

# The Protein Ontology: An Evolution

251<sup>st</sup> ACS National Meeting & Exposition

Chemistry, Data & the Semantic Web: An Important Triple to Advance Science

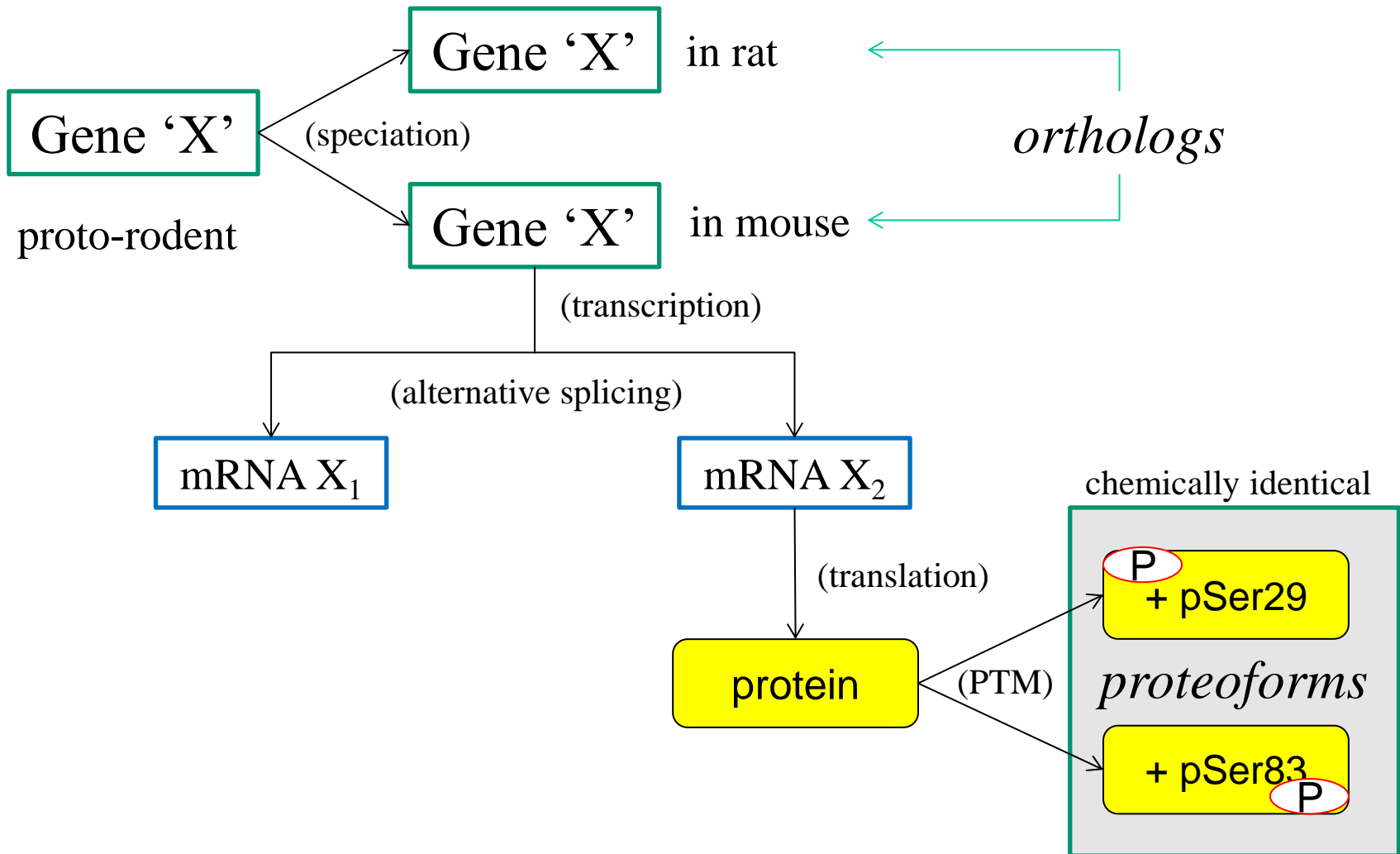
San Diego, CA

Darren A. Natale, Ph.D.

Protein Science Team Lead, PIR

Research Assistant Professor, GUMC

# Some Background





# Why create PRO: Provide precise targets for annotation

Mothers against decapentaplegic homolog 2

Smad 2

GO annotation of SMAD2\_HUMAN:

