



# EXPLAINING PROTEIN- PROTEIN INTERACTION PREDICTIONS WITH GENETIC PROGRAMMING

evo  2022

Rita T. Sousa, Sara Silva, Catia Pesquita  
LASIGE, Faculdade de Ciências da Universidade de Lisboa

WHY DO WE NEED  
**EXPLAINABLE**  
ARTIFICIAL INTELLIGENCE?

# WHY DO WE NEED **EXPLAINABLE** ARTIFICIAL INTELLIGENCE?



Users' trust

# WHY DO WE NEED **EXPLAINABLE** ARTIFICIAL INTELLIGENCE?



Users' trust



Ensure  
algorithmic  
fairness

# WHY DO WE NEED **EXPLAINABLE** ARTIFICIAL INTELLIGENCE?



Users' trust



Ensure  
algorithmic  
fairness



Identification  
of potential  
bias

# WHY DO WE NEED **EXPLAINABLE** ARTIFICIAL INTELLIGENCE?



Users' trust



Ensure  
algorithmic  
fairness

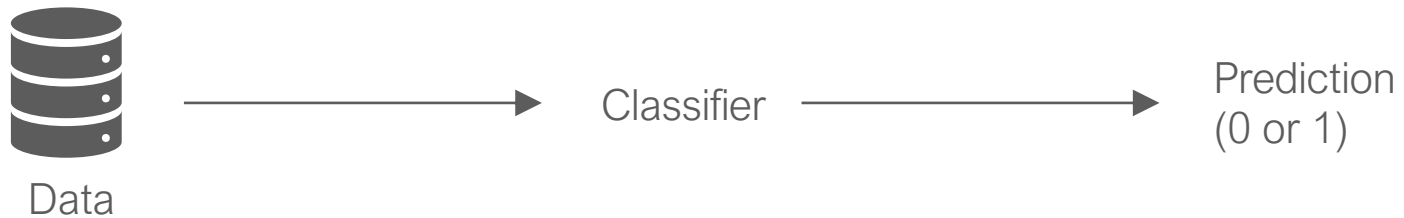
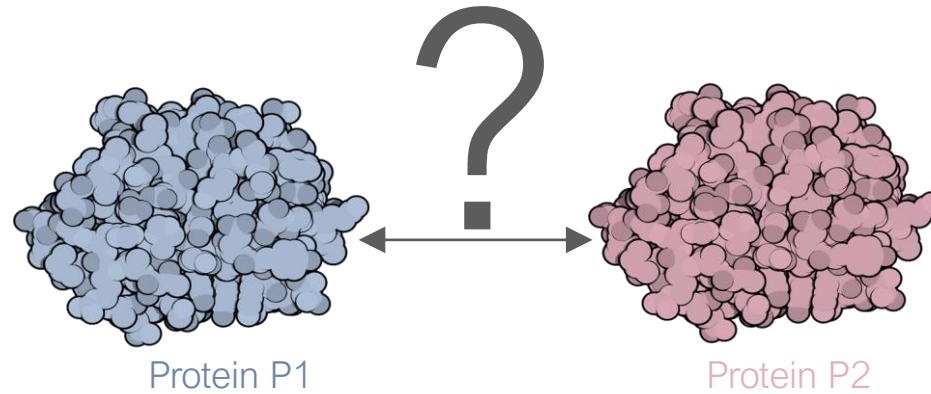


