



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Interpretability in Artificial Intelligence applications for rare diseases

Davide Cirillo

Machine Learning for Biomedical Research

Life Sciences Department

04/05/2022

Barcelona Supercomputing Center



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

The MareNostrum 4 supercomputer

Total peak performance:

13,7 Pflops/s

The MareNostrum 5 supercomputer

2022-2027

>200 Pflops/s

Disk storage
+150 PB

Tape storage
+400 PB



Access: prace-ri.eu/hpc-access



RED ESPAÑOLA DE
SUPERCOMPUTACIÓN

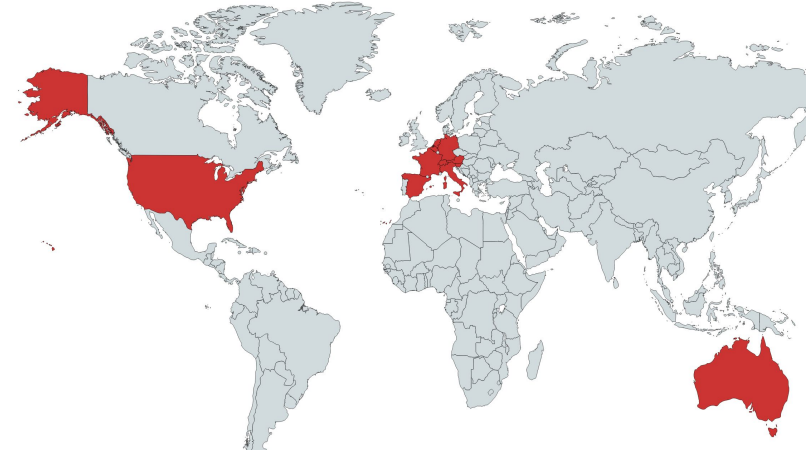
Access: bsc.es/res-intranet



Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación



21 partners from 3 continents



TECHNIKON

IBM Research



- **Cloud-based platform** to share harmonized data on paediatric cancers.
- **Predictive models** verified through **clinical trials and preclinical models**.
- **Personalized treatment recommendations** for paediatric cancer patients.



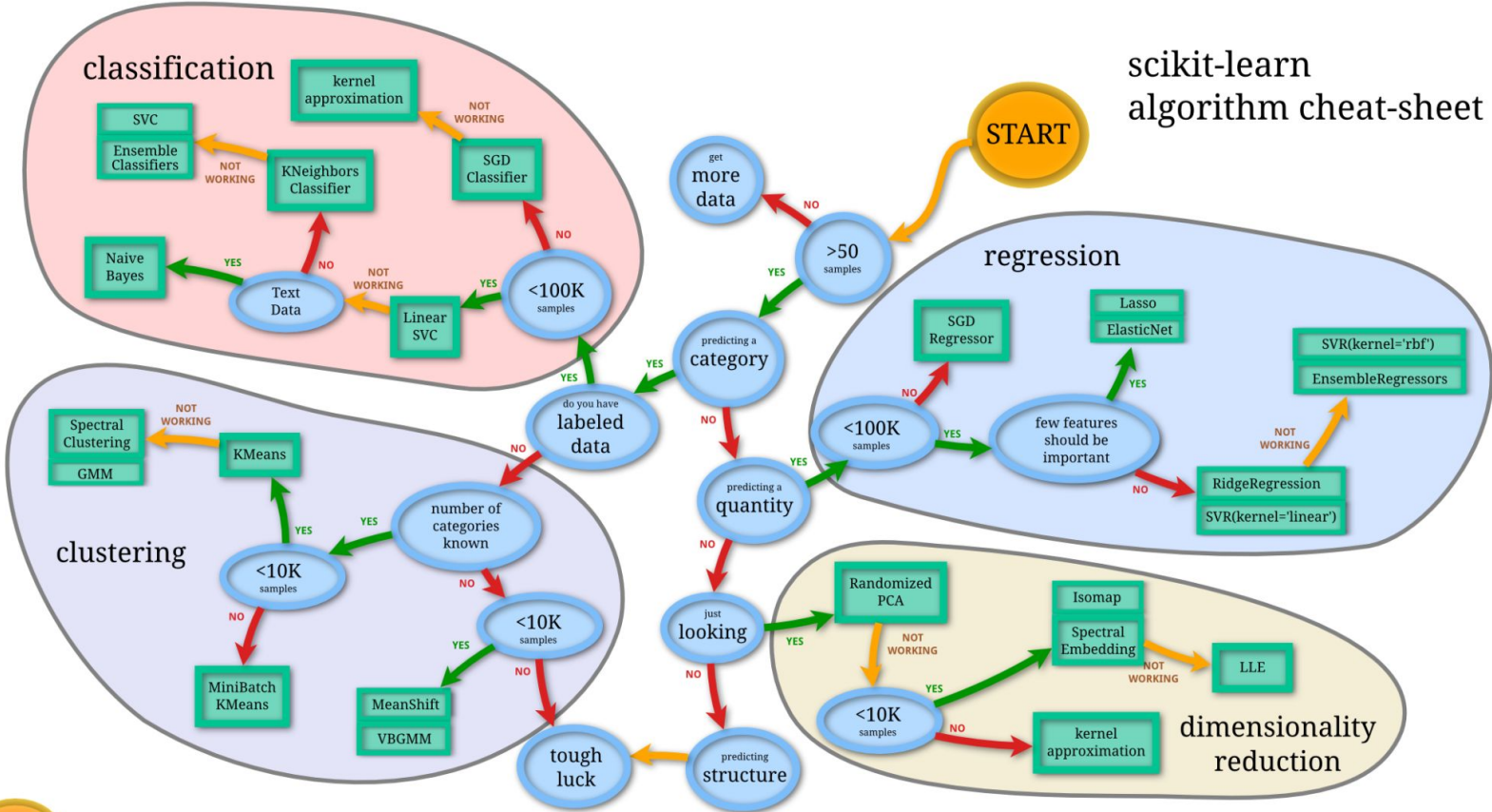
Barcelona Supercomputing Center

Centro Nacional de Supercomputación

Rare diseases

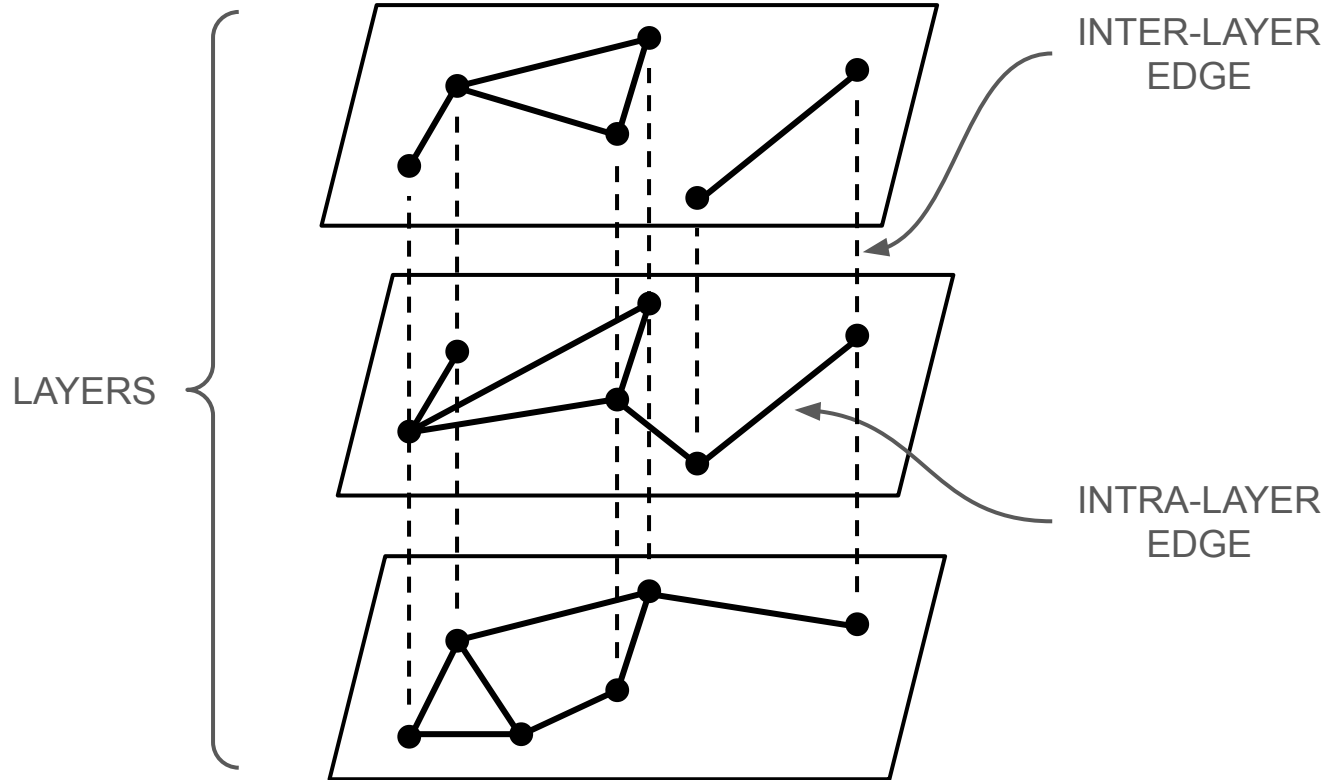


scikit-learn algorithm cheat-sheet



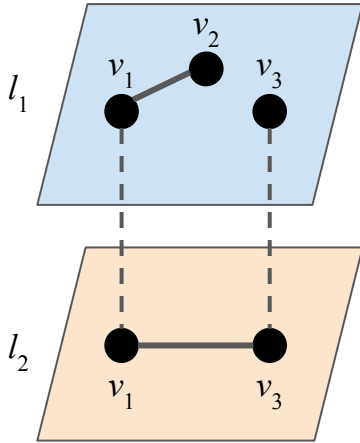


Multilayer networks



Multilayer network formalization

$$M = (V_M, E_M, V, L)$$



	v_1	v_2	v_3
l_1	1	1	1
l_2	1	0	1

	(l_1, v_1)	(l_1, v_2)	(l_1, v_3)	(l_2, v_1)	(l_2, v_2)	(l_2, v_3)
(l_1, v_1)	0	0	0	0	0	0
(l_1, v_2)	1	0	0	0	0	0
(l_1, v_3)	0	0	0	0	0	0
(l_2, v_1)	1	0	0	0	0	0
(l_2, v_2)	0	0	0	0	0	0
(l_2, v_3)	0	0	1	1	0	0