*Cellular Component:*

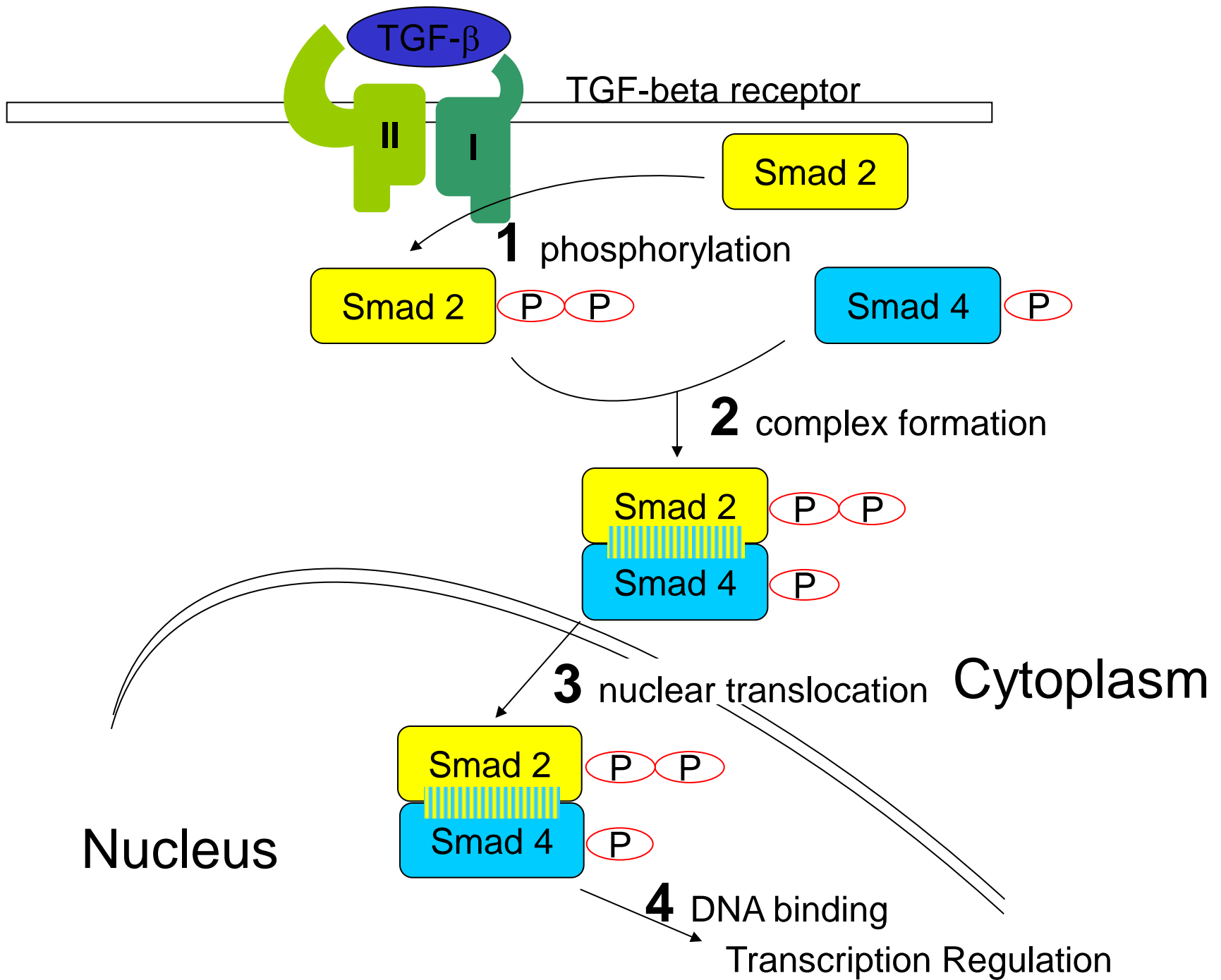
- nucleus

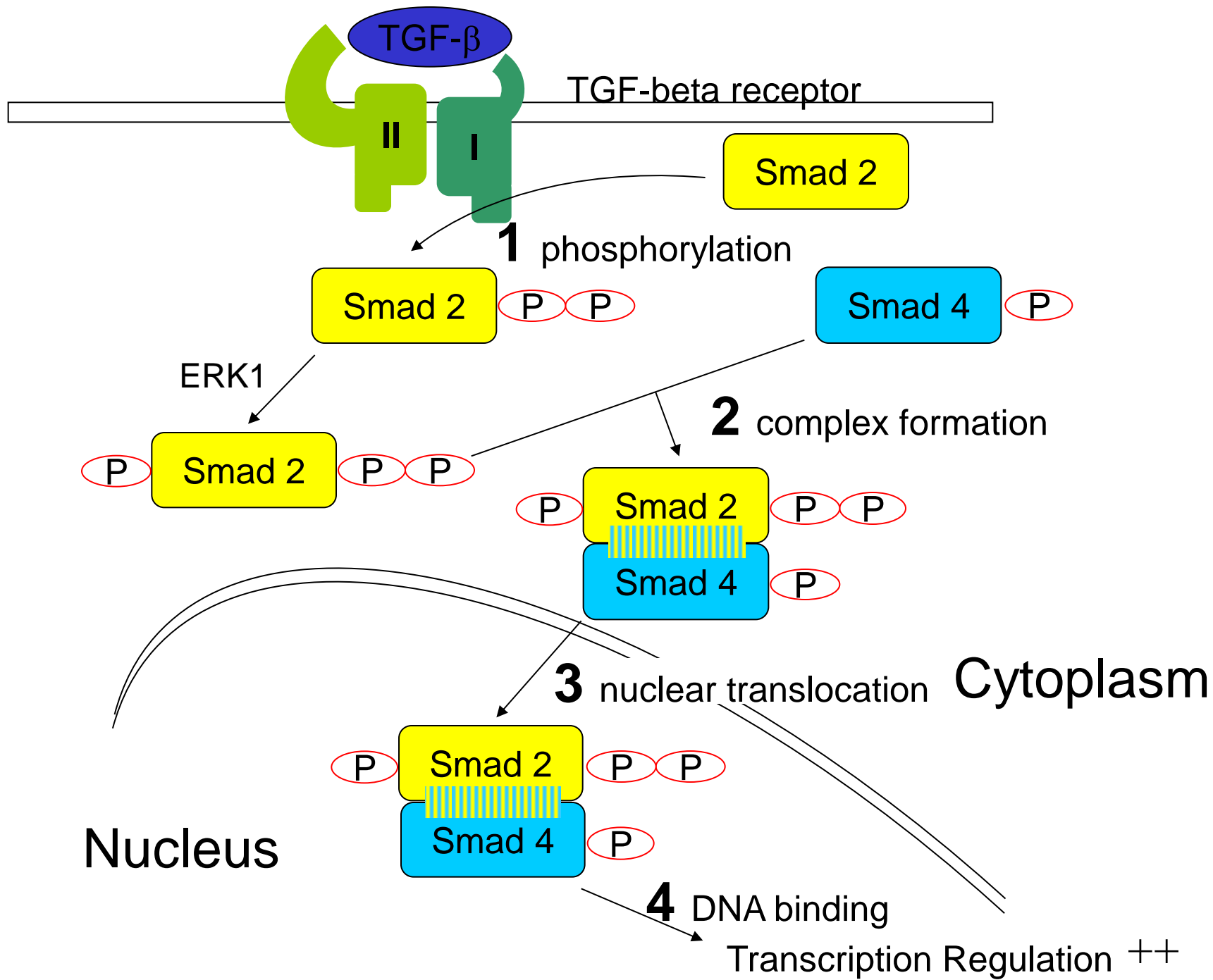
*Molecular Function:*

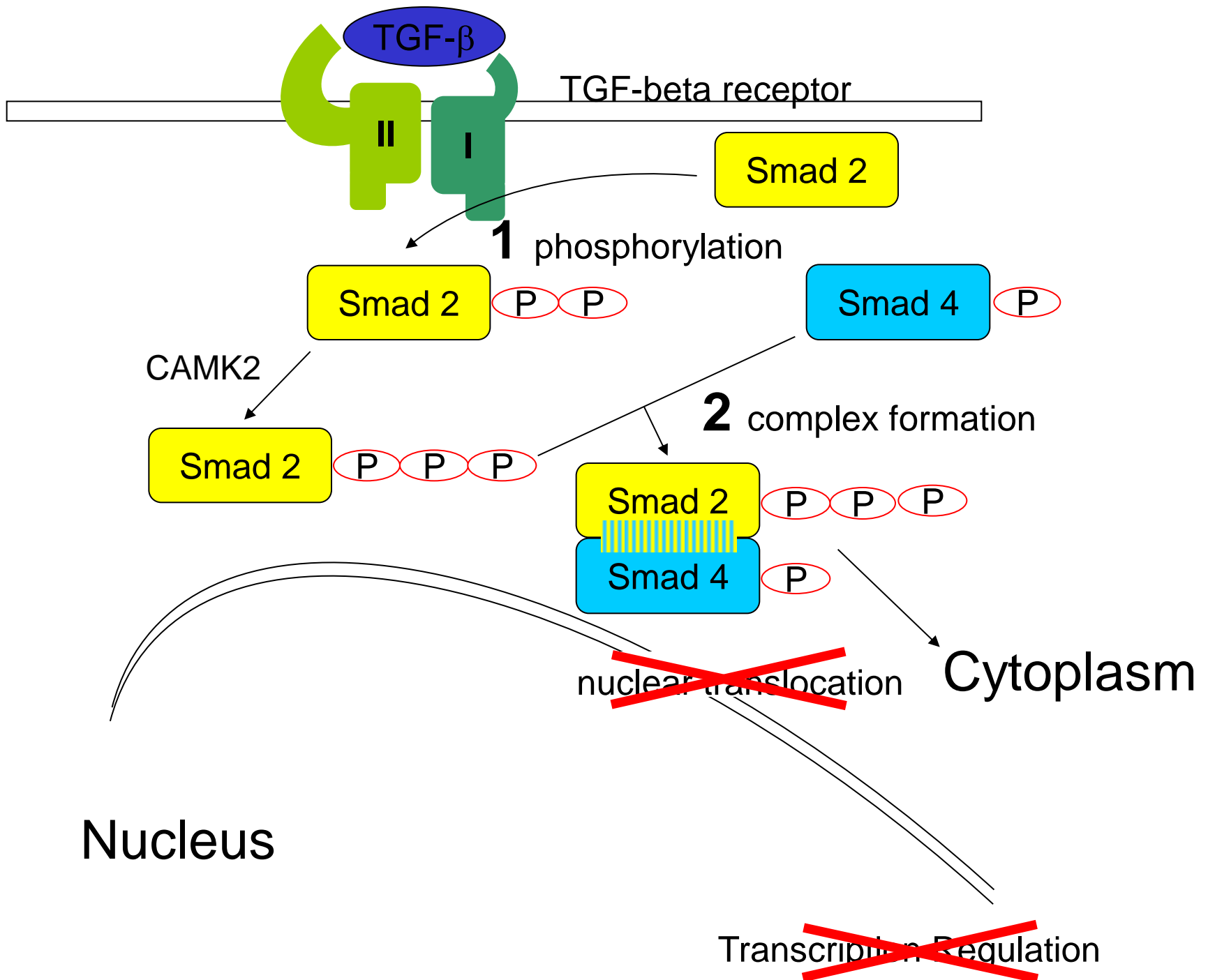
- protein binding

*Biological Process:*







- signal transduction
- regulation of transcription, DNA-dependent















Smad 2	"normal"	<ul style="list-style-type: none"> <li>•Cytoplasmic</li> </ul>	SMAD2_HUMAN
Smad 2 	TGF-β receptor phosphorylated	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Nuclear</li> <li>•Txn upregulation</li> </ul>	SMAD2_HUMAN
Smad 2 	ERK1 phosphorylated	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Nuclear</li> <li>•Txn upregulation++</li> </ul>	SMAD2_HUMAN
