 1 unmodified form							modification
	PR:000027962	plant cullin-1							gene
	PR:000027963	cullin-1 (<i>Arabidopsis thaliana</i>)							organism-gene
	PR:000027976	plant cullin-1 rubylated form							modification
		PR:000027977 cullin-1 rubylated (<i>Arabidopsis thaliana</i>)							organism-modification
	PR:000028501	plant cullin-2							gene
	PR:000028494	cullin-2 (<i>Arabidopsis thaliana</i>)							organism-gene
	PR:000028495	plant cullin-3A							gene
	PR:000028496	cullin-3A (<i>Arabidopsis thaliana</i>)							organism-gene
	PR:000028497	plant cullin-3B							gene
	PR:000028498	cullin-3B (<i>Arabidopsis thaliana</i>)							organism-gene

Cytoscape Network View: Cullin-Related Terms

Connecting protein forms and complexes with annotation => Modeling biology

- Show full set of relationships between terms, including those from other ontologies
- Traverse hierarchical (parent-child) relationships
- Find all components of a complex, including modified forms
- Find all complexes that a protein is constituent of = first (complex) neighbors of a protein
- Find all complexes that bind small molecules = those that have relationship to ChEBI



purple square complex

purple square Organism-complex

cyan circle protein

cyan circle with black outline modified protein form

blue circle Organism-protein

yellow triangle ChEBI

green diamond Taxon 11

Cytoscape & Entry Views

PRO Home Protein Ontology report for entry - PR:000028458

[Show OBO stanza](#) [Retrieve related PRO](#)

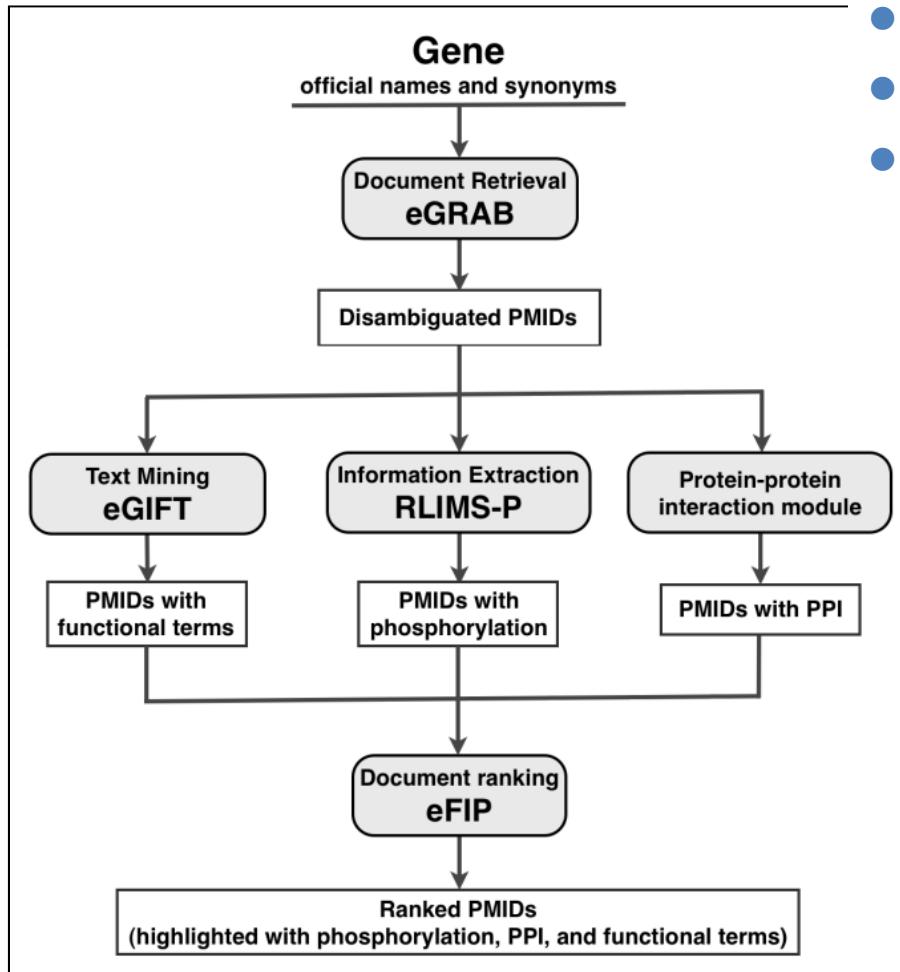
Ontology Information							
PRO ID	PR:000028458						
PRO name	SCF(TIR1) complex (Arabidopsis thaliana)						
Synonyms							
Definition	An SCF(TIR1) complex consisting of TIR1, SKP1A(ASK1), rubylated CUL1, and one of RBX1 proteins, whose component is cullin-1. [PRO:CNA, PMID:21370976]						
Comment	Category=organism-complex.						
Hierarchical relationship	Parent: PR:000028457 SCF(TIR1) complex Children: none has_part PR:000027956 {cardinality="1"} protein TRANSPORT INHIBITOR RESPONSE 1 (Arabidopsis thaliana) has_part PR:000027960 {cardinality="1"} SKP1-like protein 1A (Arabidopsis thaliana) has_part PR:000028456 {cardinality="1"} RING-box protein 1 (Arabidopsis thaliana) has_part PR:000027977 {cardinality="1"} cullin-1 rubylated (Arabidopsis thaliana) only_in_taxon NCBITaxon:3702 Arabidopsis thaliana						
Annotation							
Functional Annotation	Modifier	Relation	Ontology ID	Ontology Term	Relative_to	Interaction With	Evidence
		participates_in	GO:0009734	auxin mediated signaling pathway			PMID:21370976