$$V = \{v_1, v_2, v_3\}$$

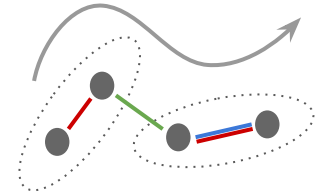
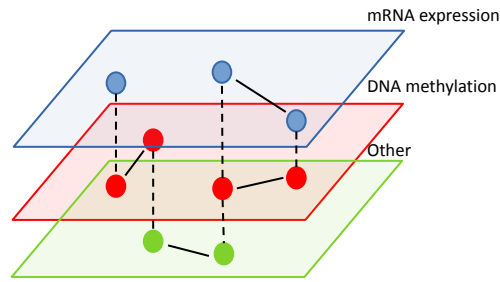
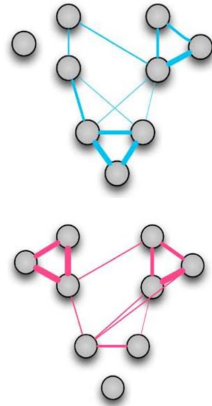
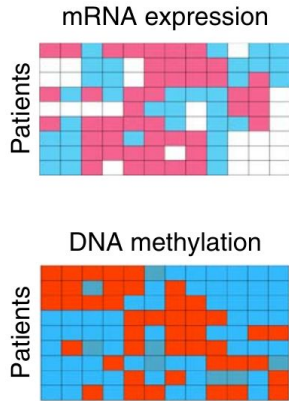
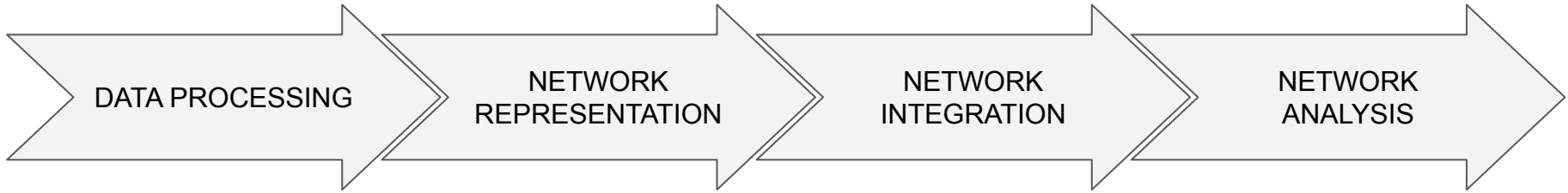
$$L = \{l_1, l_2\}$$

$$V_M \subseteq V \times L$$

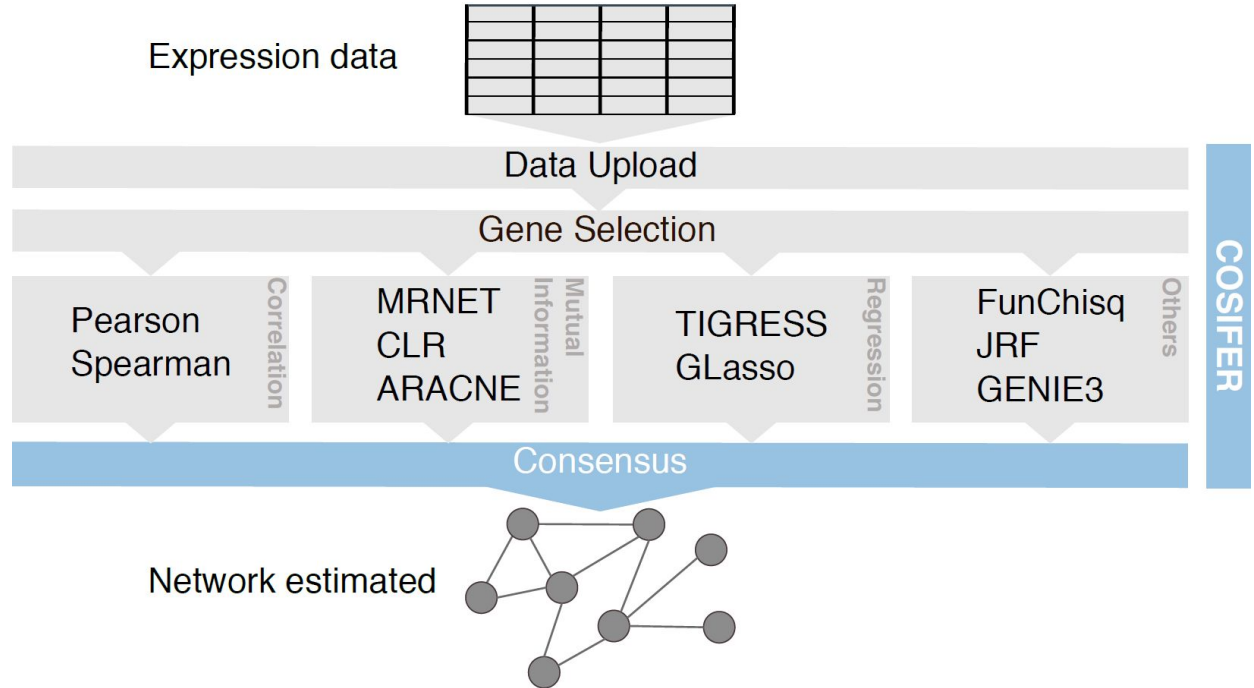
$$E_M \subseteq V_M \times V_M$$



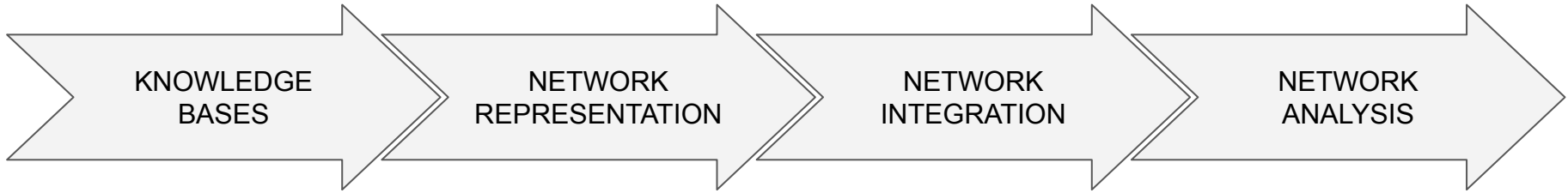
Multilayer network analysis of biomedical information



Network inference with COSIFER



Multilayer network analysis of biomedical information



GWAS Catalog
The NHGRI-EBI Catalog of published genome-wide association studies

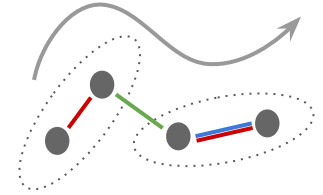
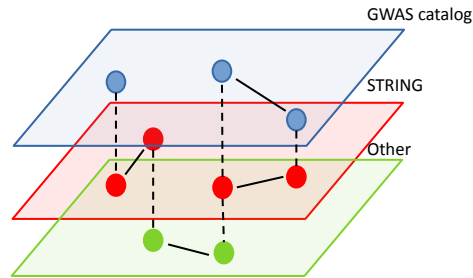
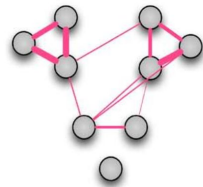
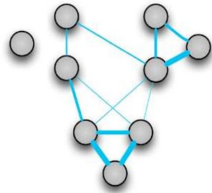
Search the catalog

Examples: breast carcinoma, rs7329174, Yao, 2q37.1, HBS1L, 6:16000000-25000000

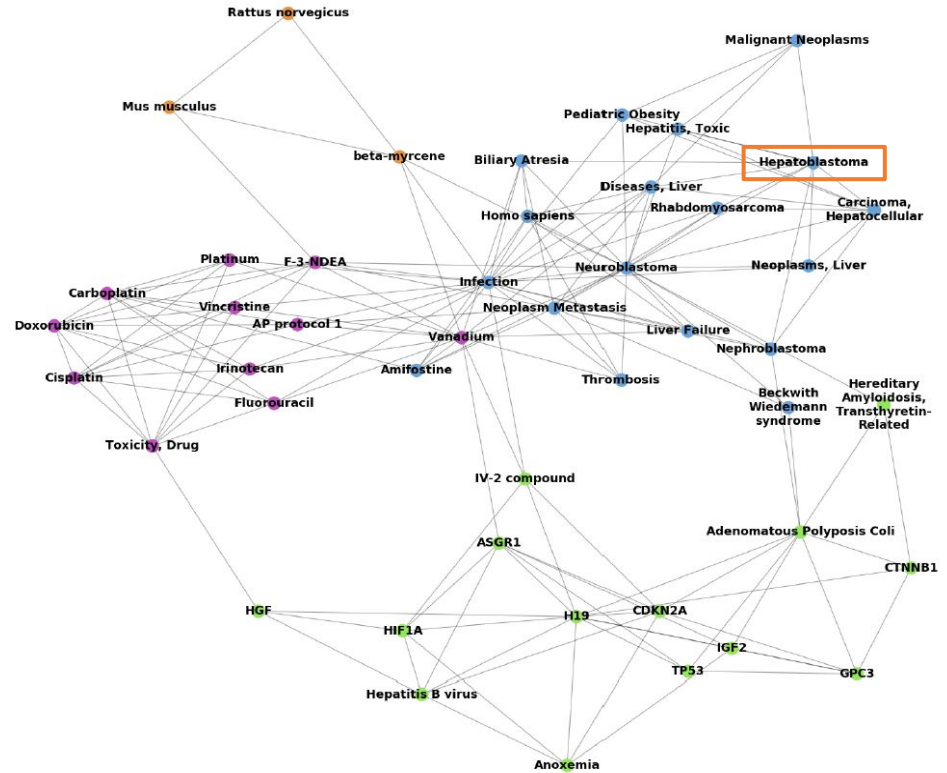
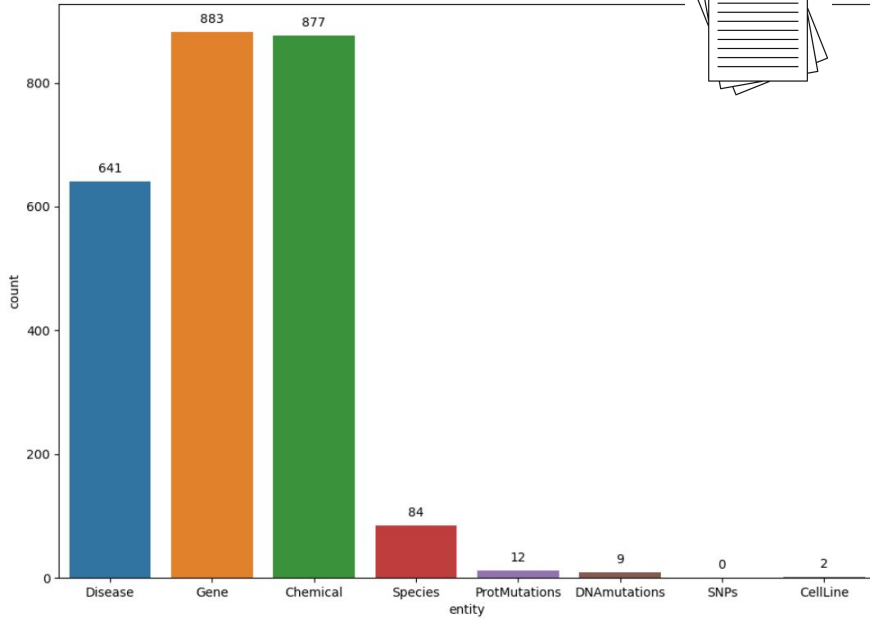
Version: 11.0