 Smad 2 	CAMK2 phosphorylated	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Cytoplasmic</li> <li>•No Txn upregulation</li> </ul>	SMAD2_HUMAN
Smad 2	alternatively spliced short form	<ul style="list-style-type: none"> <li>•Cytoplasmic</li> </ul>	SMAD2_HUMAN
Smad 2 	phosphorylated short form	<ul style="list-style-type: none"> <li>•Nuclear</li> <li>•Txn upregulation</li> </ul>	SMAD2_HUMAN
Smad 2 	point mutation (causative agent: large intestine carcinoma)	<ul style="list-style-type: none"> <li>•Doesn't form complex</li> <li>•Cytoplasmic</li> <li>•No Txn upregulation</li> </ul>	SMAD2_HUMAN



Smad 2	"normal"	<ul style="list-style-type: none"> <li>•Cytoplasmic</li> </ul>	PR:00000468
Smad 2 	TGF-β receptor phosphorylated	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Nuclear</li> <li>•Txn upregulation</li> </ul>	PR:00000650
Smad 2 	ERK1 phosphorylated	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Nuclear</li> <li>•Txn upregulation++</li> </ul>	PR:00000651
 Smad 2 	CAMK2 phosphorylated	<ul style="list-style-type: none"> <li>•Forms complex</li> <li>•Cytoplasmic</li> <li>•No Txn upregulation</li> </ul>	PR:00000652
Smad 2	alternatively spliced short form	<ul style="list-style-type: none"> <li>•Cytoplasmic</li> </ul>	PR:00000469
Smad 2 	phosphorylated short form	<ul style="list-style-type: none"> <li>•Nuclear</li> <li>•Txn upregulation</li> </ul>	PR:00000656
Smad 2 	point mutation (causative agent: large intestine carcinoma)	<ul style="list-style-type: none"> <li>•Doesn't form complex</li> <li>•Cytoplasmic</li> <li>•No Txn upregulation</li> </ul>	PR:00000470

