

# Framework for a Protein Ontology

TMBIO  
November 2006

Darren A. Natale, Ph.D.  
Protein Science Team Lead, PIR  
Research Assistant Professor, GUMC

## [-] genomic and proteomic

### [-] gene product

[+] biological process

[+] cellular component

[+] event

[x] gene product name

[+] molecular function

[+] molecule role

[+] multiple alignment

[+] pathway

### [-] protein

[+] protein covalent bond

[+] protein domain

[+] protein modification

[+] protein-protein interaction

[+] proteomics data and process provenance

[+] sequence types and features

**GO:** ontologies that pertain, in part, to the locations, the processes, and the functions of proteins

**PSI-MOD:** ontology that describe the possible modifications to protein amino acid residues

**SO:** ontology that can describe the possible causes of protein sequence, expression, or structure changes

**DO:** ontology that can describe the possible effects of protein sequence, expression, or structure changes

## Mothers against decapentaplegic homolog 2

Smad 2

GO annotation of SMAD2\_HUMAN:

*Cellular Component:*

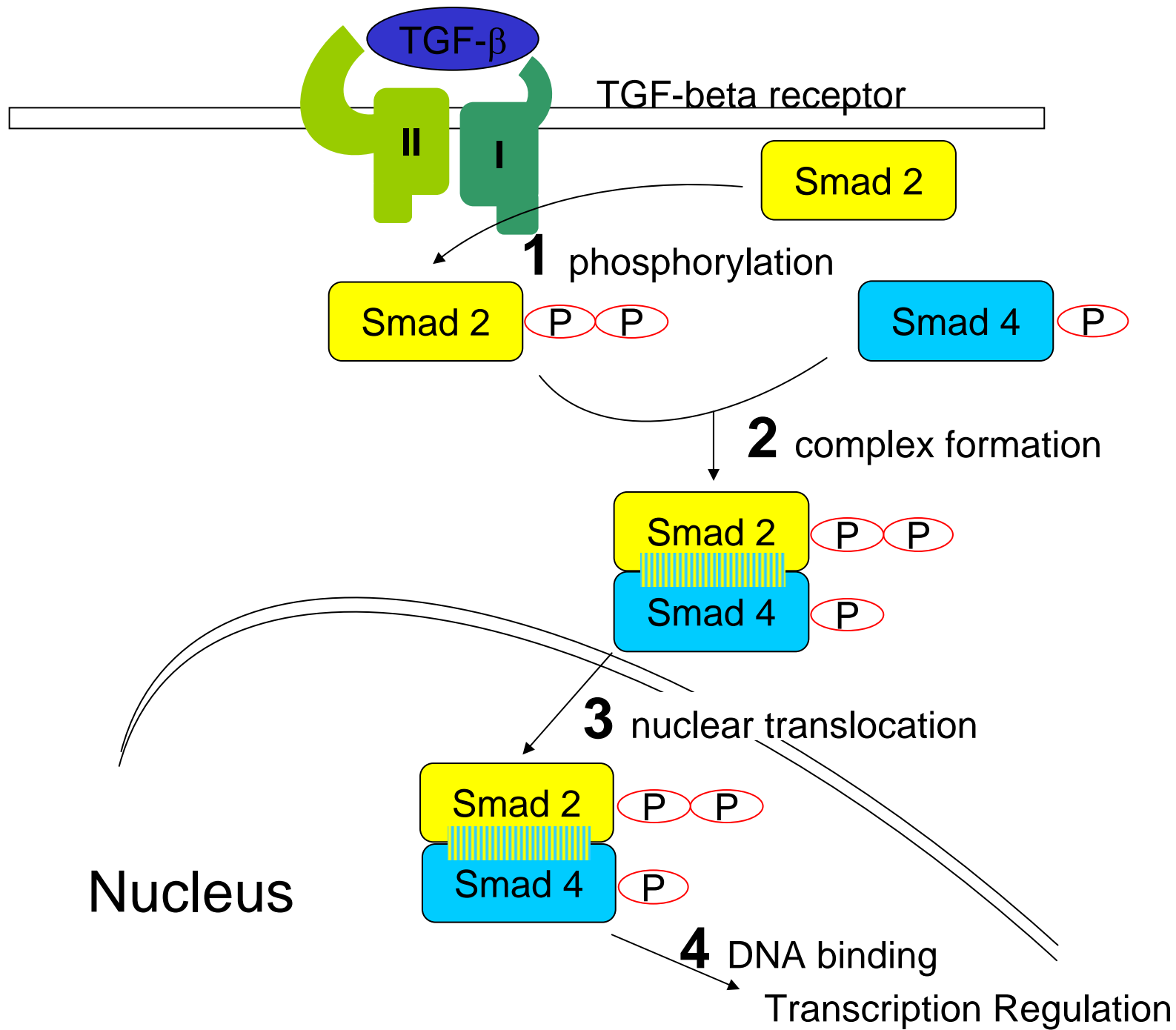
- nucleus

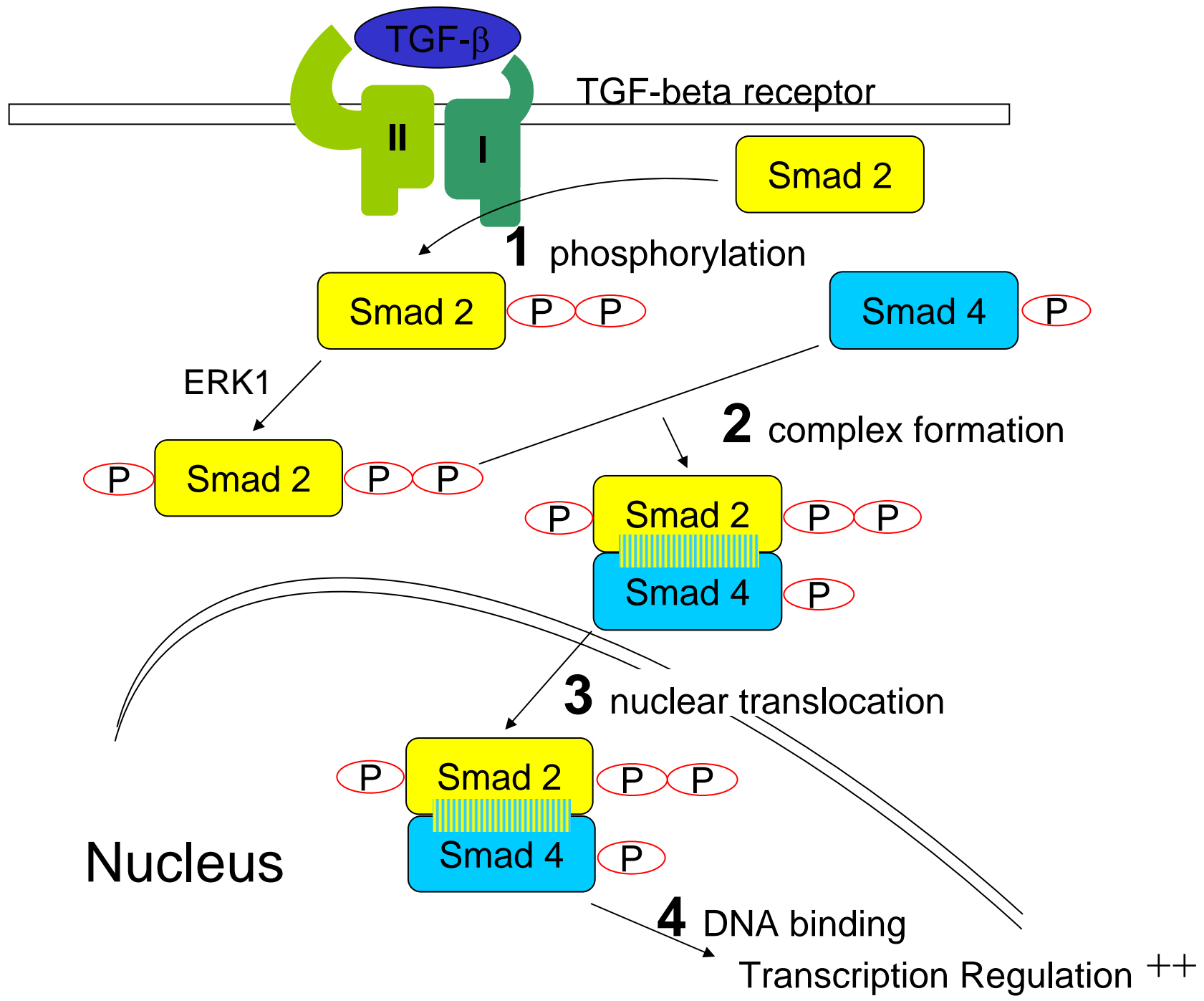
*Molecular Function:*

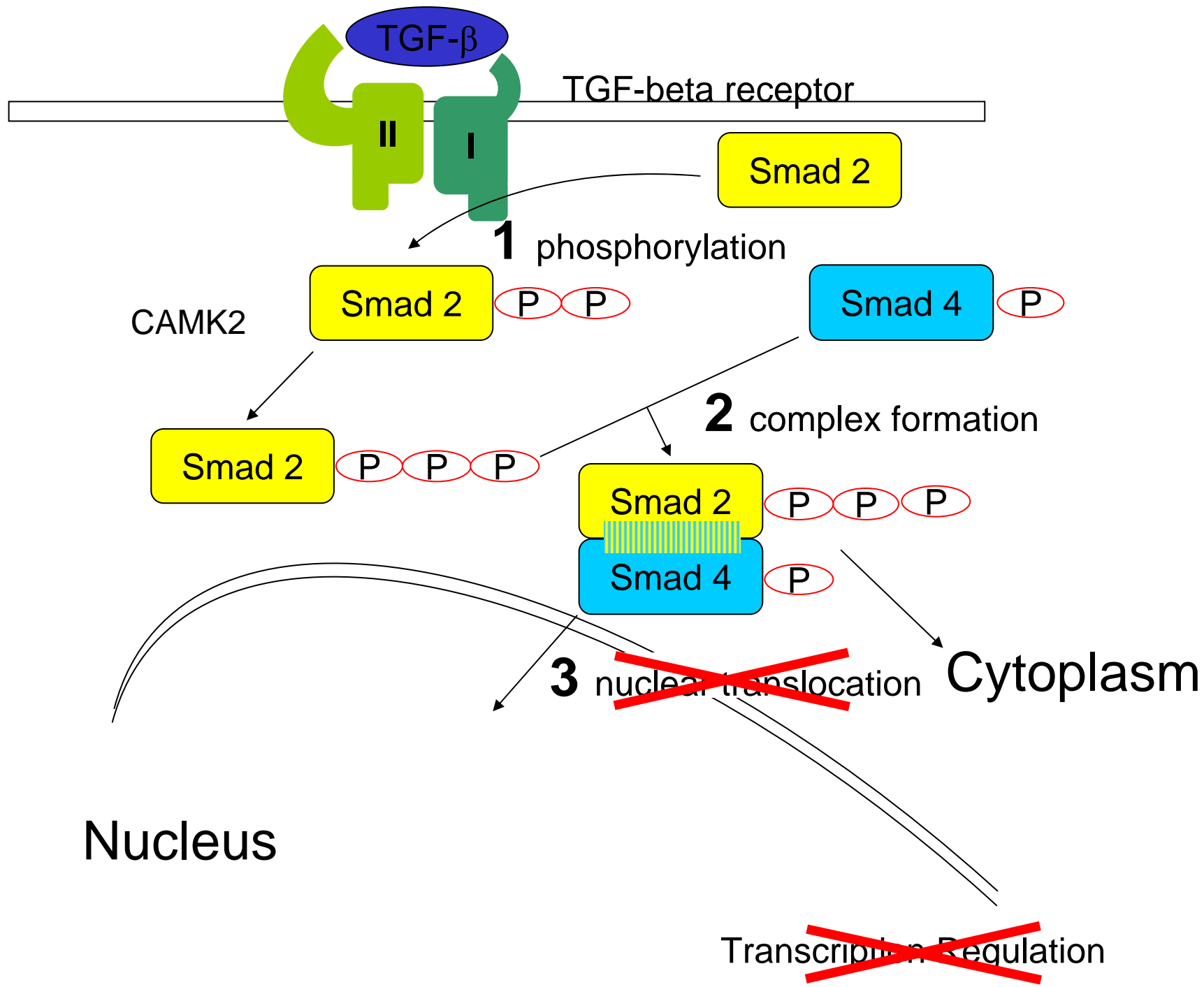
- protein binding



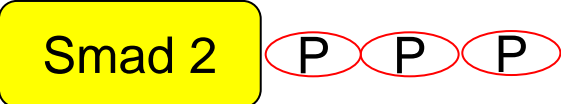
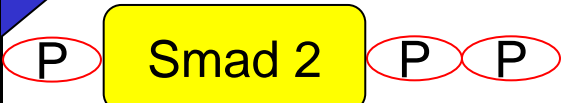



*Biological Process:*

- signal transduction
- regulation of transcription, DNA-dependent







	<p>“normal”</p>	<ul style="list-style-type: none"> <li>•Cytoplasmic</li> </ul>	<p>SMAD2_HUMAN</p>
	<p>TGF-β receptor phosphorylated</p>	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Nuclear</li> <li>•Txn upregulation</li> </ul>	<p>SMAD2_HUMAN</p>
	<p>ERK1 phosphorylated</p>	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Nuclear</li> <li>•Txn upregulation++</li> </ul>	<p>SMAD2_HUMAN</p>
	<p>CAMK2 phosphorylated</p>	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Cytoplasmic</li> <li>•No Txn upregulation</li> </ul>	<p>SMAD2_HUMAN</p>
	<p>alternatively spliced short form</p>	<ul style="list-style-type: none"> <li>•Cytoplasmic</li> </ul>	<p>SMAD2_HUMAN</p>
	<p>phosphorylated short form</p>	<ul style="list-style-type: none"> <li>•Nuclear</li> <li>•Txn upregulation</li> </ul>	<p>SMAD2_HUMAN</p>
	<p>point mutation (causative agent: large intestine carcinoma)</p>	<ul style="list-style-type: none"> <li>•Doesn't form complex</li> <li>•Cytoplasmic</li> <li>•No Txn upregulation</li> </ul>	<p>SMAD2_HUMAN</p>