Identification  
of potential  
bias



Discovery  
of new  
knowledge

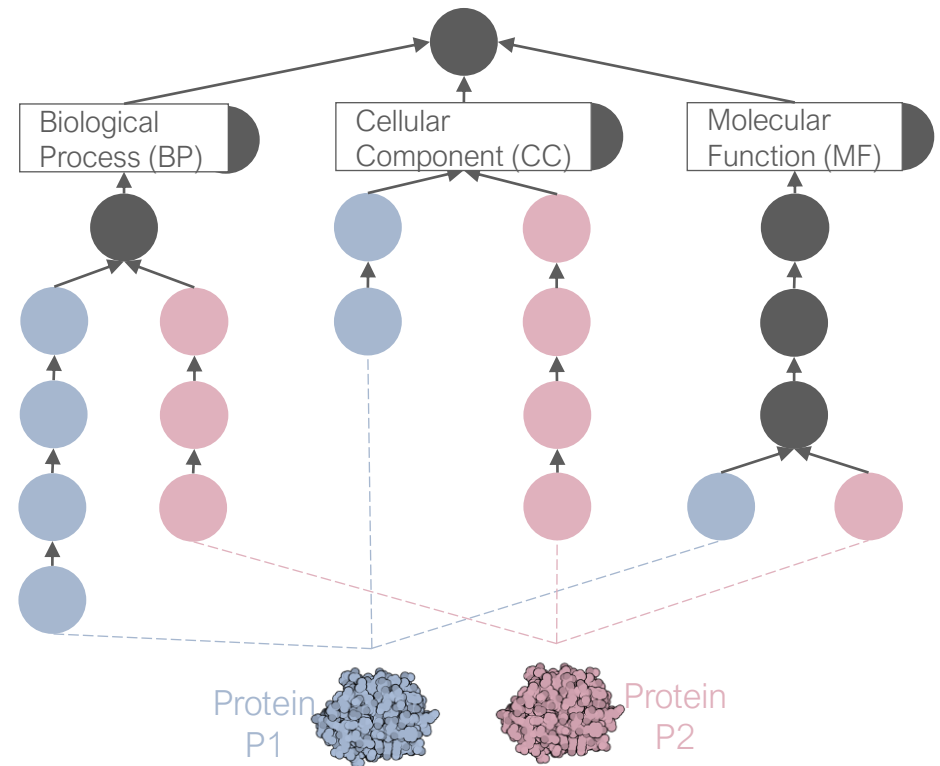
# PROTEIN-PROTEIN INTERACTION (PPI) PREDICTION



However, experts want to understand the biological mechanisms that underlie the natural phenomena they are predicting.

# ONTOLOGIES AND KNOWLEDGE GRAPHS ARE A UNIQUE OPPORTUNITY FOR EXPLAINABILITY

Ontologies and Knowledge Graphs (KGs) provide semantics (i.e., meaning) to the entities they represent through different semantic aspects.



GO KG = GO + GO annotations



# KNOWLEDGE GRAPH EMBEDDINGS ARE NOT EXPLAINABLE BY DEFAULT

An embedding is a vector representation that maps each node to a lower-dimensional space in which its graph position and the structure of its local graph neighborhood are preserved.

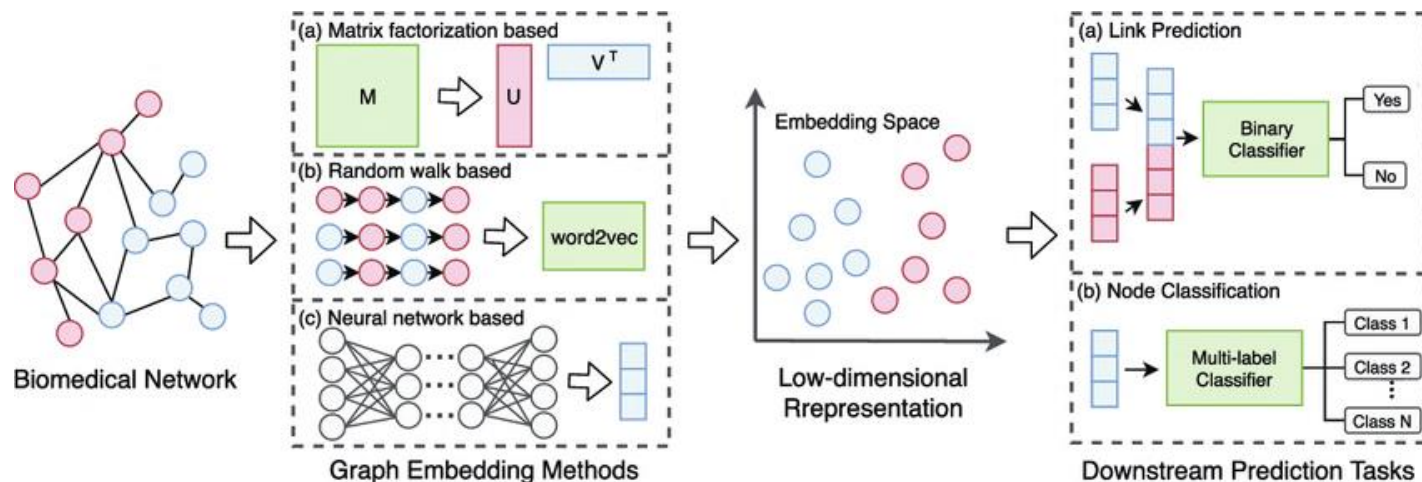
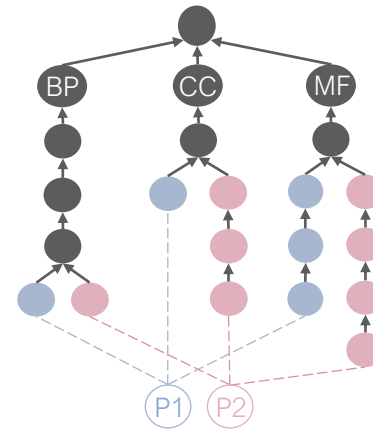