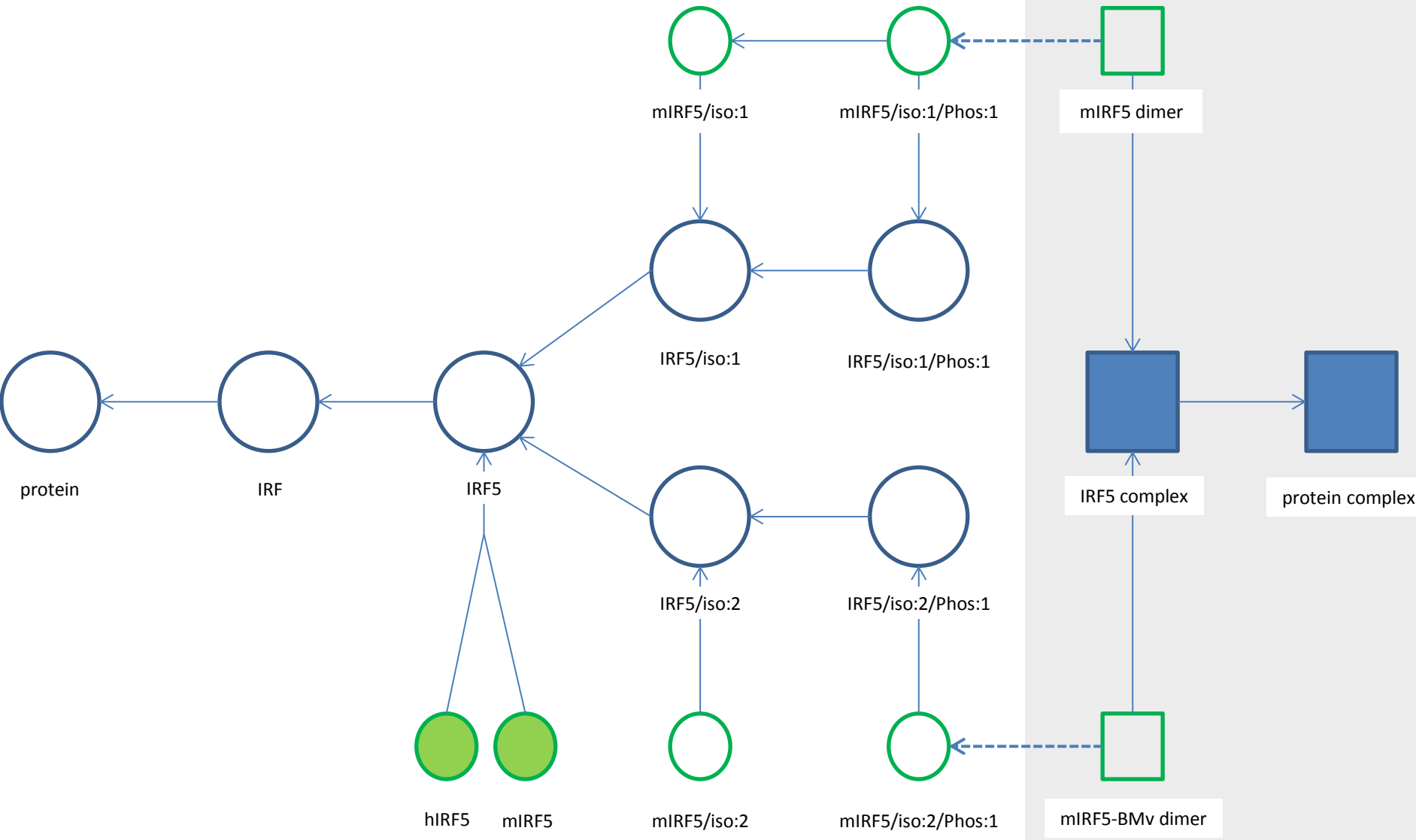
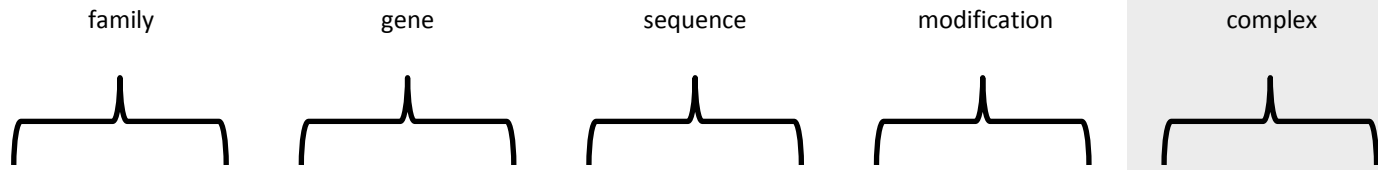


# Sometimes being precise means you need *broader* categories

BioMarker	Alternative Names	Clone	PRO	Scope
CD45	PTPRC, B220, gp180, LCA, leukocyte common antigen, Ly-5, Lyt-4, lymphocyte antigen 5, T200	HI30	PR:000001006	All forms of CD45 (human+)
CD45RA		HI100	PR:000001015	isoform RA (human+)
CD45RB		MT4 (6B6)	PR:000001016	isoform RB (human+)
CD45RO		UCHL1	PR:000001017	isoform RO (human+)
Src	SRC1, ASV, c-src	3828	PR:P12931	All forms of Src (human)
Src (pY418)		K98-37	PR:000027237	Any Src phosphorylated on Y418 (human)