Connecting protein forms and complexes with annotation

Data Panel

ID	ontology.name	Relation	Ontology_ID	Ontology_term	Evidence_source
PR:000028458	SCF(TIR1) complex (Arabidopsis thaliana)	participates_in	GO:0009734	auxin mediated signaling pathway	PMID:21370976

eFIP: an integrated system for mining Functional Impact of Phosphorylation from literature



- Entity recognition and document retrieval
- Information extraction
- Document ranking and evidence tagging

PMID 10837486 for gene BAD - Bcl2-associated agonist of cell death

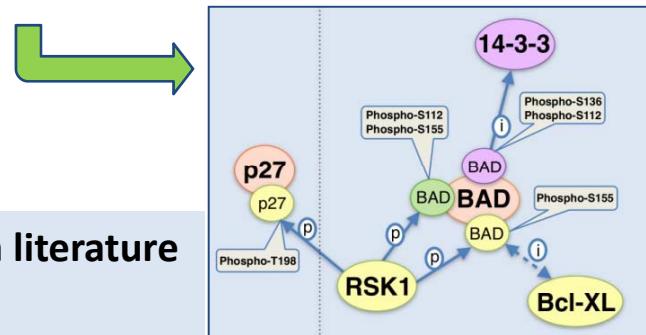
Predicted impact of phosphorylation:

Sentence #	Phosphorylation			Impact
	Substrate	Site	Kinase	
1	BAD	Ser-155	RSK1	regulates BAD/Bcl-XL interaction regulates cell survival
3,4	BAD	Ser-112 and Ser-136	N/A	promotes binding of BAD to 14-3-3 proteins
6	BAD	Ser-155	RSK1	blocking the binding of BAD to Bcl-XL
7	BAD	both Ser-112 and Ser-155	RSK1	rescues BAD-mediated cell death

Tag substrate
 Tag kinase
 Tag phosphorylation site
 Tag protein-protein interaction
 Tag functional term

Text of title and abstract:

Sentence #	Sentence
1	T1: BAD Ser-155 PHOSphorylation regulates BAD/Bcl-XL interaction and cell survival .
2	AB - The BH3 domain of BAD mediates its death-promoting activities via heterodimerization to the Bcl-XL family of death regulators .
3	Growth and survival factors inhibit the death-promoting activity of BAD by stimulating PHOSphorylation at multiple sites including Ser-112 and Ser-136 .
4	PHOSphorylation at these sites promotes binding of BAD to 14-3-3 proteins , sequestering BAD away from the mitochondrial membrane where it dimerizes with Bcl-XL to exert its killing effects .
5	We report here that the phosphorylation of BAD at Ser-155 within the BH3 domain is a second PHOSphorylation-dependent mechanism that inhibits the death-promoting activity of BAD .
6	Protein kinase A , RSK1 and survival factor signaling stimulate PHOSphorylation of BAD at Ser-155 , blocking the binding of BAD to Bcl-XL .
7	BAD phosphorylates BAD at both Ser-112 and Ser-155 and rescues BAD -mediated cell death in a manner dependent upon PHOSphorylation at both sites .