STRING



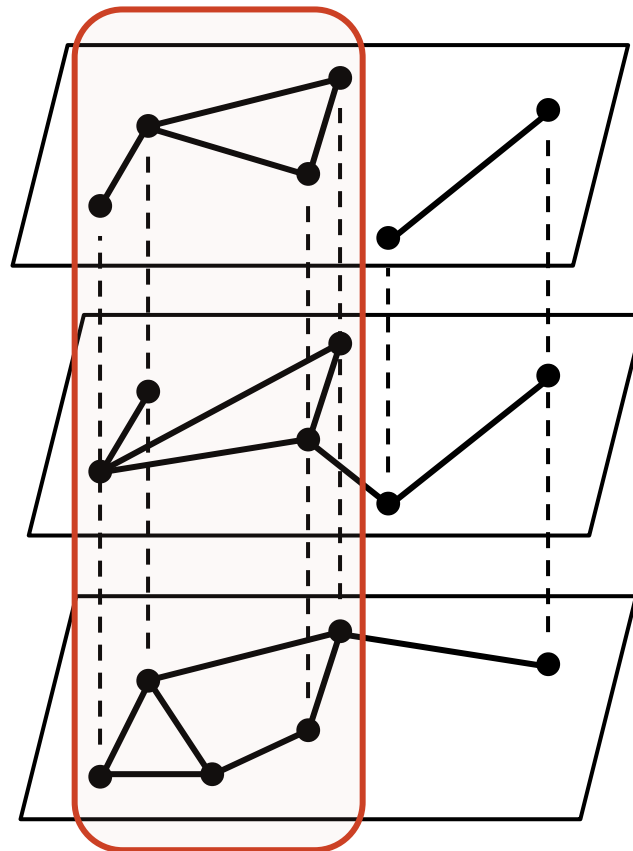
Hepatoblastoma



Multilayer community trajectories to study rare diseases



Multilayer community detection

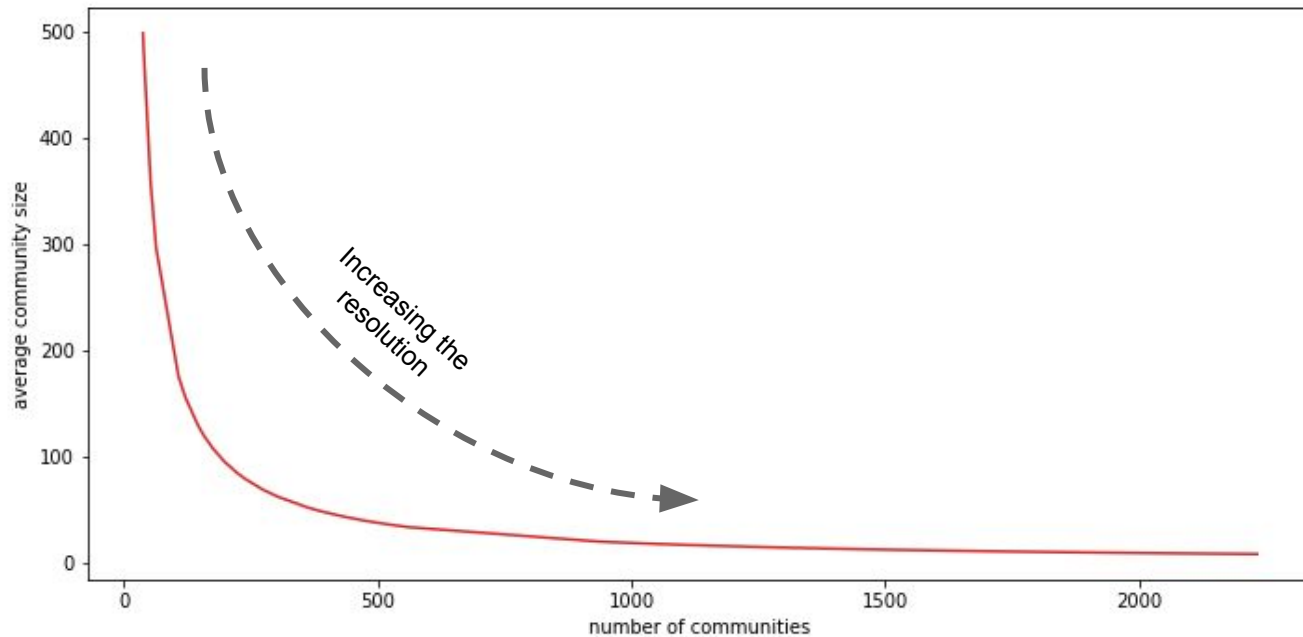
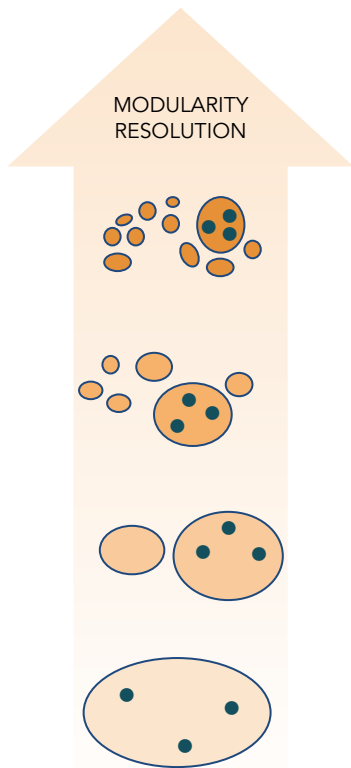


maximize $\sum_g \mathcal{Q}_\gamma(X^{(g)}, \mathbf{c})$

Newman & Girvan. *Phys Rev E*. 2004
Blondel et al. *J Stat Mech*. 2008
Didier et al. *PeerJ*. 2015
Didier et al. *F1000Res*. 2018

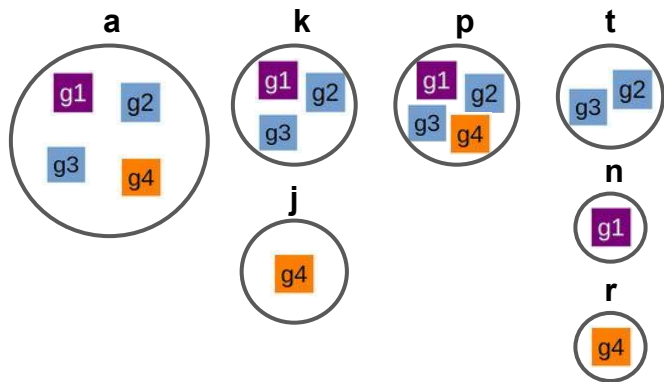
Multilayer community

Modularity resolution



Multilayer community trajectories

Modularity resolution →



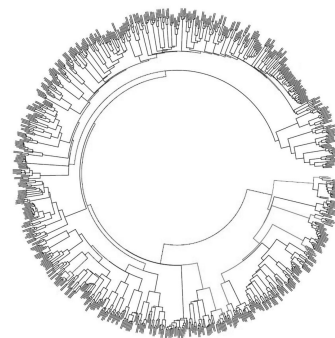
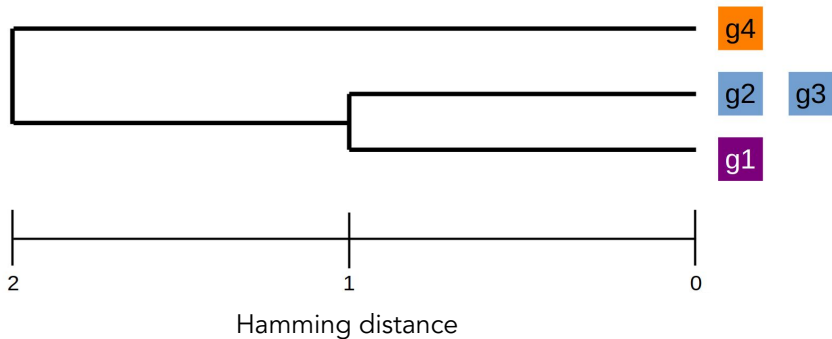
g1	a	k	p	n
g2	a	k	p	t

g2	a	k	p	t
g3	a	k	p	t

g1	a	k	p	n
g4	a	j	p	r

g3	a	k	p	t
g4	a	j	p	r

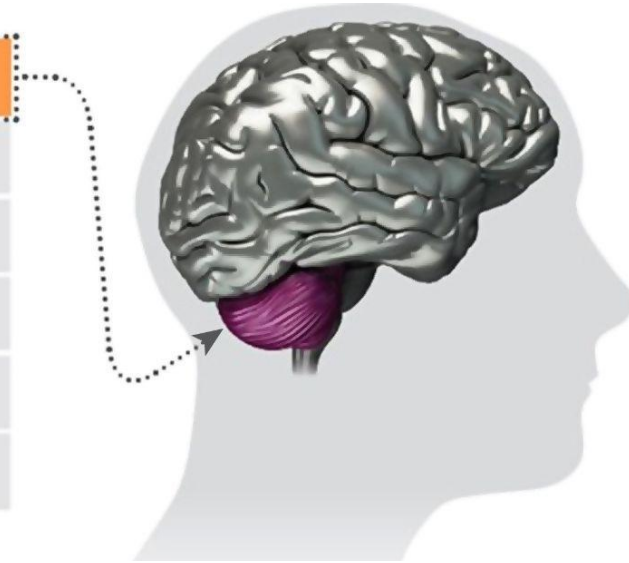
	g1	g2	g3	g4
g1	0	1	1	2
g2	1	0	0	2
g3	1	0	0	2
g4	2	2	2	0