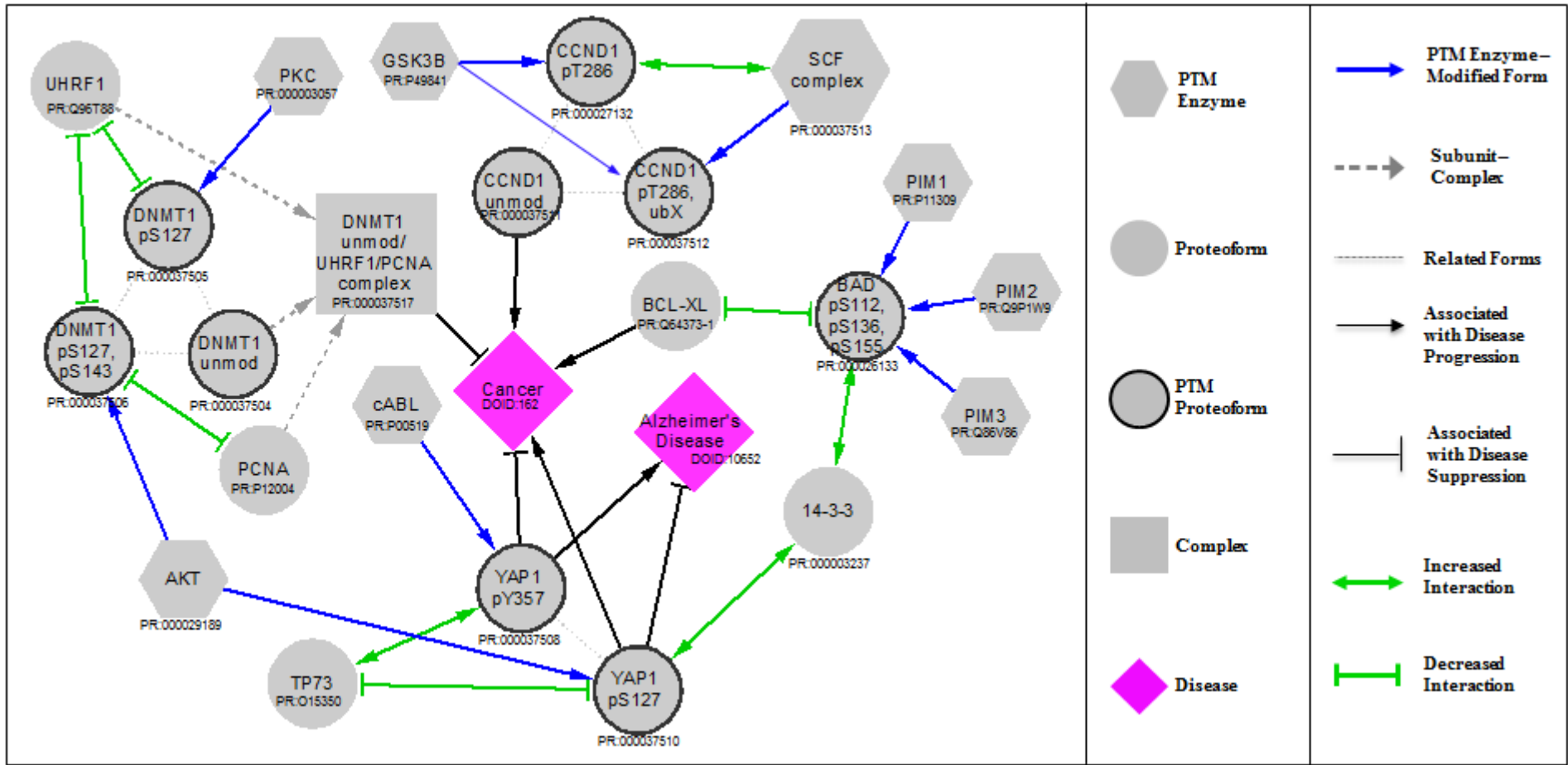
## “Smad2 Is Essential for Maintenance of the Human and Mouse Primed Pluripotent Stem Cell State”

- some exp'ts done using hSmad2 (PR:Q15796)
- other exp'ts done using mSmad2 (PR:Q62432)
- conclusion meant to apply to both (PR:000000364)