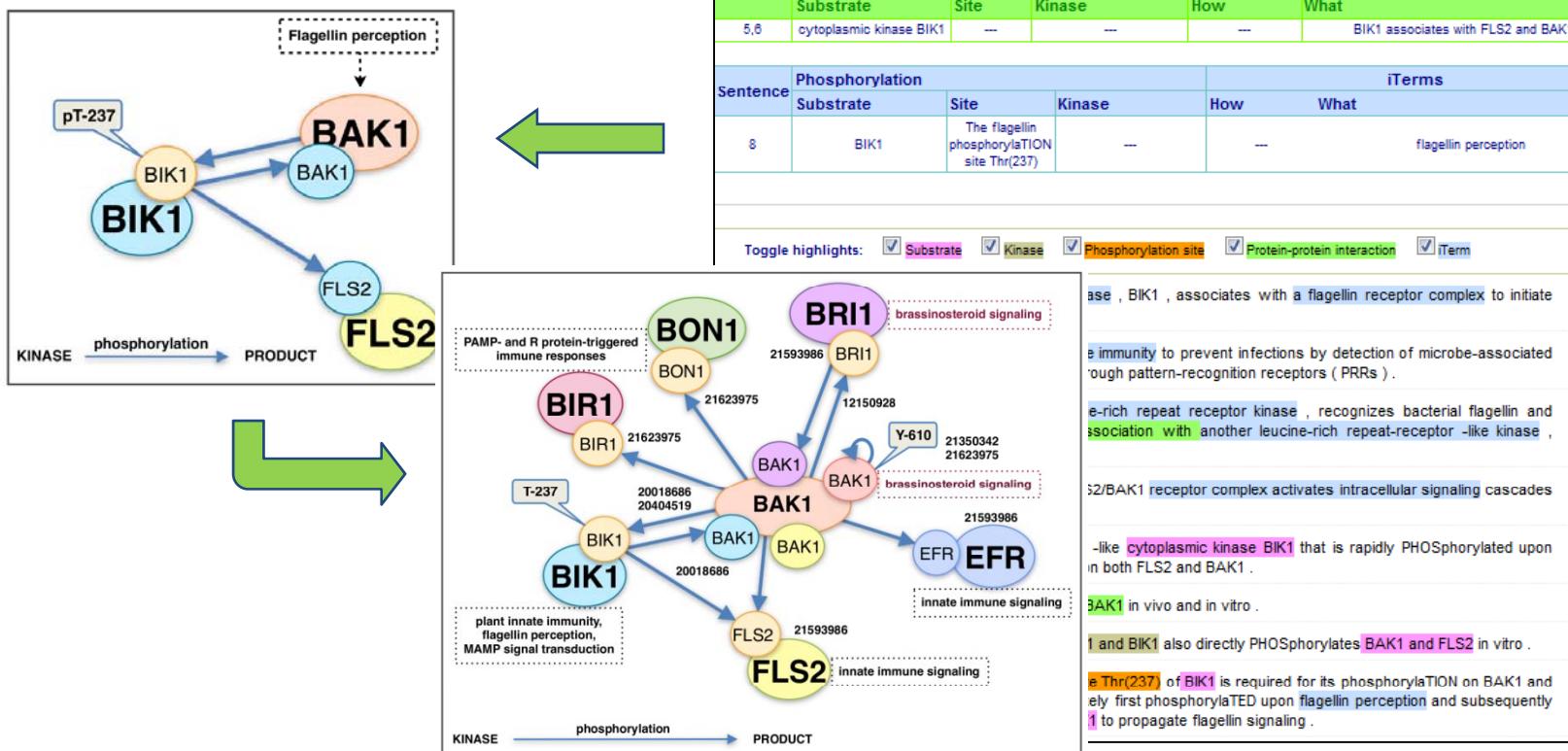


eFIP: a tool for mining functional impact of phosphorylation from literature

Arighi CN, Siu AY, Tudor CO, Nchoutmboue JA, Wu CH, Shanker VK. (2011)
Methods in Molecular Biology 694, 63-75 [PMID:21082428]

Discovery from Literature Mining

- iPTMnet: PTM network of enzyme-substrate relationships and protein-protein interactions mined from the scientific literature
- Distinct phosphorylated forms of a protein may have different interacting proteins, leading to different subcellular locations, functions and pathways
=> Knowledge captured in PRO