# Important Considerations

- Need to consider the various forms a protein might take
- Need to provide connections to established ontologies
- Need to account for the possibility that a protein might not share the traits of its parent or siblings



**%PRO:00000010** Smad2

<**PRO:00000011** Smad2 sequence 1 (long form)

>**PRO:00000012** Smad2 sequence 1 phosphorylated form

**%PRO:00000013** Smad2 sequence 1, TGF- $\beta$  receptor I-phosphorylated

**%PRO:00000014** Smad2 sequence 1, TGF- $\beta$  receptor I and ERK1-phosphorylated

**has\_modification MOD:**O-phosphorylated L-serine

**has\_modification MOD:**O-phosphorylated L-threonine

**has\_function GO:** TGF- $\beta$  receptor, pathway-specific cytoplasmic mediator activity

**has\_function GO:**SMAD binding

**has\_function GO:**transcription coactivator activity

**participates\_in GO:**signal transduction

**participates\_in GO:**SMAD protein heteromerization

**participates\_in GO:**regulation of transcription, DNA-dependent

**located\_in GO:**nucleus

**part\_of GO:**transcription factor complex

**%PRO:00000015** Smad2 sequence 1, TGF- $\beta$  receptor I and CAMK2-phosphorylated

<**PRO:00000016** Smad2 sequence 2 (short form) - splice variant

>**PRO:00000017** Smad2 sequence 2 phosphorylated form

**%PRO:00000018** Smad2 sequence 2, TGF- $\beta$  receptor I-phosphorylated

<**PRO:00000019** Smad2 sequence 3 - genetic variant related to colorectal carcinoma

**has\_agent SO:** amino\_acid\_substitution

**lacks\_modification MOD:** phosphorylated residue

**lacks\_function GO:** transcription coactivator activity

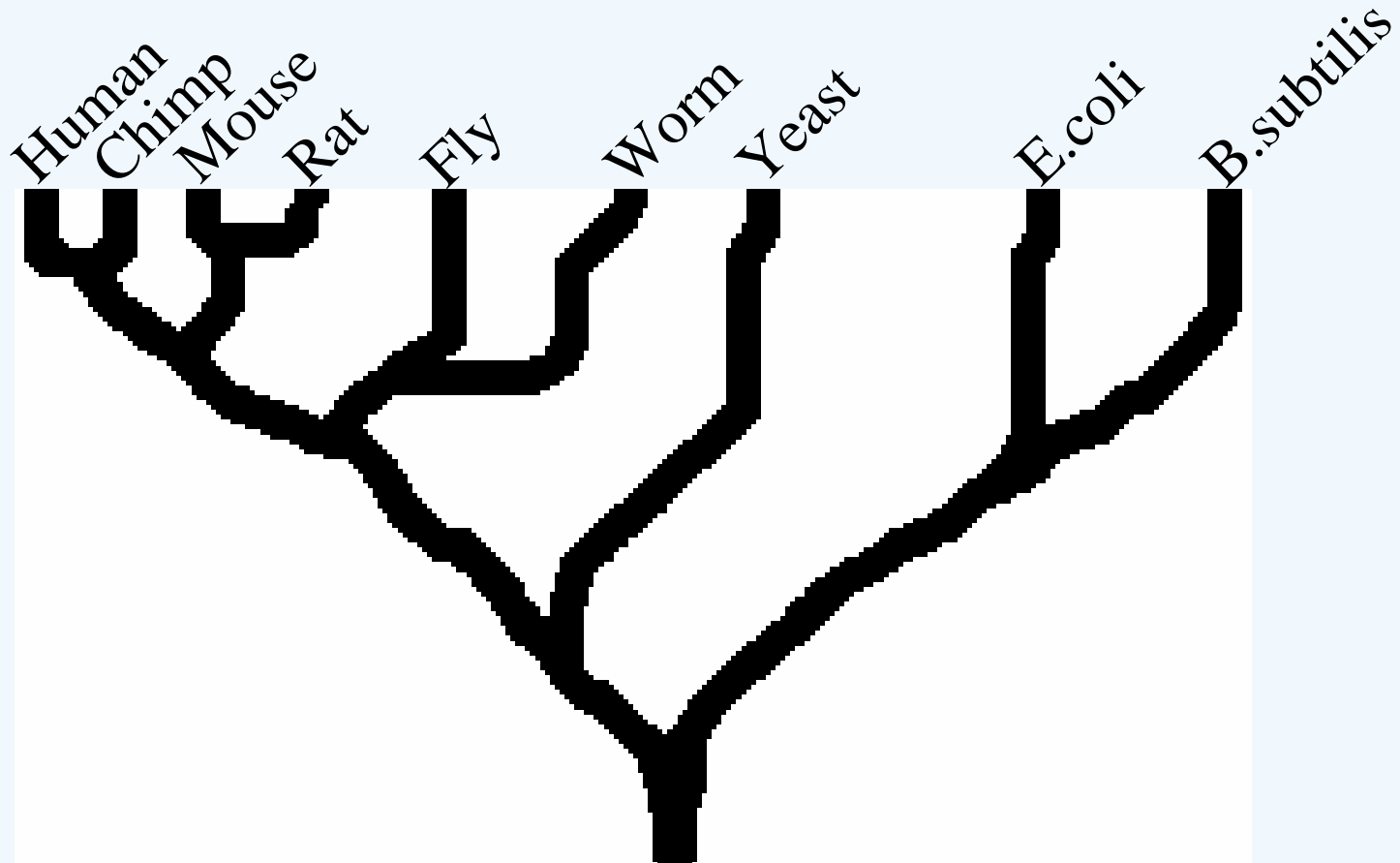
**agent\_of DO:** carcinoma of the large intestine