Image credit: Yue, X., Wang, Z., Huang, J., Parthasarathy, S., Moosavinasab, S., Huang, Y., ... & Sun, H. (2020). Graph embedding on biomedical networks: methods, applications and evaluations. *Bioinformatics*, 36(4), 1241-1251.

# USING SEMANTIC SIMILARITY AS AN ALTERNATIVE EXPLANATORY STRATEGY

- KGs provide the scaffolding for comparing entities at a higher level of complexity by comparing the ontology classes with which they are annotated.
- Semantic similarity computed using different portions of the KG to reflect different semantic aspects (SA) can provide more granular explanations with higher information content.



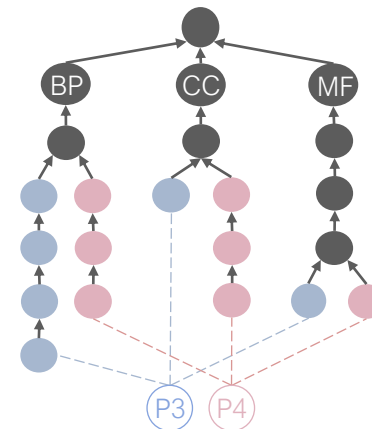
P1 is vacuolar protein-sorting-associated protein 36.  
P2 is vacuolar-sorting protein SNF8.

P1 and P2 interact.

$$SS_{BP} = 0.793$$

$$SS_{CC} = 0.326$$

$$SS_{MF} = 0.061$$



P3 is eukaryotic translation initiation factor 5B, isoform B.  
P4 is Mig-2-like GTPase Mtl.

P3 and P4 do not interact.

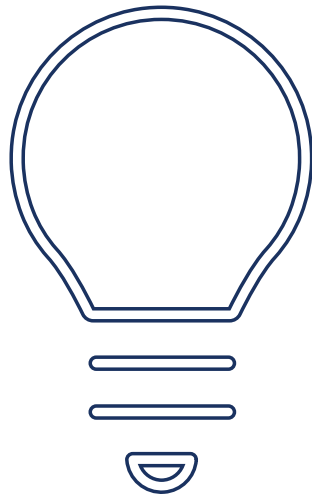
$$SS_{BP} = 0.006$$

$$SS_{CC} = 0.202$$

$$SS_{MF} = 0.713$$

HOW CAN WE GENERATE  
**GLOBAL AND INTERPRETABLE**  
EXPLANATION FOR PPI PREDICTION USING  
SEMANTIC SIMILARITY AS REPRESENTATION?