User Annotation: RACE-PRO

Capture knowledge of protein forms/complexes of interest to support integrated analysis and computer-based reasoning

- ❖ Obtain a PRO ID for the protein objects of interest
- ❖ Define a protein object (based on literature, experimental data)
- ❖ Add annotation to that protein object
- ❖ How does it work?
 - ❖ Input your personal information (only for internal use)
 - ❖ Complete form with sequence information and annotation
 - ❖ Submit when ready (otherwise you can save for later)
 - ❖ PRO curation team will take the data, revise it, and create the corresponding PRO node in the ontology
 - ❖ User will be informed via email about the new PRO IDs and when they will be public

RACE-PRO Annotation Tool

Definition of the Protein Object

1. Enter a UniProtKB identifier (?)

OR, click here to insert a different sequence:

```
MTRDFKPGDL IFAKMKGYPH WPARVDEVPD GAVKPPPTNKL PIFFFGTHET AFLGPKDIFP 60
YSENKEKYKGK PNKRKGFGNEG LWEIDNNNPKV KFSSQQAAATK QSNASSDVEE EEEKETSVSKE 120
DTDHEEKASN EDVTKAADVIT TPKAARRGRK RKAEKQVETE EAGVVTTATA SVNLKVSPKR 180
GRPAATEVKI PKPRGRPKMV KQCPSESDI ITEEDKSKKK GQEEKQPKIQ PRKDEEGQKE 240
EDKPRKREPDK KEKGKKEVESK RKNLAKTGV TSTDSEEEGD DQEGEKKRIG GRNFQTAHRR 300
NMLIKQHEKE AADRRRKQEE QMTEHQTC NLQ
```

2. Specify sequence region

Full-length Region: from to

3. Indicate post-translational modifications (add amino acid number re)

Amino acid number: -choose PTM-

4. Protein Object name (separate multiple names using ",")

LEDGF/p38; DN85

5. Evidence Source (separate multiple IDs using ",")

Db name: PubMed IDs: 18708362

A-Enter accession or paste sequence

Organism: HOMO SAPIENS

B-Define protein region and/or PTMs

C-Enter protein form name

D- Data source

E- Annotation

Annotation of the Protein Object

Domain [more] [less] Link to PFAM

Modifier	Relation	Pfam ID	Pfam name	PMIDs
NOT	has_part	PF00855	PWWP domain	18708362

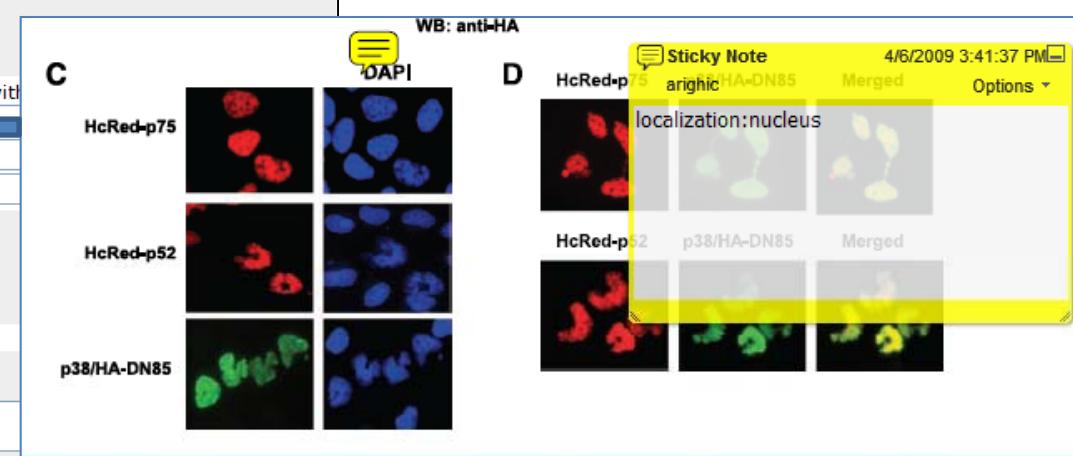
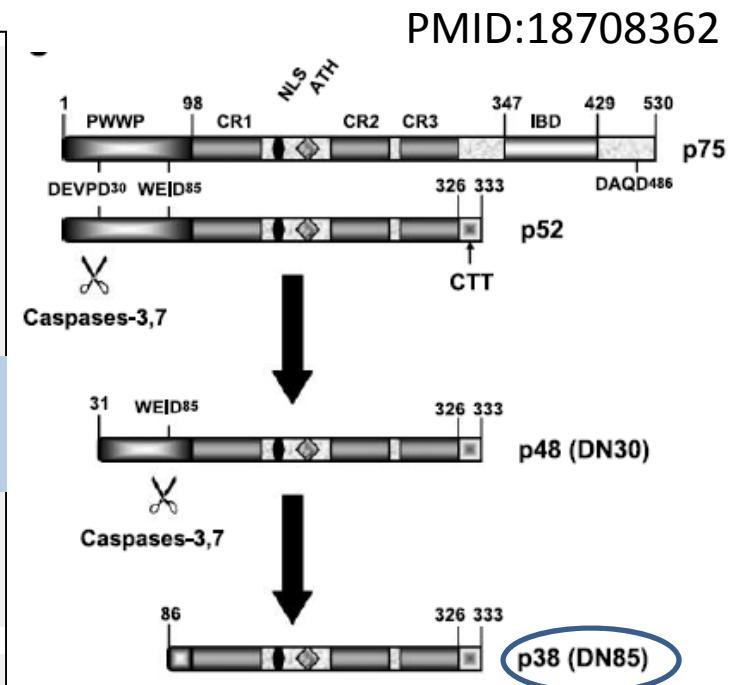
Functional Annotation [more] [less] Link to GO