- distance in the tree
- community composition

Medulloblastoma

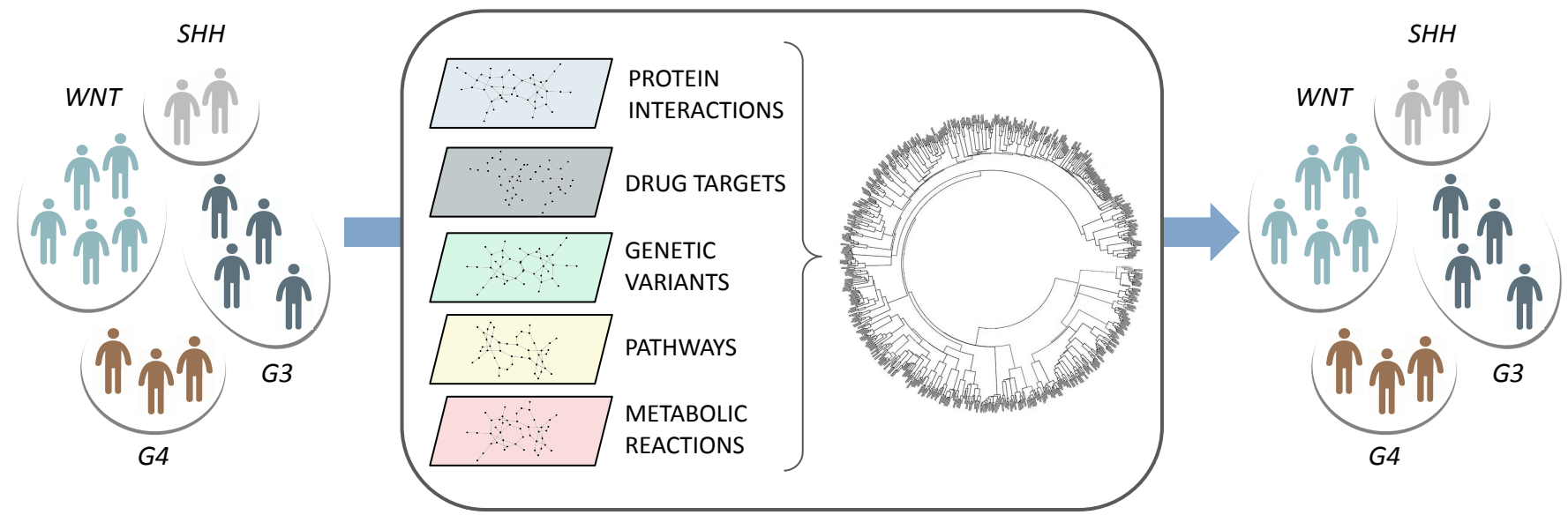
Molecular subtype	Wnt group	Shh group	Group 3	Group 4
Prognosis	Very good	Good in infant and intermediate in others	Poor	Intermediate
Main signaling pathway	Wnt	Shh, PI3K	TGF- β , photoreceptor/ GABAergic	NF- κ B
Metastasis	Rare	Uncommon	Very frequent	Frequent
Characteristic feature	<i>CTNNB1</i> mutation	<i>SMO/PTCH/SUFU</i> mutation	<i>MYC</i> amplification	<i>CDK6</i> amplification
<i>MYC</i> status	<i>MYC</i> ⁺	<i>MYCN</i> ⁺	<i>MYC</i> ⁺⁺⁺	Minimal <i>MYC/MYCN</i>



Patient stratification in medulloblastoma

PROTEOGENOMIC DATA
(~14,000 altered genes per patient)

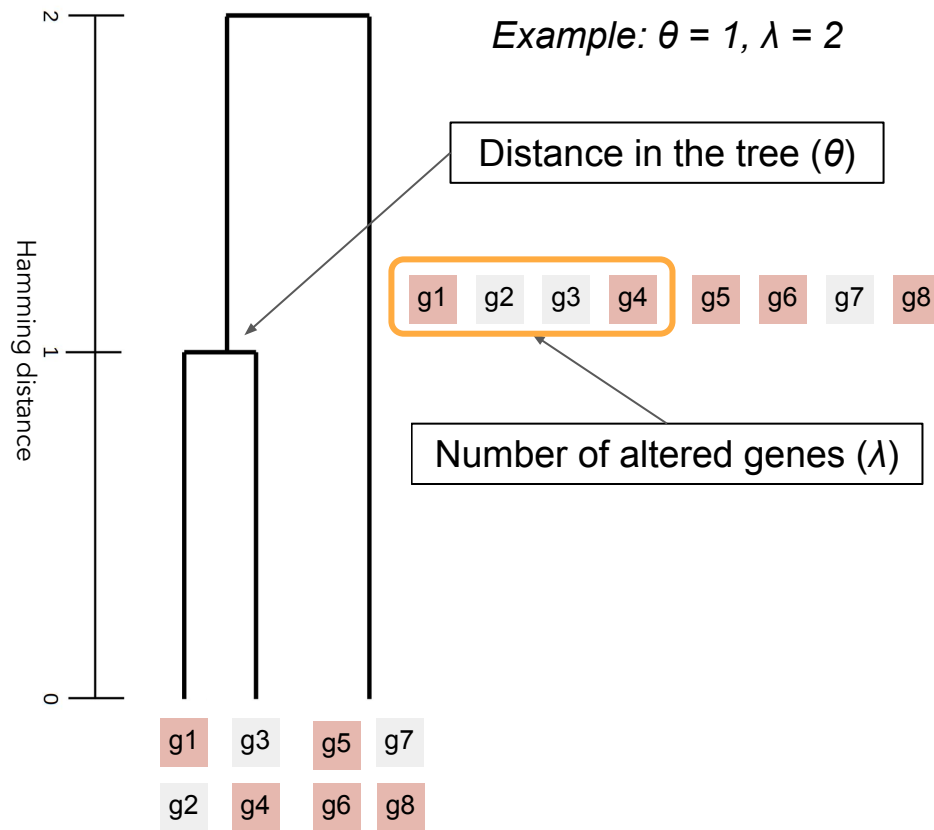
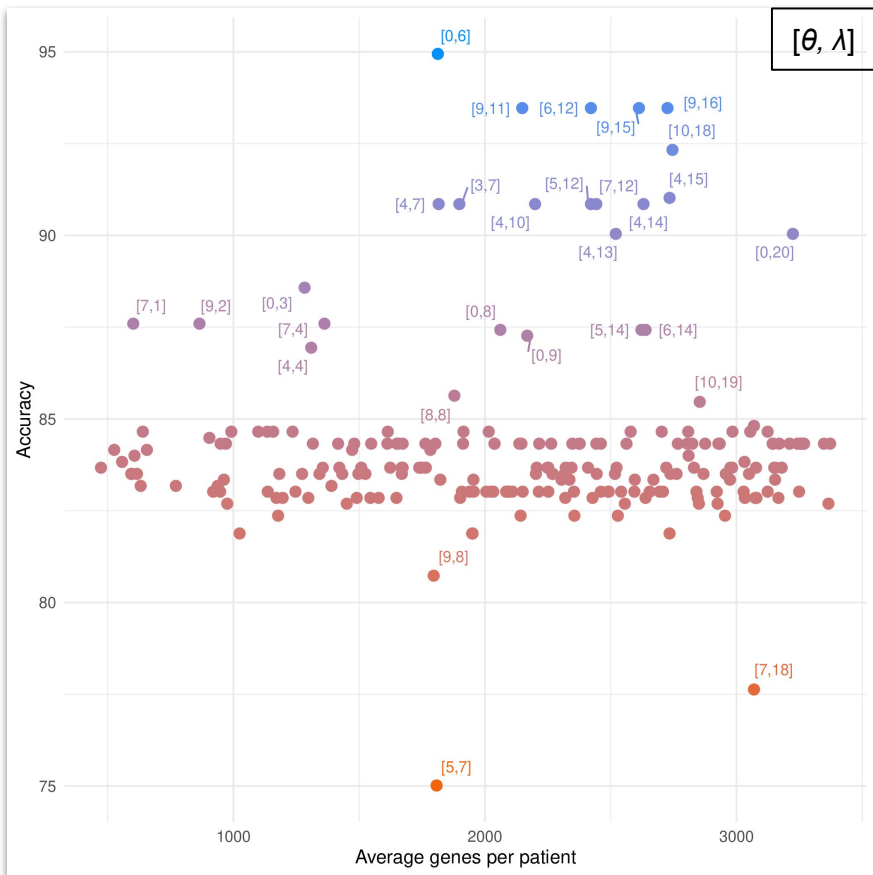
MINIMAL SET OF GENES
(~1,800 altered genes per patient)



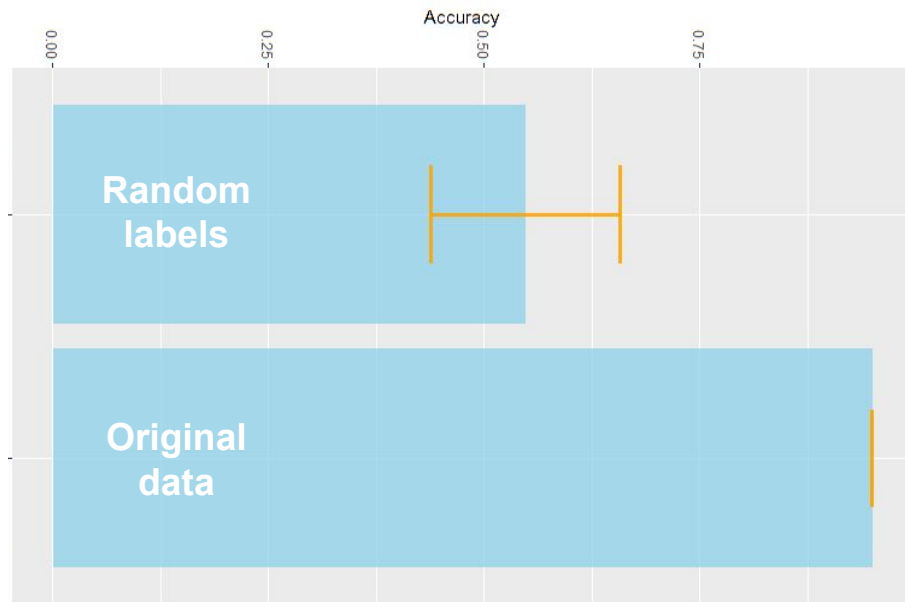
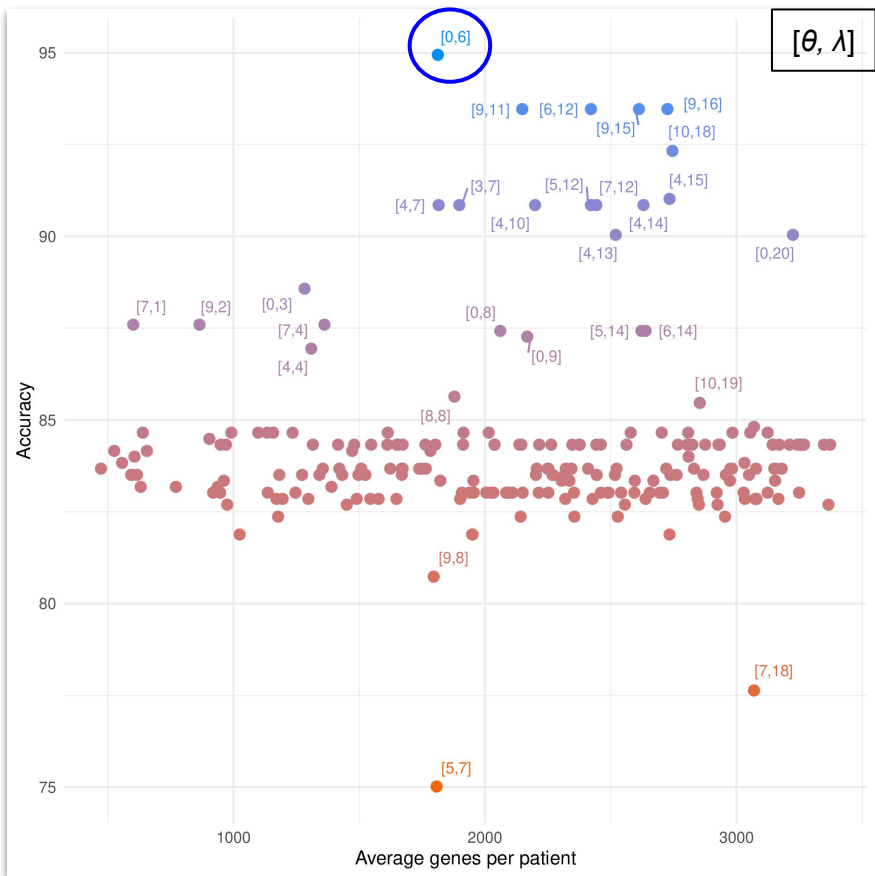
95% accuracy