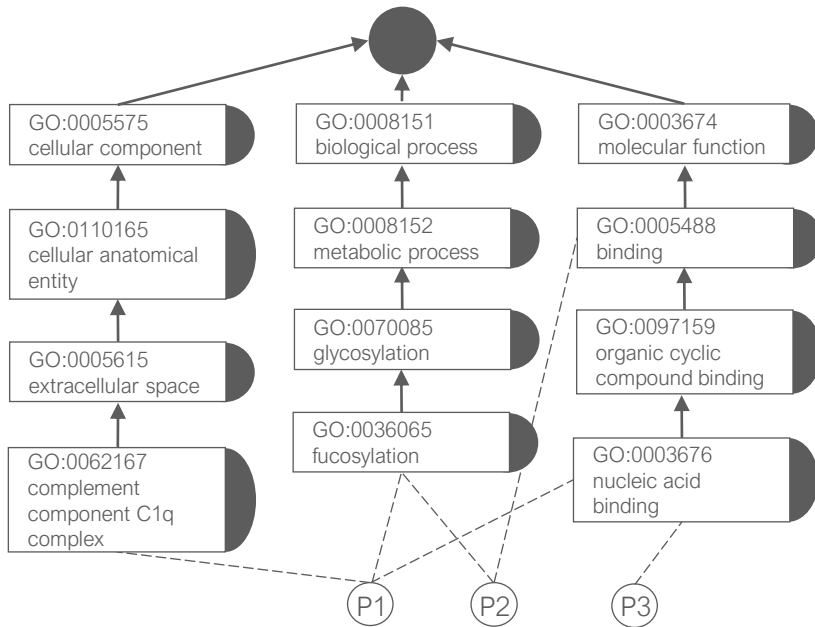
# HOW CAN WE GENERATE GLOBAL AND INTERPRETABLE EXPLANATION FOR PPI PREDICTION USING SEMANTIC SIMILARITY AS REPRESENTATION?



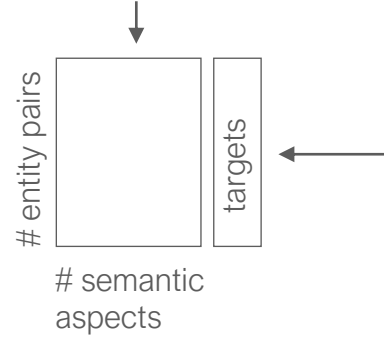
Using [Genetic Programming \(GP\)](#) over a set of semantic similarity values, each describing a semantic aspect represented in the knowledge graph.

# OUR APPROACH

Gene Ontology KG



Computing KG-based semantic similarity between entity pairs for each semantic aspect