Modifier	Relation	GO ID	GO term	Interaction with
	located_in	GO:00056	nucleus	
	participates_in	GO:001648	negative regulation of tra	

Sequence Ontology [add] Link to SO

Disease [add] Link to MIM

Comments:



PRO Dissemination

- PRO Website (<http://www.proconsortium.org>)
 - Searching, browsing, downloading
- PRO Views
 - Entry view
 - Table summary
 - OBO stanza, OWL
 - Ontology hierarchy
 - Cytoscape network
- PRO Link: Persistent URL: http://purl.obolibrary.org/obo/PR_xxxxxxxxxx
- OBO Foundry (<http://www.obofoundry.org/>)
- NCBO Bioportal (<http://bioportal.bioontology.org/>)

PRO Communities

- Ontology Developers
 - GO ontology: Interfaces of GO/PRO complexes; GO definition (e.g., GO:0005109)
 - GO annotation: precise annotation of protein forms in PomBase
 - Dendritic Cell Ontology: Define cell types based on +/- protein types [PMID:19243617]
 - Annotation Ontology for annotating scientific documents on the web [PMID:21624159]
 - Brucellosis Ontology (IDOBRU), extension of the Infectious Disease Ontology (IDO) [PMID:22041276]
- Semantic Resources
 - Royal Society of Chemistry (RSC); Science Collaboration Framework; Semantic Web Applications in Neuromedicine (SWAN); Neuroscience Information Framework (NIF)
- Pathway/Process-Modeling Resources:
 - Reactome, MouseCyc, EcoCyc/BioCyc, Center for Molecular Immunology (Duke)
- Molecule-Modeling Resources: Int'l Union of Basic and Clinical Pharmacology (IUPhar)
- Pharma/Clinical Communities: Drug Discovery & Disease Biomarker
 - Alzforum
 - Salivaomics KB/SALO (Saliva Ontology): Saliva Biomarkers
 - Pistoia

PRO Consortium Team (current)

Protein Information Resource (PIR) [Georgetown U & U Delaware]

Cathy Wu, Cecilia Arighi, Darren Natale, Natalia Roberts
Hongzhan Huang, Jian Zhang



The Jackson Lab – Mouse Genome Informatics (MGI)

Judith Blake, Carol Bult, Harold Drabkin, Alexei Evsikov



University at Buffalo-SUNY

Barry Smith, Alan Ruttenberg, Alexander Diehl

NYU School of Medicine – Reactome

Peter D'Eustachio, Michael Caudy



AlzForum

Elizabeth Wu

New Consortium Member Welcome!



1R01GM080646-01
3R01GM080646-04S2
2R01GM080646-06

PRO: A Protein Ontology in OBO Foundry for Integration of Biomedical Knowledge
(http://projectreporter.nih.gov/project_info_description.cfm?aid=8187900&icde=10696318)