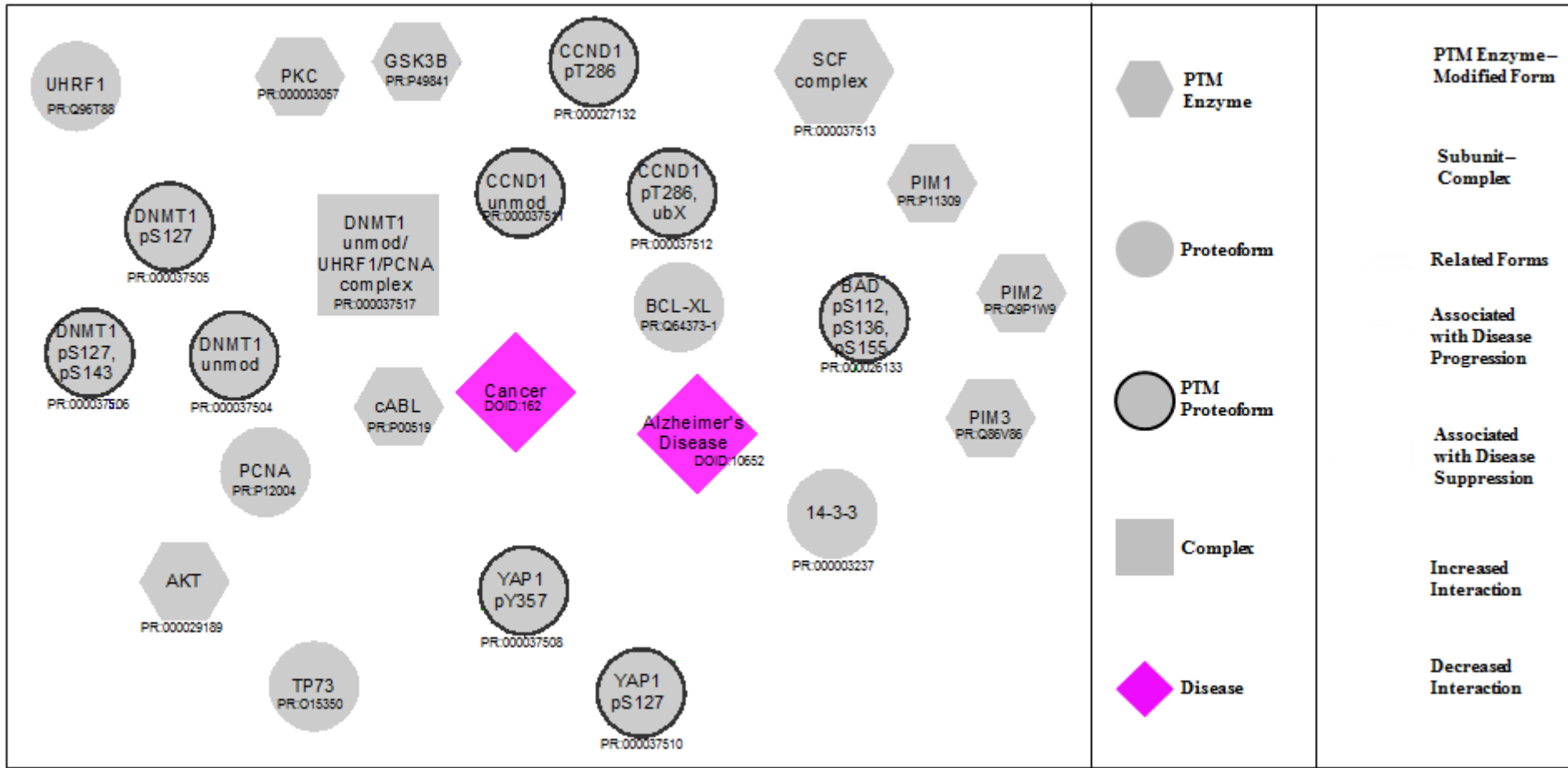


# But Why An Ontology?





# But Why An Ontology?



Ontology provides opportunity for rich relations









# Reasoning...

Examine statements to (automatically) find new relationships or detect conflicts





# The Benefit, Demonstrated

human enzymatic glycoprotein

- = *is\_a* protein
- + *capable\_of* catalytic activity
- + *has\_part* glycan
- + *in\_taxon* *Homo sapiens*

*hELANE is\_a* protein ✓

*hELANE capable\_of* peptidase activity (→ *is\_a* catalytic...) ✓

~~*hELANE has\_part* GlcNAc (→ *is\_a* glycan)~~

*hELANE in\_taxon* *Homo sapiens* ✓

**What if it's discovered that *hELANE* isn't glycosylated?**



# The Benefit, Demonstrated

human enzymatic glycoprotein

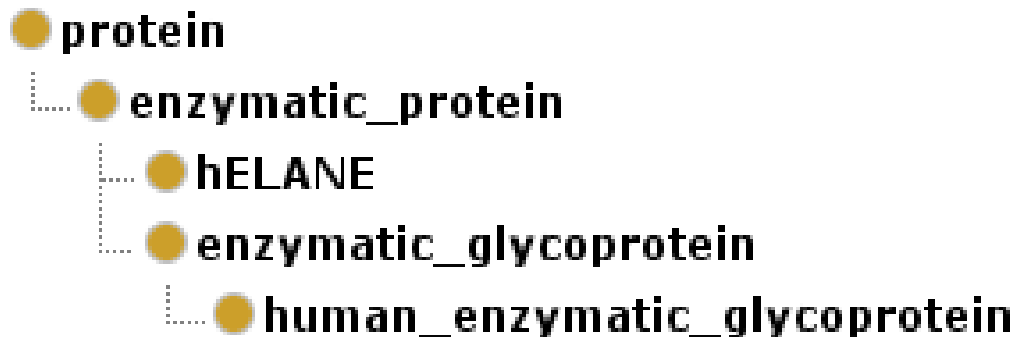
- = *is\_a* protein
- + *capable\_of* catalytic activity
- + *has\_part* glycan
- + *in\_taxon* *Homo sapiens*