PPI from STRING Database



+

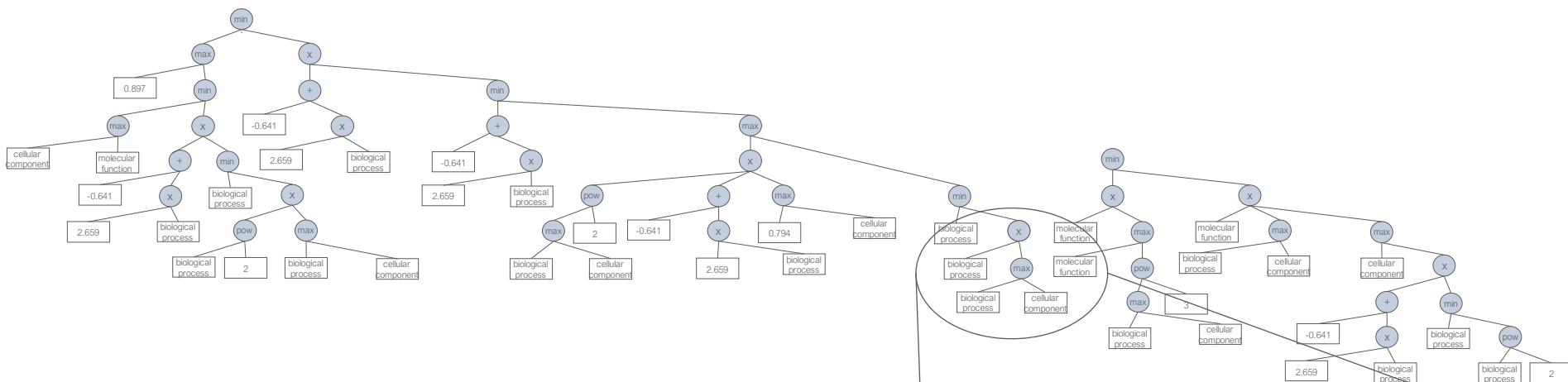
Random Negative Sampling



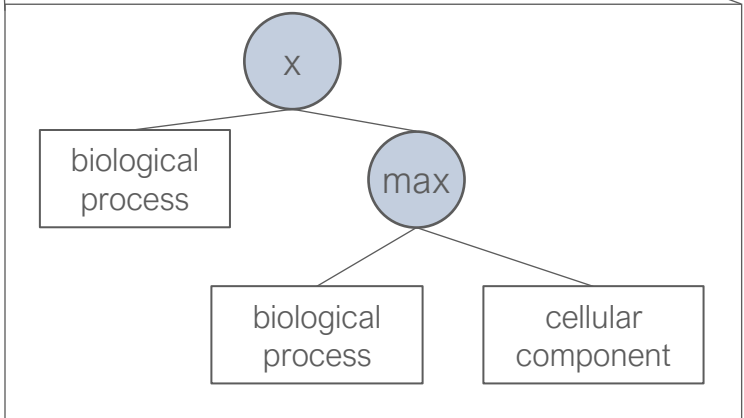
Evolving a GP model

Predicting on unseen data using the GP model

# GP IS NOT ALWAYS INTERPRETABLE



- The solutions may grow exponentially with each generation, and the interpretability is lost.
- Some operators such as multiplication and division are not interpretable in the biological context.



# GP VERSUS GP6X

## GP

- No depth penalization
- 6 operators: multiplication, division, maximum, minimum, addition and subtraction