87% dimensionality reduction

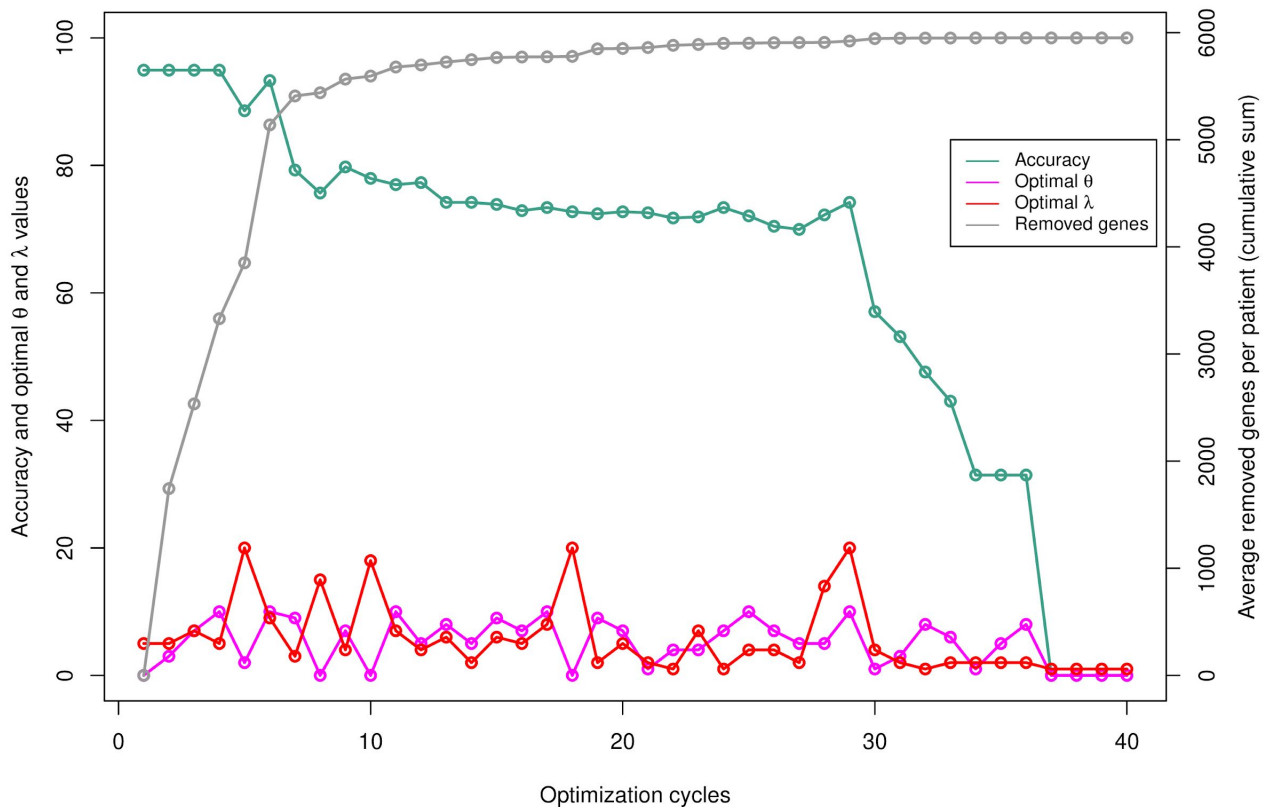
An optimization problem



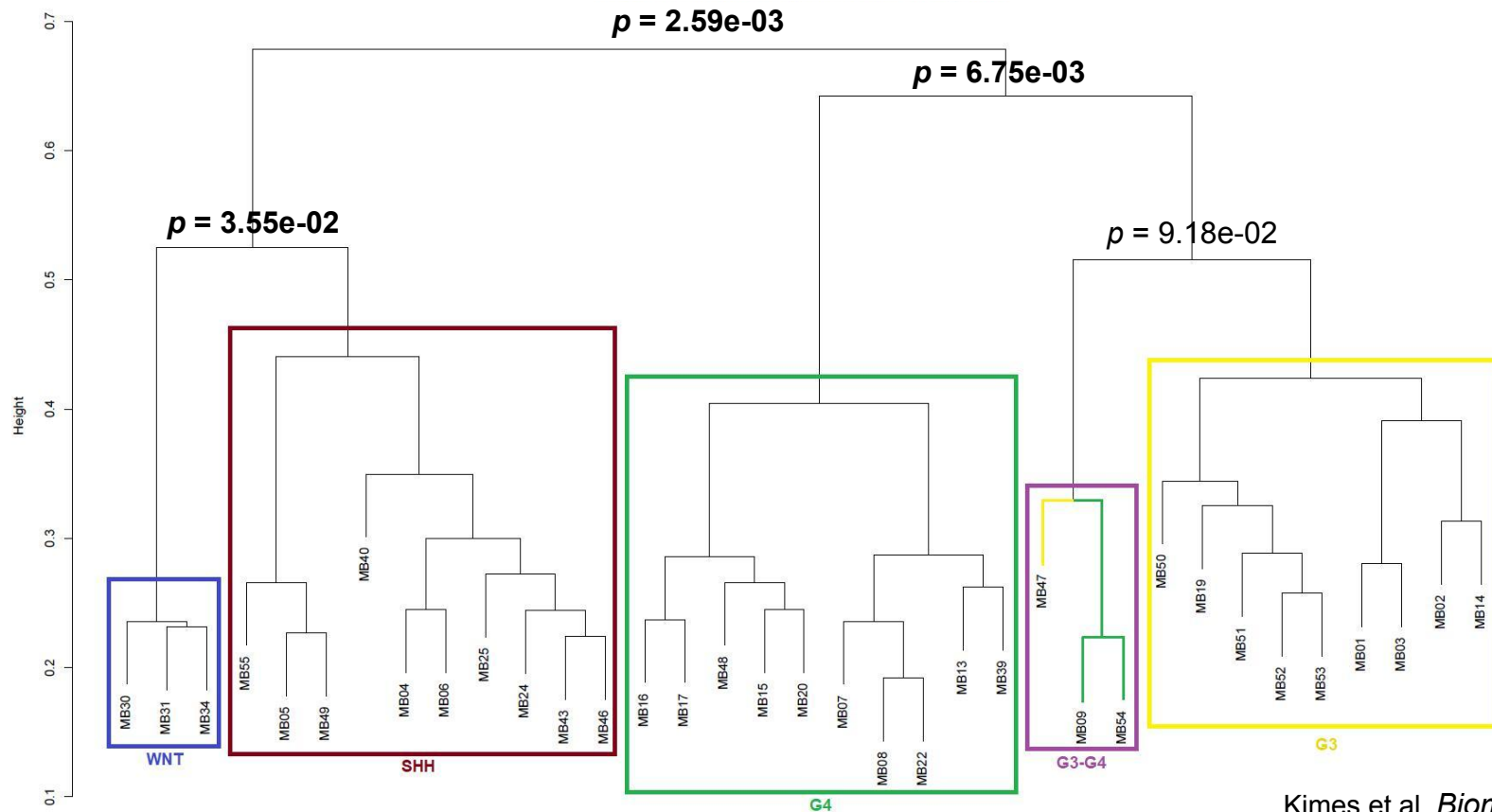
An optimization problem



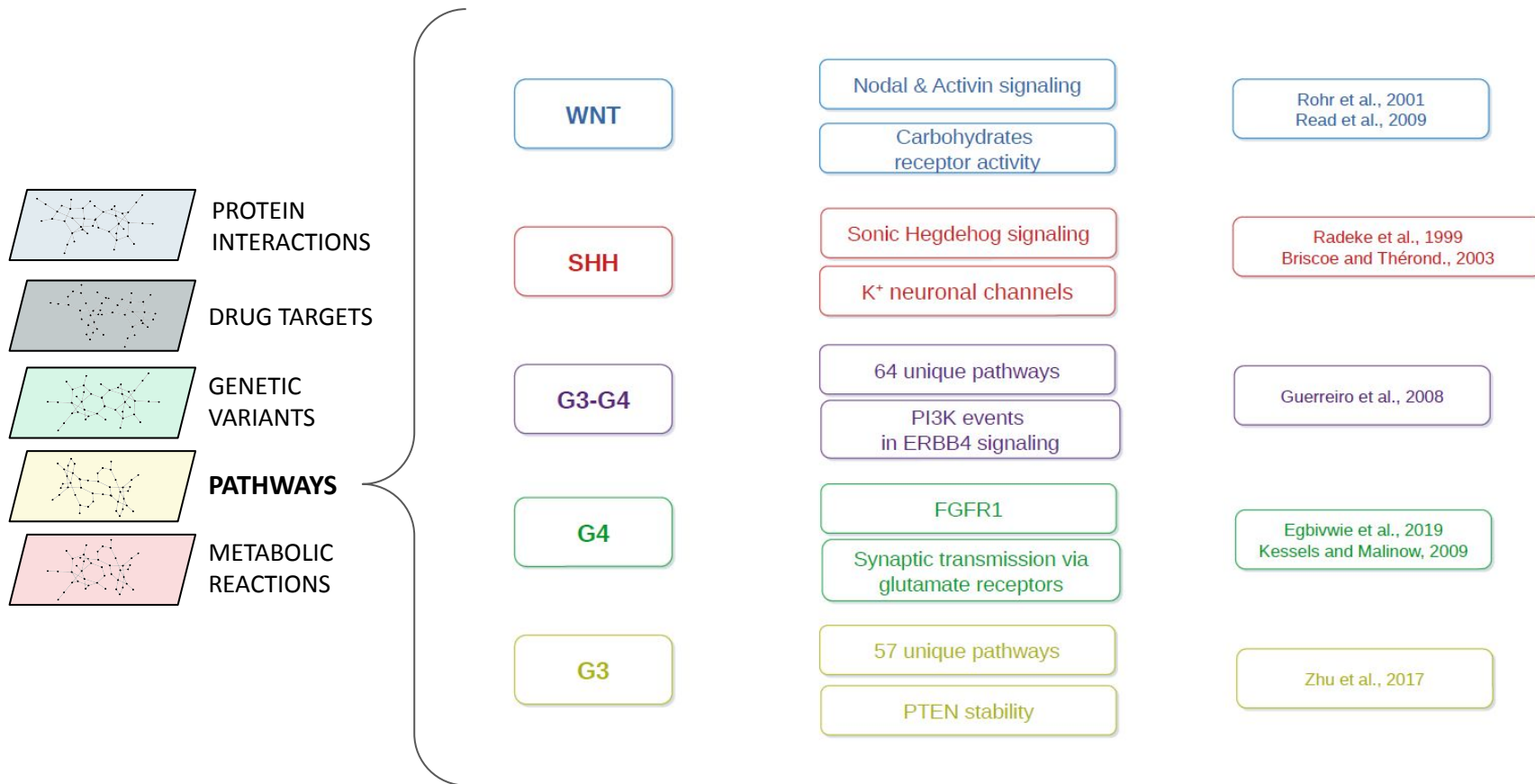
Sequential exclusion test



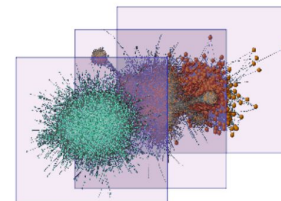
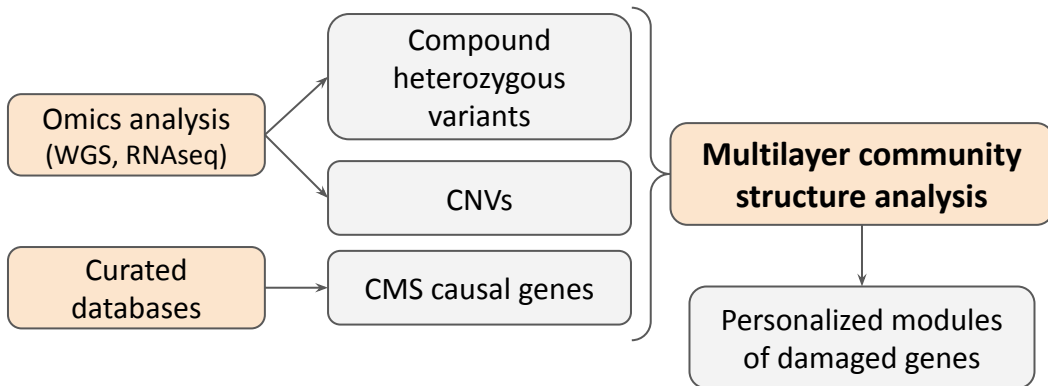
Patient stratification



Provenance of the associations



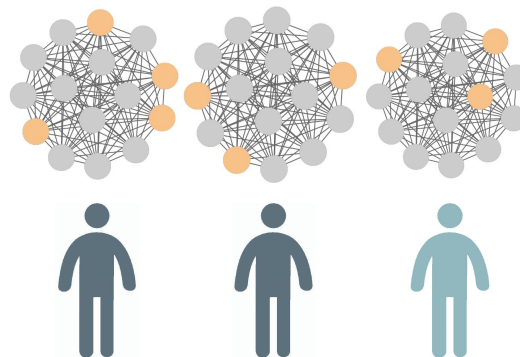
