

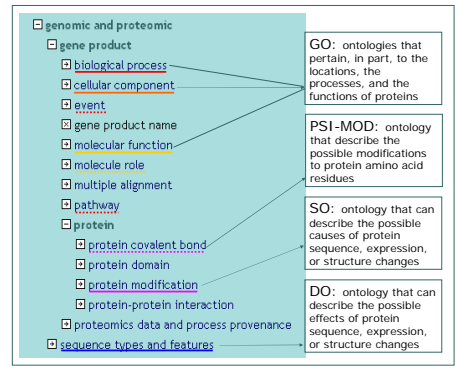
Darren A. Natale<sup>1</sup>, Cecilia N. Arighi<sup>1</sup>, Lai-Su L. Yeh<sup>1</sup>, Winona C. Barker<sup>1</sup>, Judith Blake<sup>2</sup>, Ti-Cheng Chang<sup>1</sup>, Zhangzhi Hu<sup>1</sup>, Hongfang Liu<sup>3</sup>, Barry Smith<sup>4</sup>, and Cathy H. Wu<sup>1</sup>

<sup>1</sup>Protein Information Resource, <sup>2</sup>The Jackson Laboratory, <sup>3</sup>Department of Biostatistics, Bioinformatics, and Biomathematics, Georgetown University, Washington, <sup>4</sup>Department of Philosophy, State University of New York at Buffalo

**Introduction**

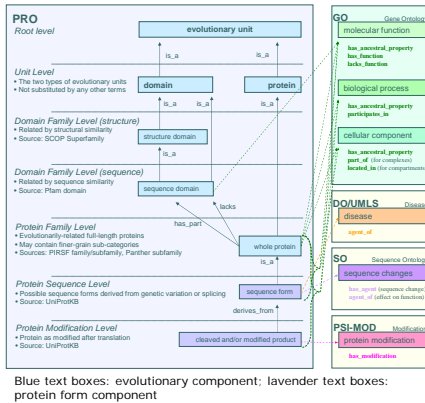
Biological ontologies are part of caBIG initiatives critical for cancer data integration and sharing through the grid environment. A number of these ontologies describe the properties that can be attributed to proteins (See right figure). But the lacking of an appropriate ontology for protein objects and classes has impeded the genome/proteome annotation, data integration and analysis. **PROtein Ontology (PRO)** is designed to provide this ontological framework by:

- Providing a structure to support formal, computer-based inferences of shared attributes among homologous proteins
- Delineating the multiple protein forms of a gene locus
- Interconnecting existing OBO foundry ontologies
- Providing an enabling technology for knowledge discovery



**PRO Framework**

The components of PRO extend from the classification of proteins on the basis of evolutionary relationships to the representation of the multiple protein forms of a gene (products generated by genetic variation, alternative splicing, proteolytic cleavage, and other post-translational modification). The figure in this panel shows the current working model and a subset of the possible connections to other ontologies.



**PRO Prototype: the transforming growth factor (TGF) beta signaling pathway**

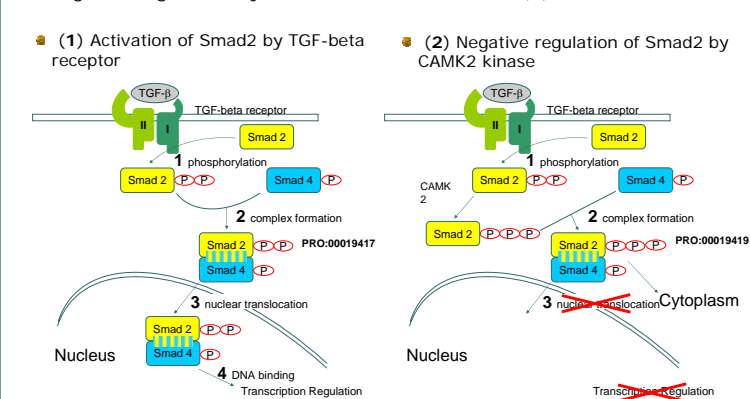
PRO connections are first generated by automated extraction of existing annotations from the PIRSF database, for the evolutionary component, and the UniProtKB/Swiss-Prot and MGI entries to get information on the protein forms. Manual curation involves the verification/creation of the protein form nodes and the assignment of the appropriate experimentally verified ontological/controlled vocabulary terms to the sequence forms.

Here we describe the initial development of PRO, illustrated using human and mouse Smad proteins that are essential components of the TGF-beta signaling pathway. These proteins are regulated by phosphorylation and the phosphorylation state dictates cellular location and activity.

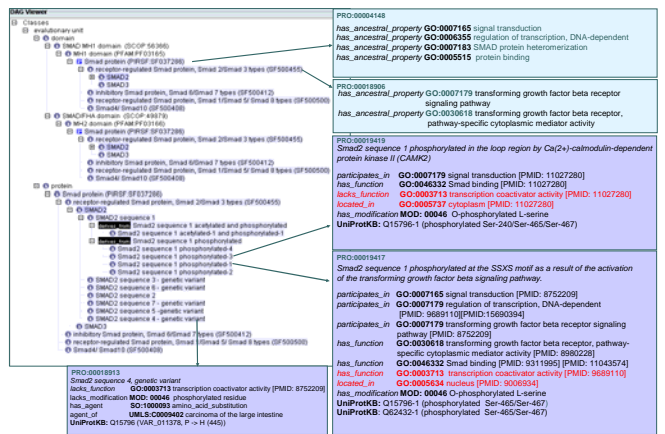
**Smad2 in the transforming growth factor beta signaling pathway**

The figure below (1) depicts the TGF-beta signaling pathway, focusing on the Smad 2 component. The steps shown are preceded by TGF-beta binding to the receptor, and receptor phosphorylation.

- Step 1: Phosphorylation of Smad2 by activated TGF beta receptor 1.
- Step 2: Complex formation of Smad2 and Smad4.
- Step 3: Nuclear import of Smad2:Smad4.
- Step 4: Binding of Smad2:Smad 4 complex coactivator to responsive element.



PRO can describe the different forms of Smad2 attaching the specific ontology terms. As an example, the two phosphorylated forms depicted in the graph on the left (**PRO:00019417** and **PRO:00019419**) are shown in detail with the differential properties marked in red. In addition, the evolutionary component not only provides the relationship of Smad2 with other Smad protein families (blue boxes), but also may allow to infer properties of homologous proteins in other organisms.



PRO also describes the properties of disease-related genetic variants

**Conclusions**

- PRO is designed to be a formal, well-principle and extensible OBO Foundry ontology for proteins
- PRO will provide research technology to answer new scientific questions. For example, the transfer of described function/phenotypes of protein from model organisms to human orthologs may highlight potential candidates to explain human disease
- PRO is expected to create a cycle of improvement for both the ontology and the protein knowledgebases from which the initial information is extracted
- PRO will promote cancer data integration and sharing via the caBIG grid