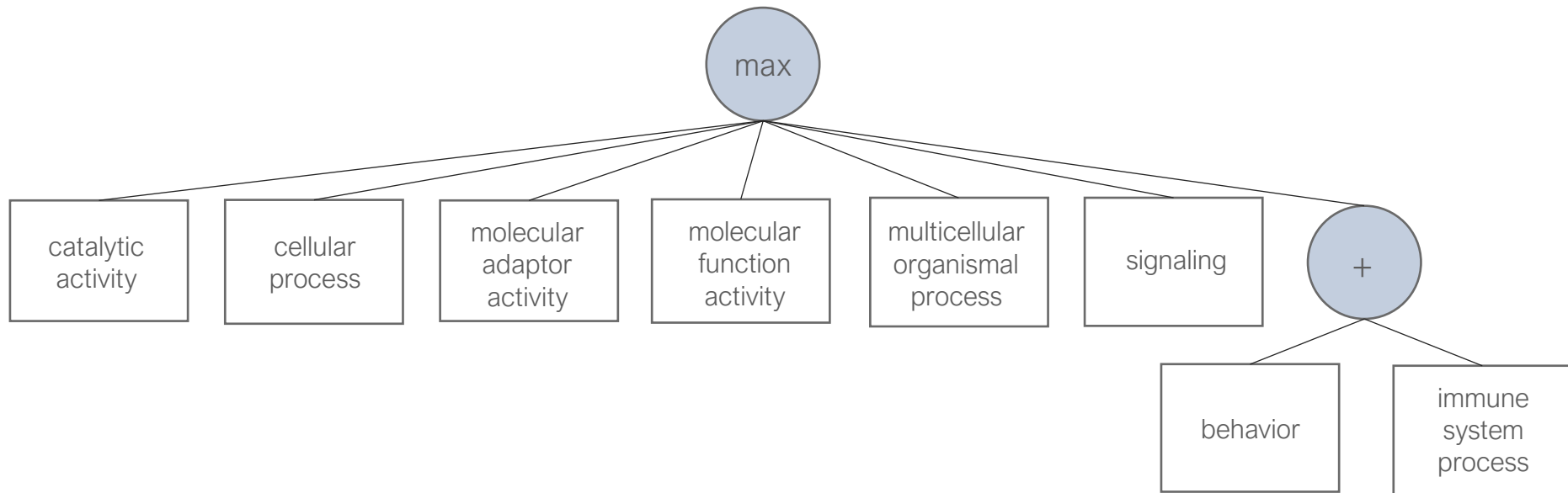
## GP6x

- Penalize solutions with a depth greater than 6
- Interpretable operators, namely maximum, minimum, addition and subtraction

Results for 10-fold cross-validation :

	Median weighted average of F-measures (WAF)	Median number of nodes
GP	0.875	49
GP6x	0.866	17

# GP6X MODEL ANALYSIS



$$\max(SS_{\text{catalytic activity}}, SS_{\text{cellular process}}, SS_{\text{molecular adaptor activity}}, SS_{\text{molecular function regulator}}, SS_{\text{multicellular organismal process}}, SS_{\text{signaling}}, SS_{\text{behavior}} + SS_{\text{immune system process}})$$

Two proteins that interact usually participate in the same biological processes.

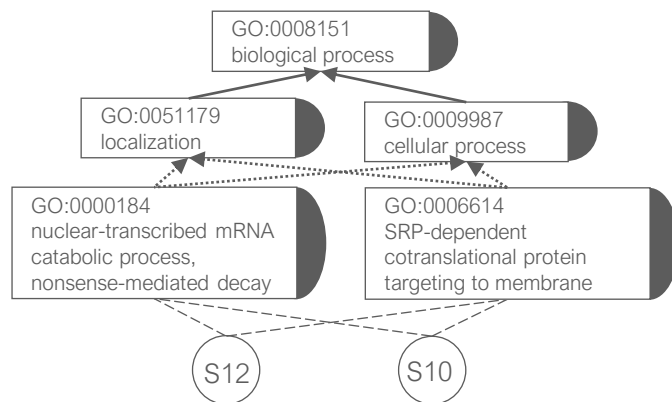
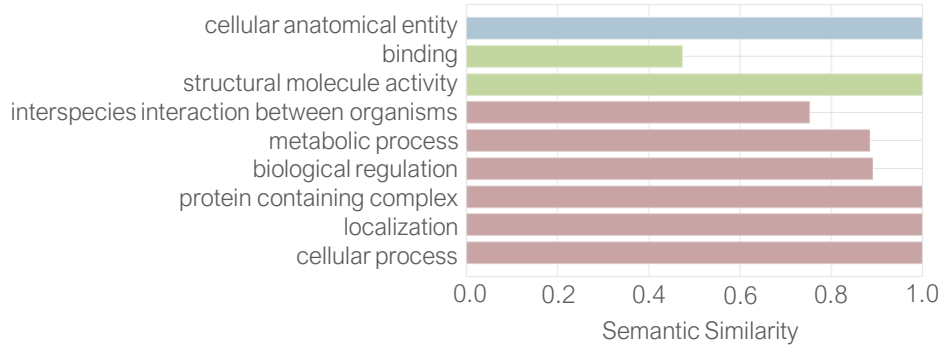


# WHEN GP CORRECTLY CLASSIFIES

$$\max(SS_{\text{catalytic activity}}, SS_{\text{cellular process}}, SS_{\text{molecular adaptor activity}}, SS_{\text{molecular function regulator}}, SS_{\text{multicellular organismal process}}, SS_{\text{signaling}}, SS_{\text{behavior}} + SS_{\text{immune system process}})$$