Severity in Congenital Myasthenic Syndromes



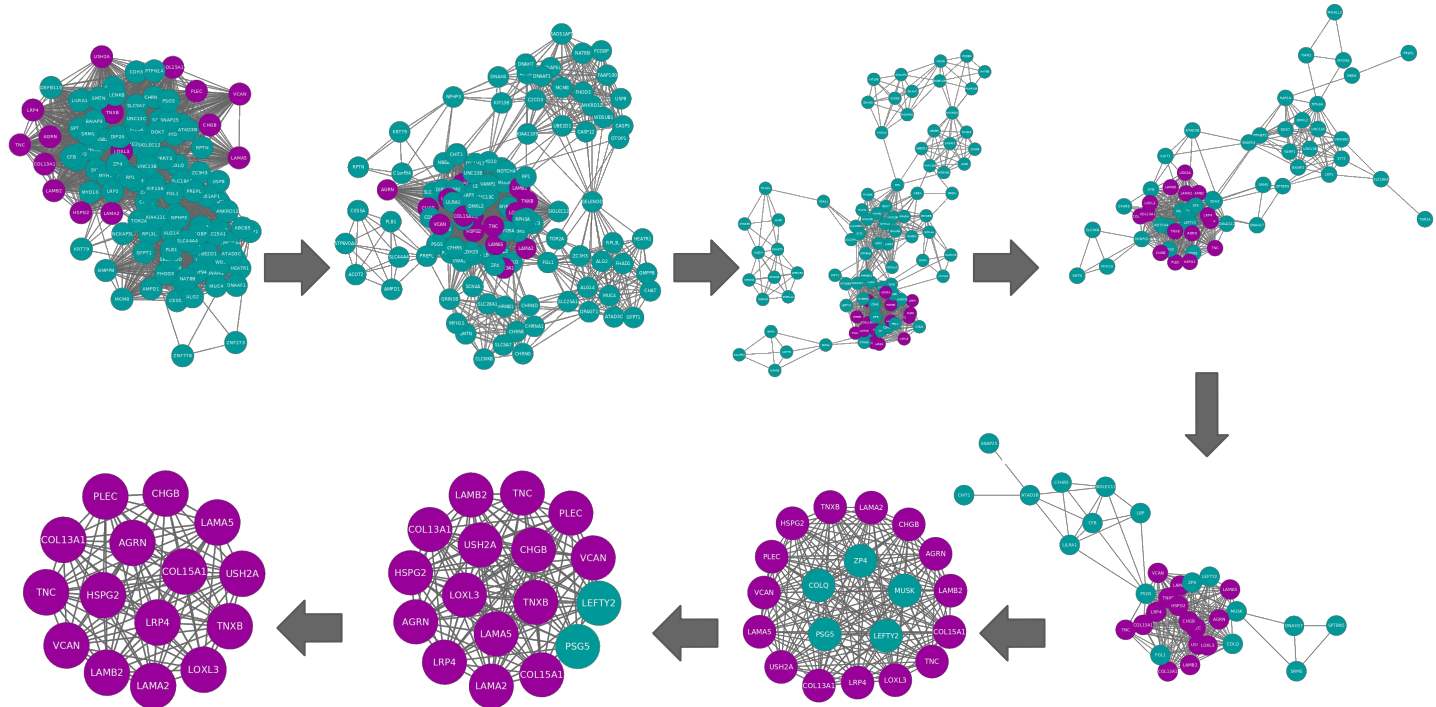
*Interactome
Metabolome
Reactome*

CHRNE c.1327delG

Severe Non-severe

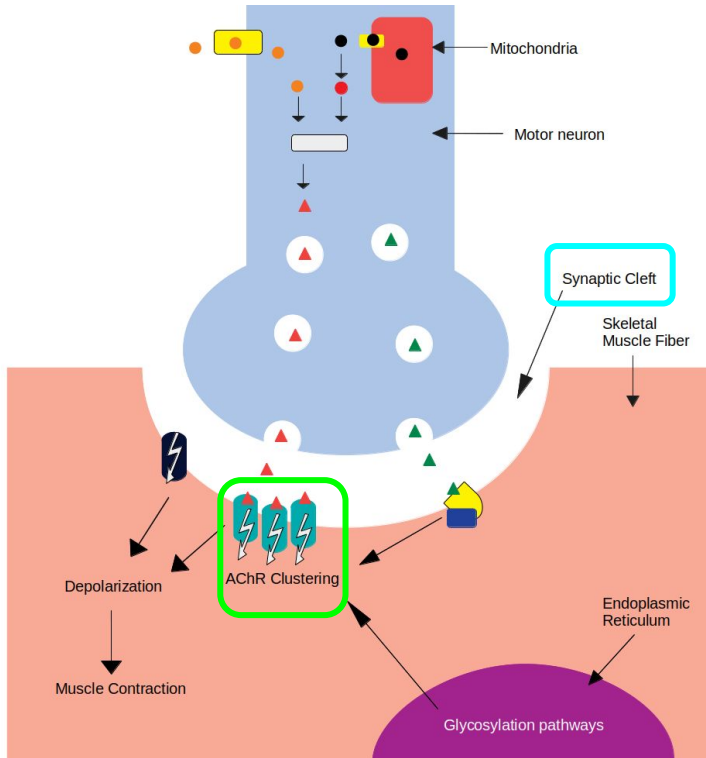


Severity in Congenital Myasthenic Syndromes



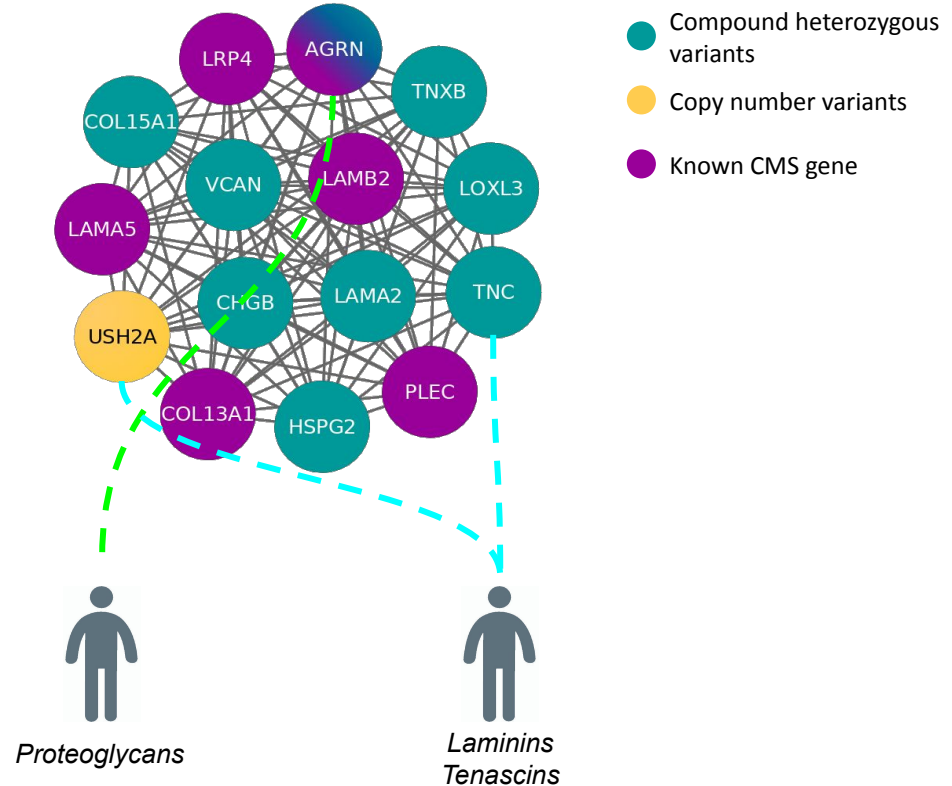
Severity in Congenital Myasthenic Syndromes

Neuromuscular junction

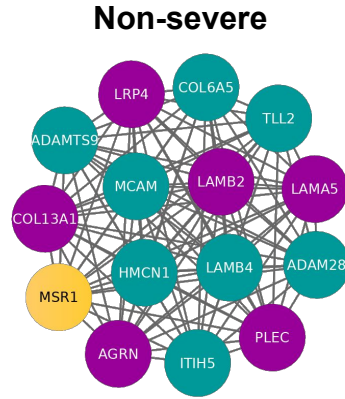
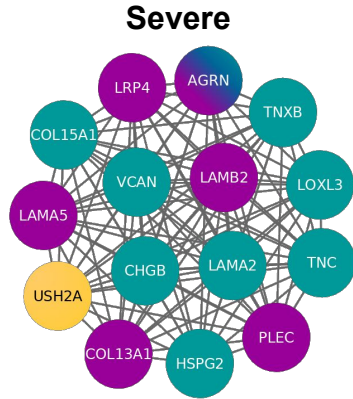