*hELANE is\_a* protein ✓

*hELANE capable\_of* peptidase activity (→ *is\_a* catalytic...) ✓

~~*hELANE has\_part* GlcNAc (→ *is\_a* glycan)~~

*hELANE in\_taxon* *Homo sapiens* ✓





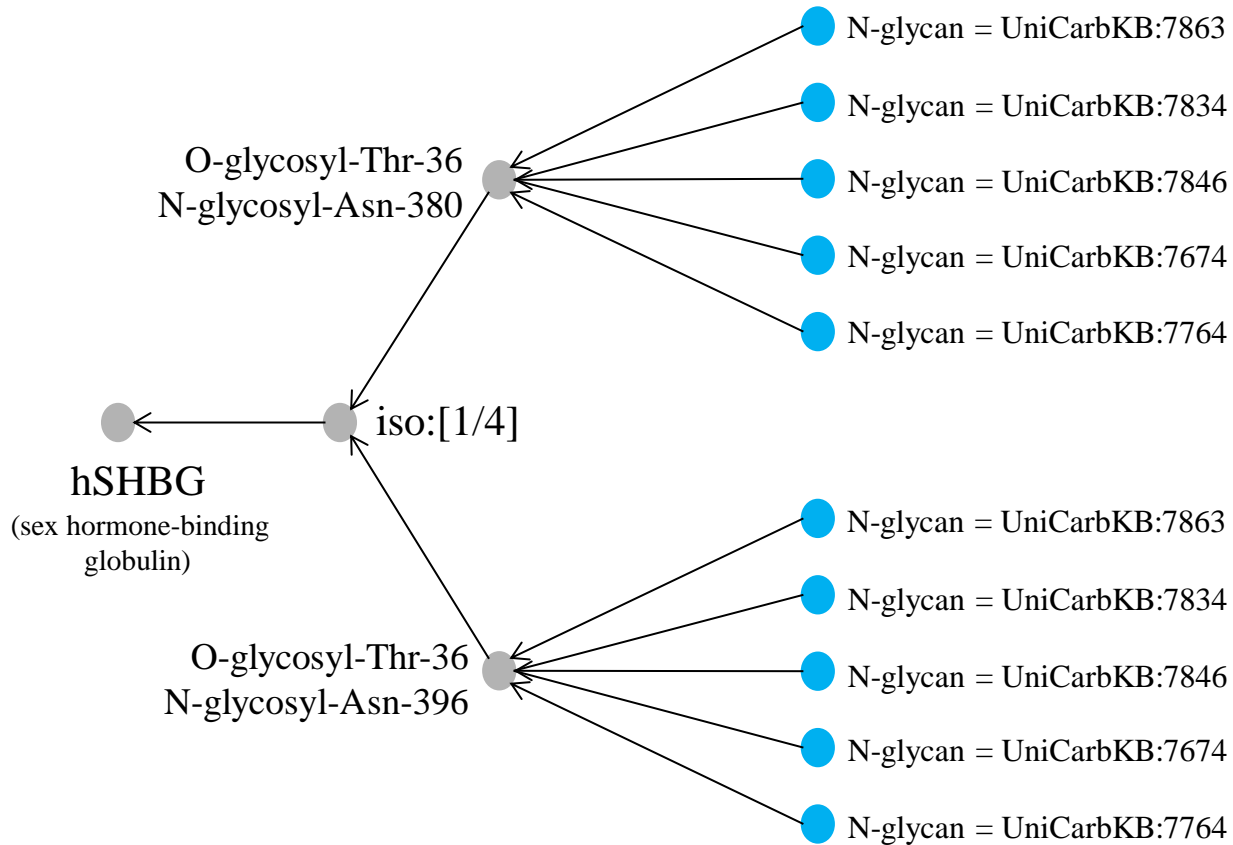


# Lessons Learned...

- It is better to make statements on a more elemental level than make statements that are more complex
  - can add and remove as needed
  - easier to understand
  - enables reasoning
- It is better to make statements about evidence rather than make statements about conclusions
  - can import needed information
- It is beneficial to note mutually-exclusive conditions
  - helps flag inconsistent cases



# Future Directions: Glycobiology



“Oh the pain, the pain of it all”



- Reaching consensus
- Make it right, or make it useful
- Improvement vs content completeness



# The PRO Consortium

## Collaborating Institutions:



Protein Information Resource (PIR)



The Jackson Laboratory



NYS Center of Excellence  
in Bioinformatics and Life Sciences



Reactome

## Principle Investigators:

Cathy Wu, Director of PIR

Judith Blake, Associate Professor, Jackson Laboratory

Barry Smith, Professor, University at Buffalo

Peter D'Eustachio, Associate Professor, NYU Med Ctr

## Funding & Sponsors:

PRO is funded by NIGMS / NIH grants 1R01GM080646-01,  
3R01GM080646-04S1, 3R01GM080646-04S2, 5R01GM080646-05