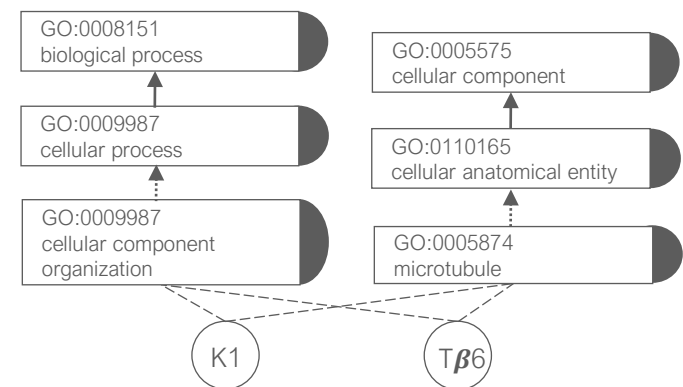
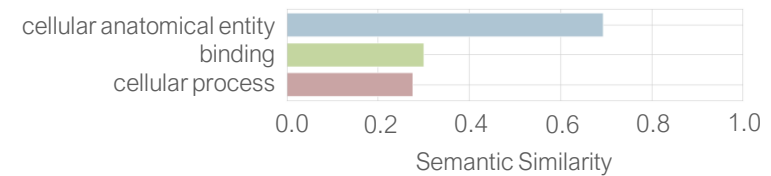
## True Positive (+/+)

40S ribosomal protein S12 and  
40S ribosomal protein S10



## True Negative (-/-)

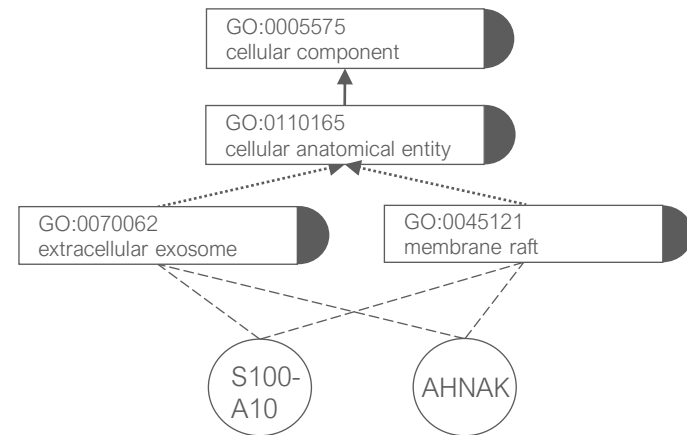
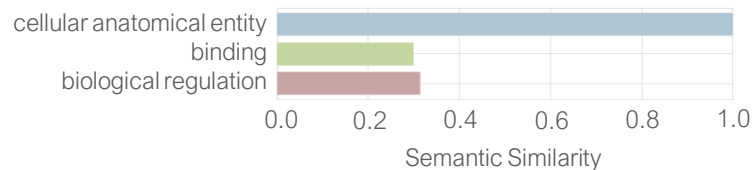
Kinetochores-associated protein 1 and  
Tubulin  $\beta$ -6 chain



# WHEN GP FAILS

## False Negative (+/-)

S100-A10 protein and neuroblast differentiation-associated protein AHNAK



Protein S100-A10 works together with neuroblast differentiation-associated protein AHNAK in the development of the intracellular membrane.



*J Cell Biol.* 2004 Jan 5; 164(1): 133–144.  
doi: [10.1083/jcb.200307098](https://doi.org/10.1083/jcb.200307098)

PMCID: PMC2171952  
PMID: [14699089](https://pubmed.ncbi.nlm.nih.gov/14699089/)