%	is_a
<	variant_of
>	derives_from

# Important Considerations

- Need to consider the various forms a protein might take
- Need to provide connections to established ontologies
- Need to account for the possibility that a protein might not share the traits of its parent or siblings
- Need to take advantage of model organism data to generate hypotheses about human biology

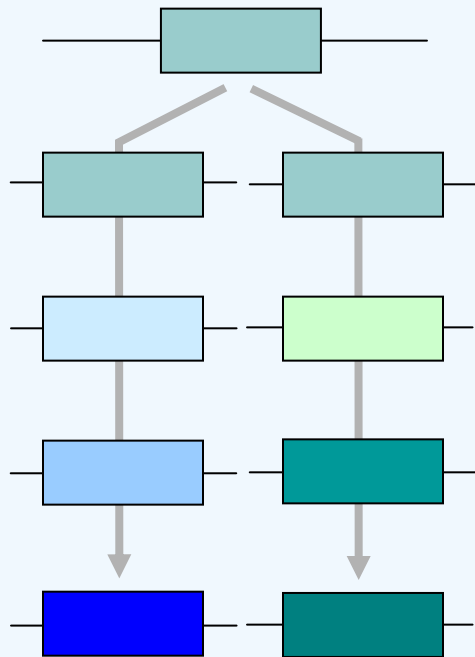
# Implications of Protein Evolution



- Conclusions from experiments performed on proteins from one organism are often applicable to the homologous protein from another organism.
- Information learned about existing proteins allows us to infer the properties of ancestral proteins.

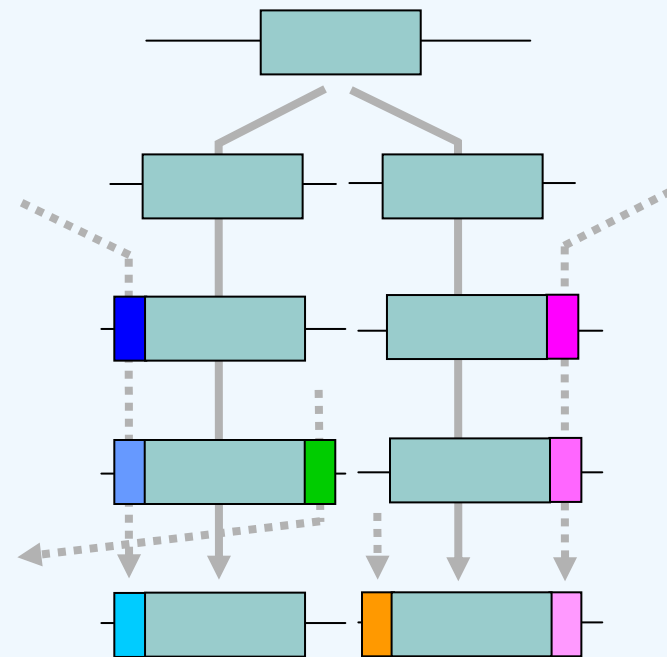
# Protein Evolution

## Sequence changes



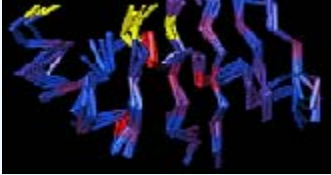

With enough similarity, one can trace back to a common origin

## Domain shuffling



What about these?