Nuñez et al. (in preparation)

Gene module of severe cases



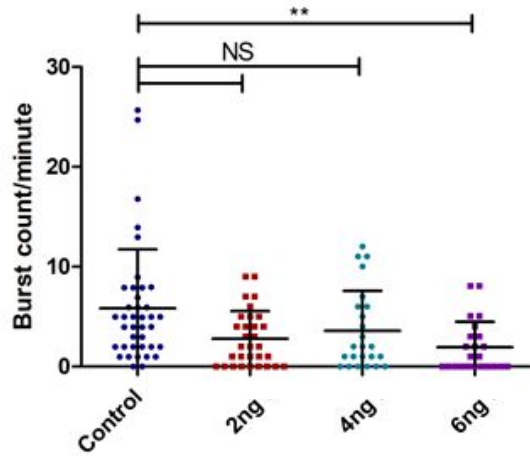
Severity in Congenital Myasthenic Syndromes



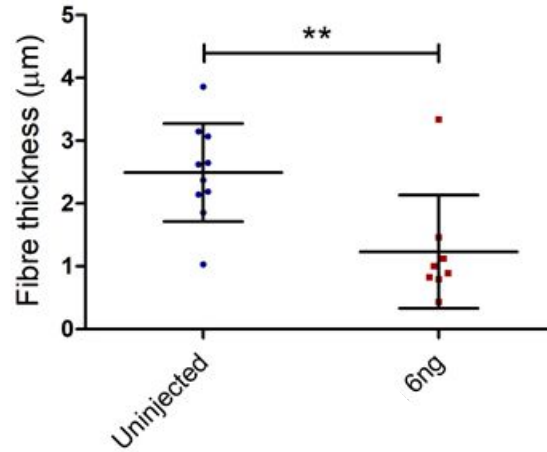
Activity localization	Class	CMS causal gene	Phenotype group		Function	Synaptic localization (Manual curation)	Localization (UniProt)
			Not-severe	Severe			
ECM (ECM)	Proteoglycans	AGRN	-	AGRN	Cell hydration and growth factor trapping	Pre- and post-synaptic (PMID:29462312)	Synaptic basal lamina / ECM
		-	-	HSPG2		Basement membrane (PMID:30453502)	Basement membrane / ECM
		-	-	VCAN		ECM (PMID:29211034)	ECM
	Collagens	COL13A1	-	-	Structural support	Basement membrane, post-synaptic (PMID:30768864)	Post-synaptic cell membrane
		-	COL6A5	-		Basement membrane (PMID:23869615)	Extracellular matrix
	Laminins	LAMA5	-	-	Web-like structures	Pre-synaptic (PMID:28544784)	Basement membrane / ECM
		LAMB2	-	-		Basement membrane (PMID:27614294)	Basement membrane / ECM / Synaptic cleft
		-	LAMB4	-		Myenteric plexus basement membrane (PMID:28595269)	Basement membrane / ECM
		-	-	LAMA2		Pre-synaptic (PMID:9396756)	Basement membrane / ECM
		-	-	USH2A		Neuronal projection of stereocilia (PMID:19023448)	Stereocilium membrane / Secreted (Extracellular region)
	Fibulins	-	HMCHL1	-	Scaffolding	Glomerular Extracellular matrix (PMID:29488390)	Basement membrane / ECM
	Tenascins	-	-	-	Anti-adhesion	Basement membrane (PMID:29466693)	ECM / Perisynaptic ECM (Ensembl)
		-	-	THC		Basement membrane (PMID:23768946)	ECM
	Enzymes	-	-	-	Collagen assembly	Basement membrane (PMID:26954549)	Secreted (extracellular region)
		-	ADAMTS9	-		Secreted to ECM (PMID:30626608)	ECM
		-	ADAM28	-		ECM (PMID:24613731)	Cell membrane / Secreted (extracellular region)
	Neuropeptides	-	-	CHGB	Regulatory peptides precursor	Pre- and post-synaptic (PMID:7526287)	Secreted (extracellular region)
Others	-	ITIH5	-	Hyaluronic acid binding	ECM (PMID:27143355)	Secreted (extracellular region)	
Cell surface	Receptors	-	MSR1	Proteoglycan and collagen binding	Macrophage surface Scavenger Receptor (PMID:12488451)	Plasma membrane	
		-	MCAM		Plasma membrane (PMID:28923978)	Plasma membrane	
		LRP4	-	Laminin binding	Post-synaptic (PMID:26319686)	Post-synaptic cell membrane	
Cytoplasm	Cytoskeleton	PLEC	-	Structural support	Post-synaptic (PMID:20624679)	Post-synaptic cytoskeleton	



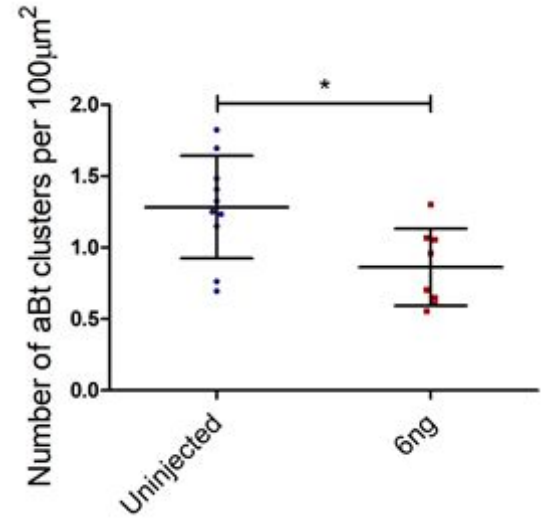
Decreased chorion movement
(1 day post fertilisation)



Decreased muscle fibre thickness
(5 days fish)

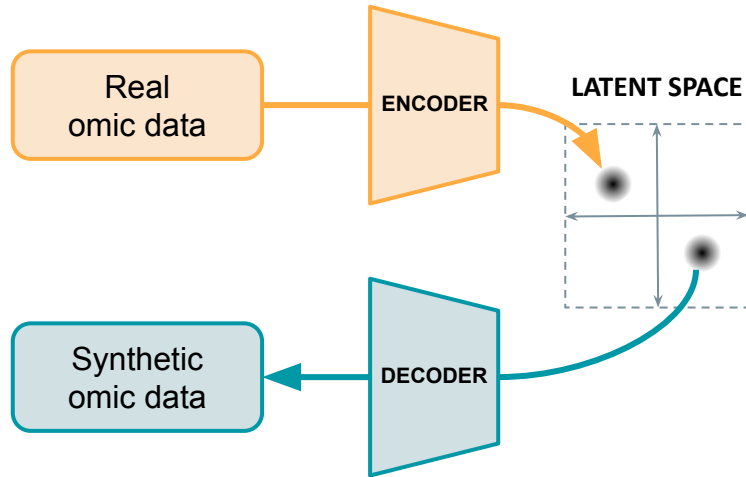


Decreased number of AChR clusters
(5 days fish)

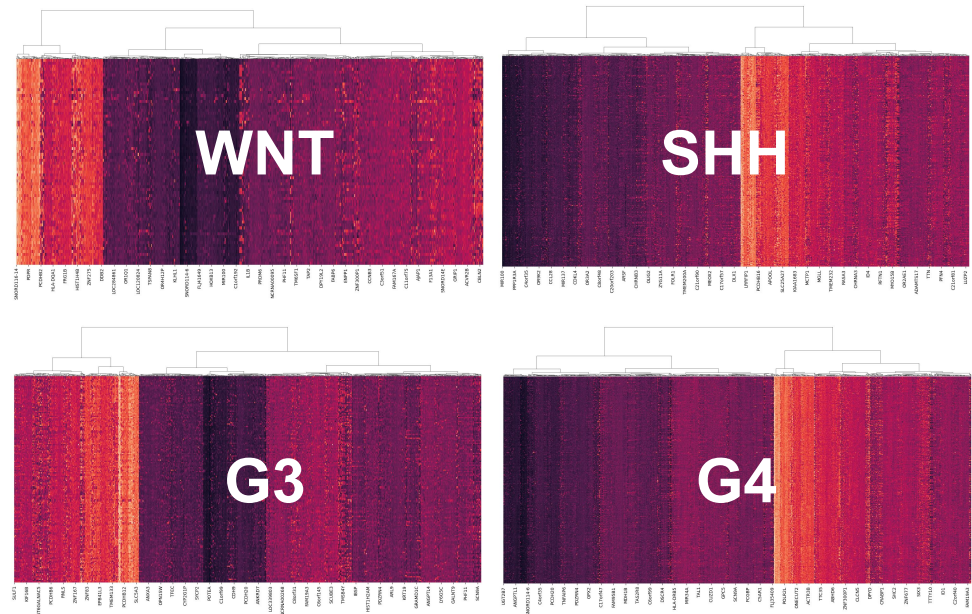


Explainable synthetic data generation for paediatric cancer

A **Variational Autoencoder** (VAE) can learn representations of **real data** of patients and therefore generate **synthetic data**.

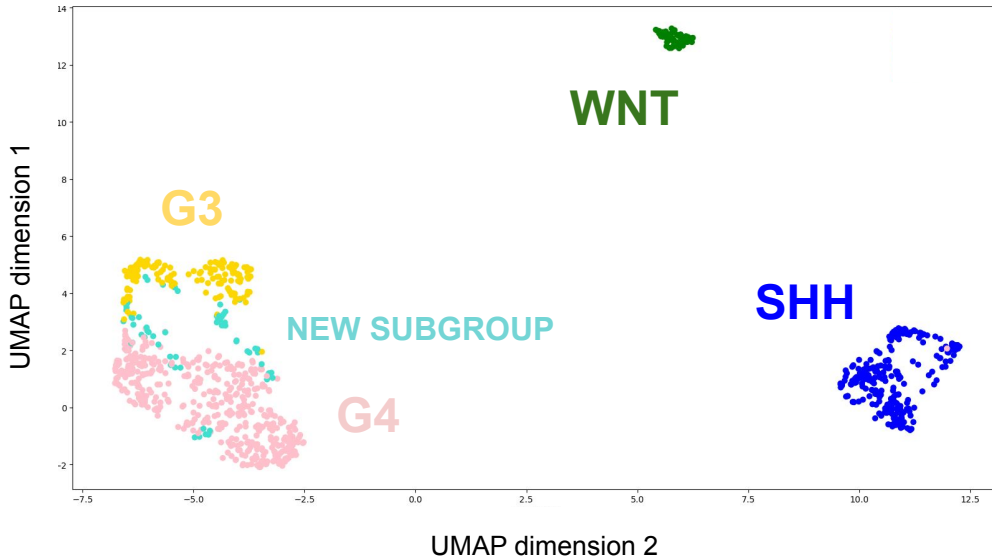


By studying **how** the VAE generates synthetic data we identified four distinct **omic signatures** of a **childhood brain tumor** (medulloblastoma).