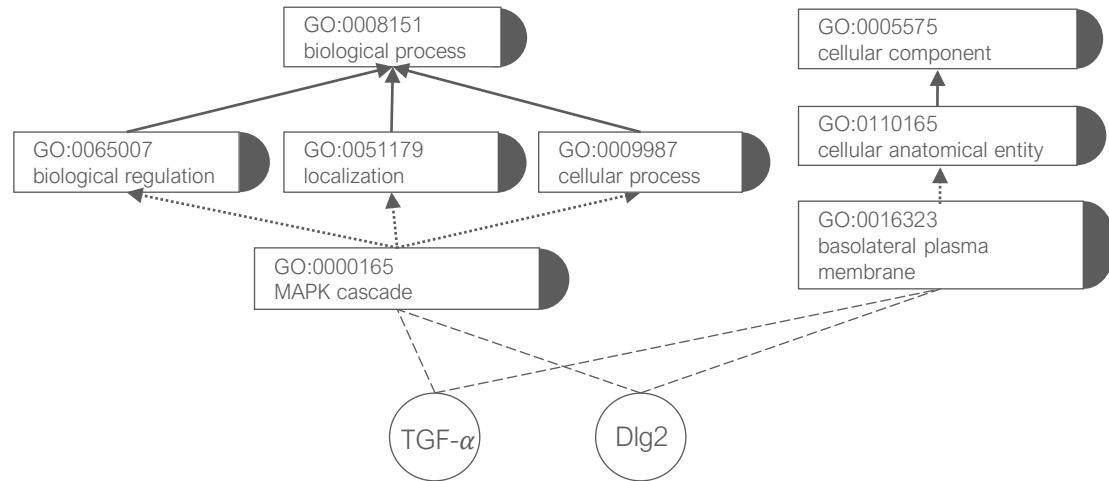
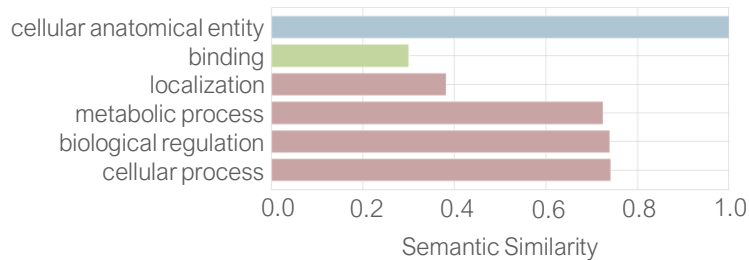
AHNAK interaction with the annexin 2/S100A10 complex regulates cell membrane cytoarchitecture

Christelle Benaud,<sup>1</sup> Benoît J. Gentil,<sup>1</sup> Nicole Assard,<sup>1</sup> Magalie Court,<sup>2</sup> Jerome Garin,<sup>2</sup> Christian Delphin,<sup>1</sup> and Jacques Baudier<sup>1</sup>

# WHEN MISCLASSIFICATIONS ARE NOT MISTAKES

## False Positive (-/+)

Protransforming growth factor  $\alpha$  and Disks large homolog 2



The literature describes interactions between proteins of the same family of the pair, indicating that this is likely a true but still unknown interaction.



FEBS Letters

Volume 587, Issue 11, 5 June 2013, Pages 1624-1629



Dlg5 interacts with the TGF- $\beta$  receptor and promotes its degradation

Edited by Gianni Cesareni

Takuhito Sezaki <sup>a</sup>, Lucia Tomiyama <sup>a, b</sup>, Yasuhisa Kimura <sup>a</sup>, Kazumitsu Ueda <sup>a, b</sup>, Noriyuki Kioka <sup>a, b, c</sup>

# CLOSING REMARKS



- Explainability can be key to **uncover issues with the underlying data** and even **pose new hypothesis**.
- When we use GP, the explanation is the model itself, avoiding the need for local explanations or post-hoc techniques. Moreover, operators can be tailored to the domain.
- The performance of the more interpretable GP methods is not substantially lower, but what little they sacrifice in performance is more than gained in explainability.

# ACKNOWLEDGEMENTS




Catia Pesquita



Sara Silva

 risousa@ciencias.ulisboa.pt

 @RitaTorresSousa

Read the short paper: [https://github.com/liseda-lab/ExplainingPPIpredictions/blob/main/Explaining\\_PPI\\_predictions\\_with\\_GP.pdf](https://github.com/liseda-lab/ExplainingPPIpredictions/blob/main/Explaining_PPI_predictions_with_GP.pdf)

This work was funded by FCT through LASIGE Research Unit (UIDB/00408/2020, UIDP/00408/2020); projects GADgET (DSAIPA/DS/0022/2018) and BINDER (PTDC/CCI-INF/29168/2017); PhD grant SFRH/BD/145377/2019.