# Levels of Protein Classification

<i>Evolutionary Divergence</i>	<i>End-to-End Alignment</i>	<i>Example</i>	<i>Conservation</i>	<i>Database</i>
<b>Very Ancient</b>	<b>Domain structure</b>		<b>Very general biochemical activity</b>	<b>SCOP Superfamily</b>
<b>Ancient</b>	<b>Domain sequence</b>	<pre>CVSGSGNNTSITATGGVVDLQSSSAVKVRSTK CYKSG---IQVRLGEDNINVVEGNEQFISASK *  ...      .:      . :::  ...      :  ::*</pre>	<b>General biochemical activity</b>	<b>Pfam</b>
<b>Recent</b>	<b>Protein sequence Domain architecture</b>		<b>Specific biochemical activity</b>	<b>PIRSF</b>
<b>Functionally Specialized</b>	<b>Protein sequence Domain architecture</b>	<pre>CYKSRIQVRLGEHNIDVLEGNEQFINAAKIIT CYKSGIQVRLGEDNINVVEGNEQFISASKSIV ****  ***** .*:*:***** .*:* *</pre>	<b>Specific biological function</b>	<b>PANTHER</b>

# TGM3 & EPB42

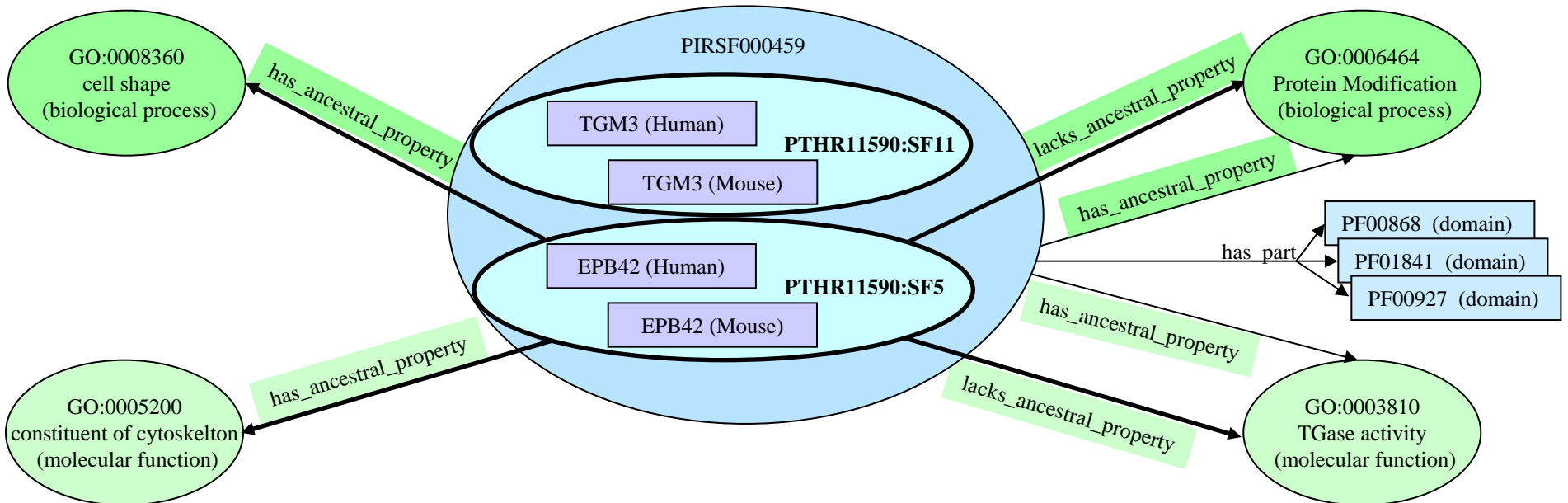
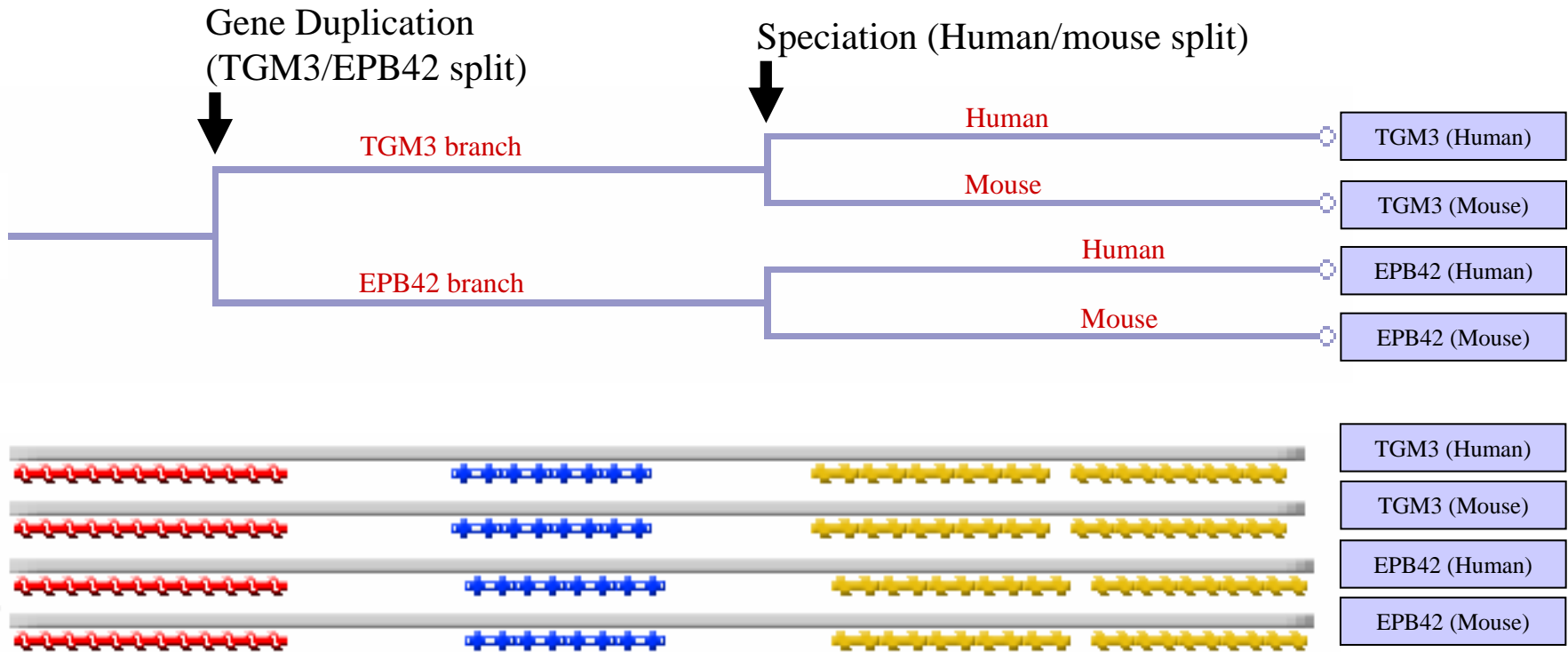
TGM3 = Protein-glutamine gamma-glutamyltransferase  
(Transglutaminase; involved in protein modification)

EBP42 = Erythrocyte membrane protein band 4.2  
(Constituent of cytoskeleton; involved in cell shape)

Q: Are these related?

Q: What is known?

Q: How to capture?



# PRO

Root level

Unit Level

- The two types of evolutionary units
- Not substituted by any other terms

Domain Family Level (structure)

- Related by structural similarity
- Source: SCOP Superfamily

Domain Family Level (sequence)

- Related by sequence similarity
- Source: Pfam domain

Protein Family Level

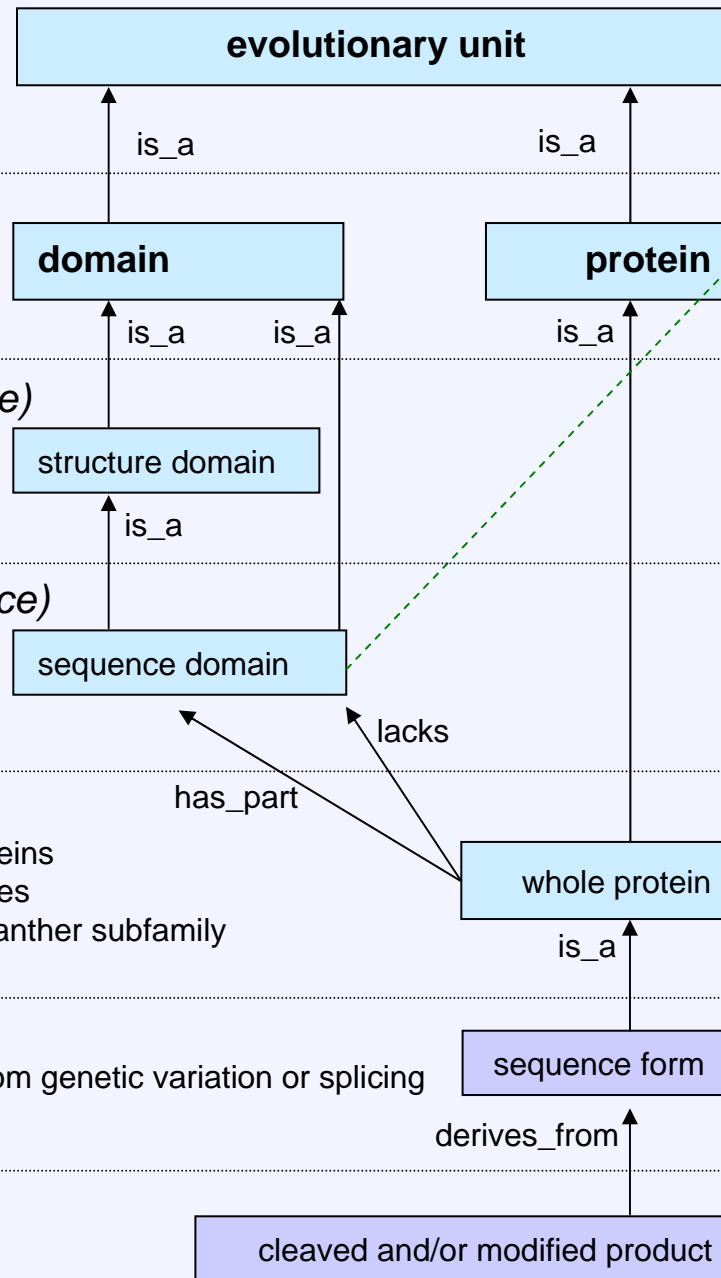
- Evolutionarily-related full-length proteins
- May contain finer-grain sub-categories
- Sources: PIRSF family/subfamily, Panther subfamily

Protein Sequence Level

- Possible sequence forms derived from genetic variation or splicing
- Source: UniProtKB

Protein Modification Level

- Protein as modified after translation
- Source: UniProtKB



# GO

Gene Ontology

molecular function

- has\_ancestral\_property
- has\_function
- lacks\_function

biological process

- has\_ancestral\_property
- participates\_in

cellular component

- has\_ancestral\_property
- part\_of (for complexes)
- located\_in (for compartments)

# DO/UMLS

Disease

disease

- agent\_of

# SO

Sequence Ontology

sequence changes

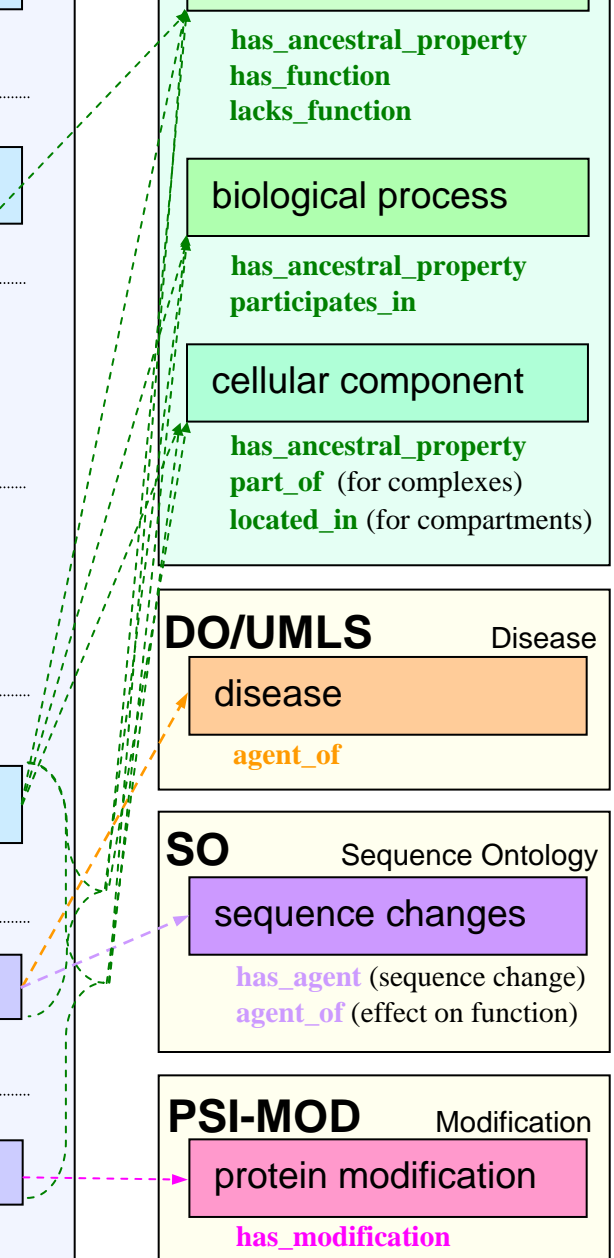
- has\_agent (sequence change)
- agent\_of (effect on function)

# PSI-MOD

Modification

protein modification

- has\_modification





# PRO Team (so far...)

## •Principle Investigators

Cathy Wu	(PIR at GUMC)
Judith Blake	(The Jackson Laboratory)
Hongfang Liu	(GUMC)
Barry Smith	(SUNY Buffalo)

## •Curators

Darren Natale	(PIR at GUMC)
Cecilia Arighi	(PIR at GUMC)
Winona Barker	(PIR at GUMC)
Zhang-zhi Hu	(PIR at GUMC)