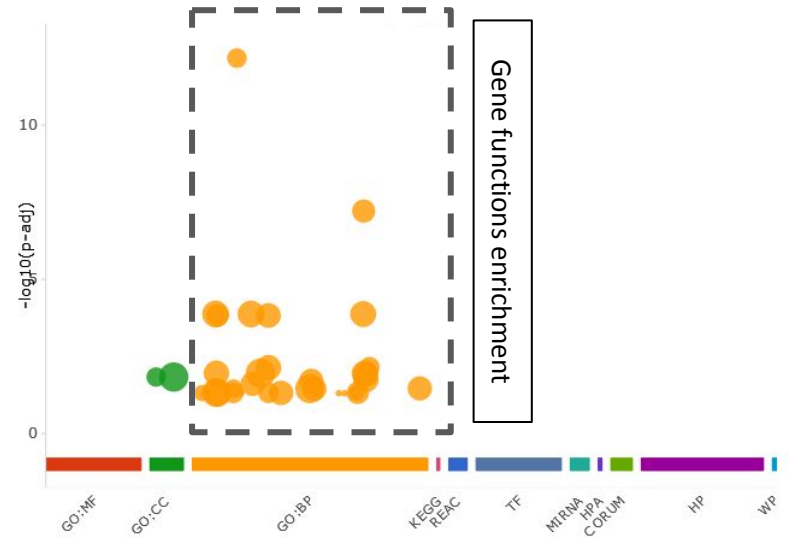


Explainable synthetic data generation for paediatric cancer

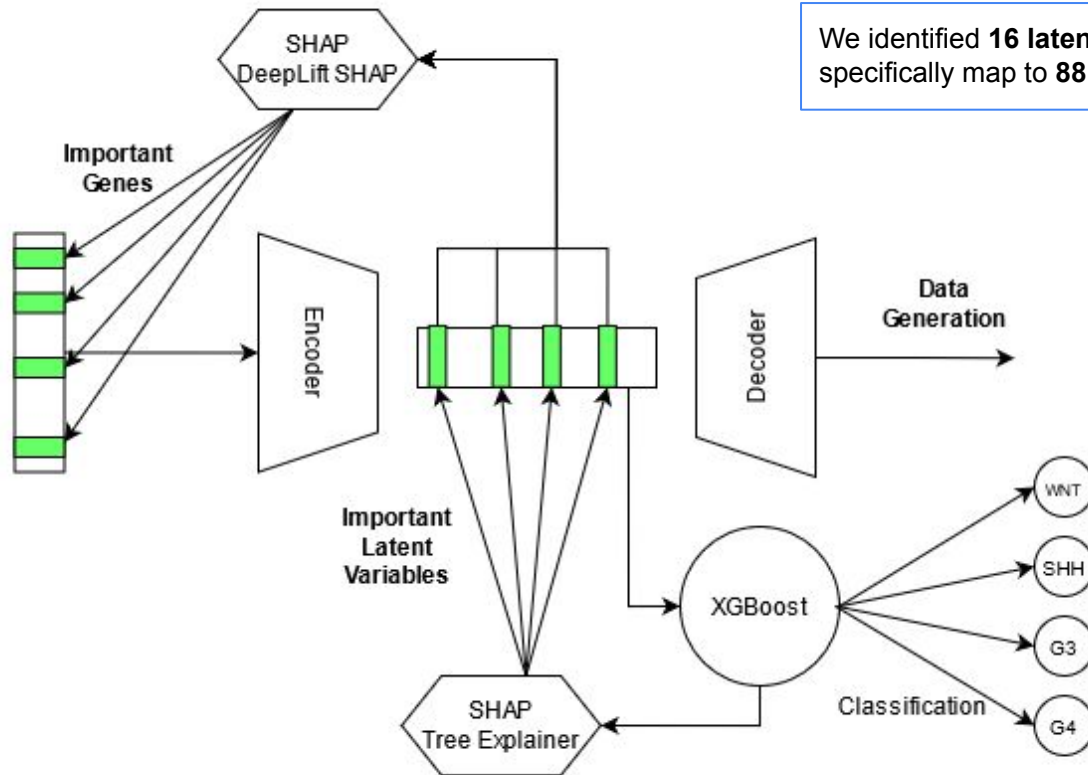
The VAE allowed us to discover **a new subgroup** of the childhood brain tumor characterized by **specific genes**.



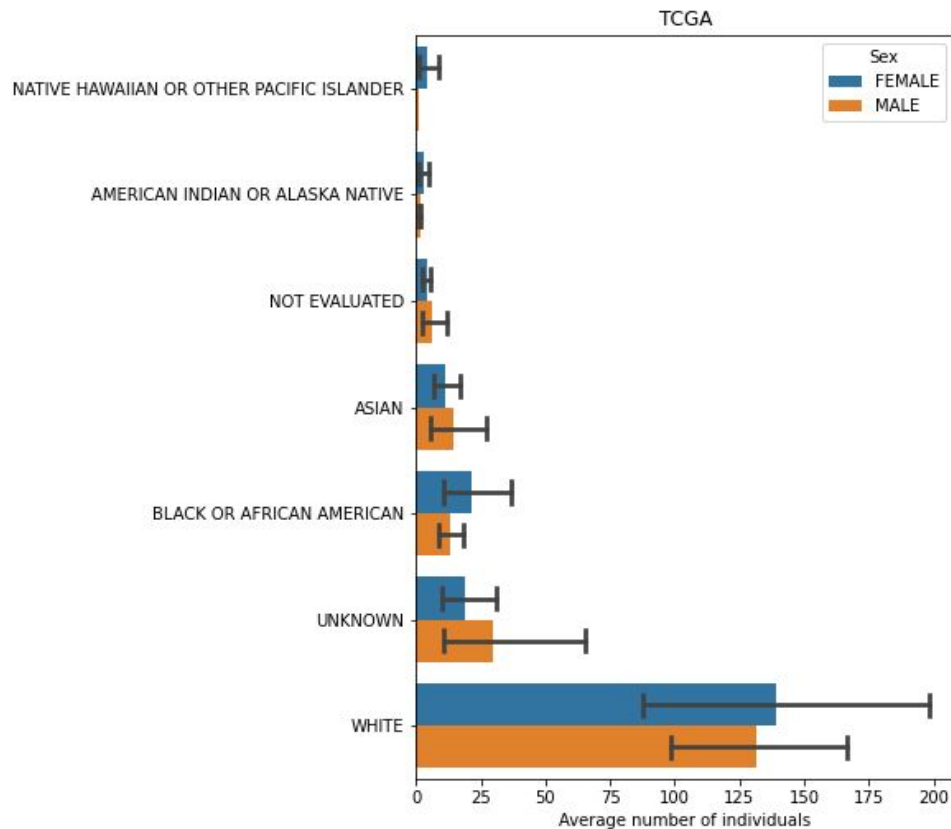
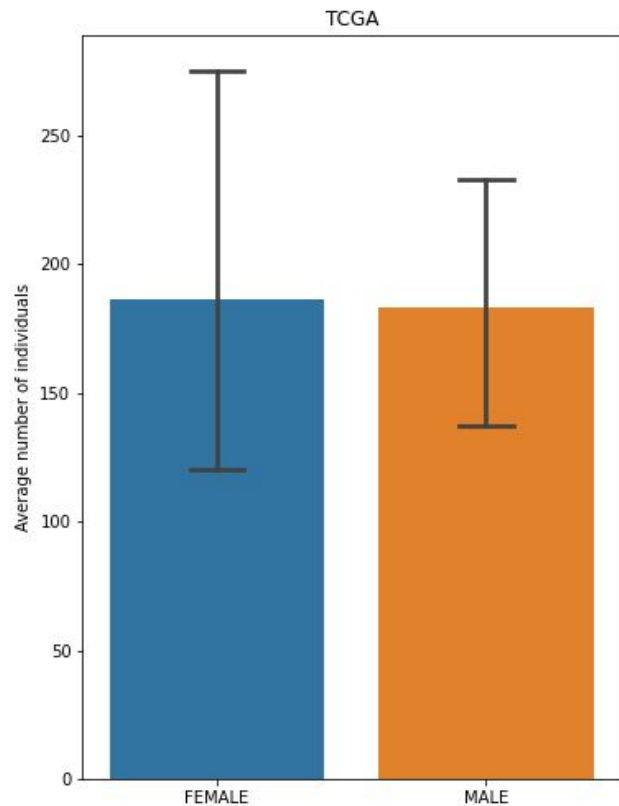
The functions of these genes are enriched in **synaptic signaling** and **nervous system development**.



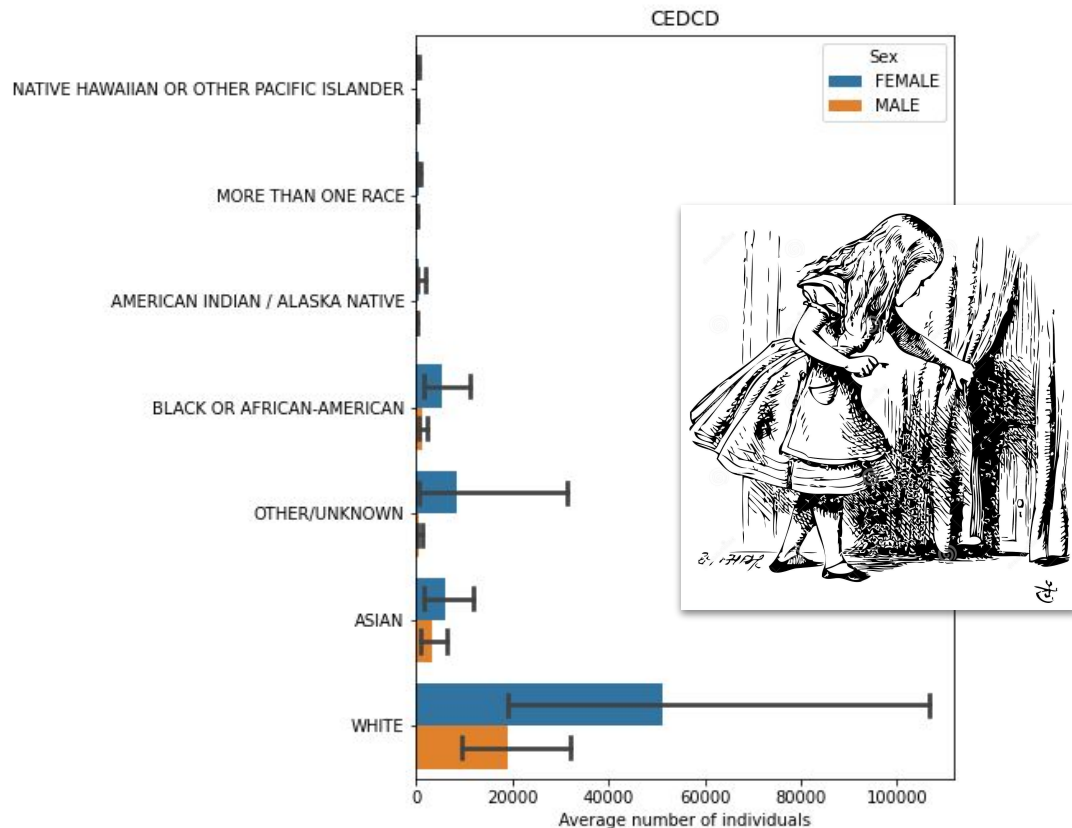
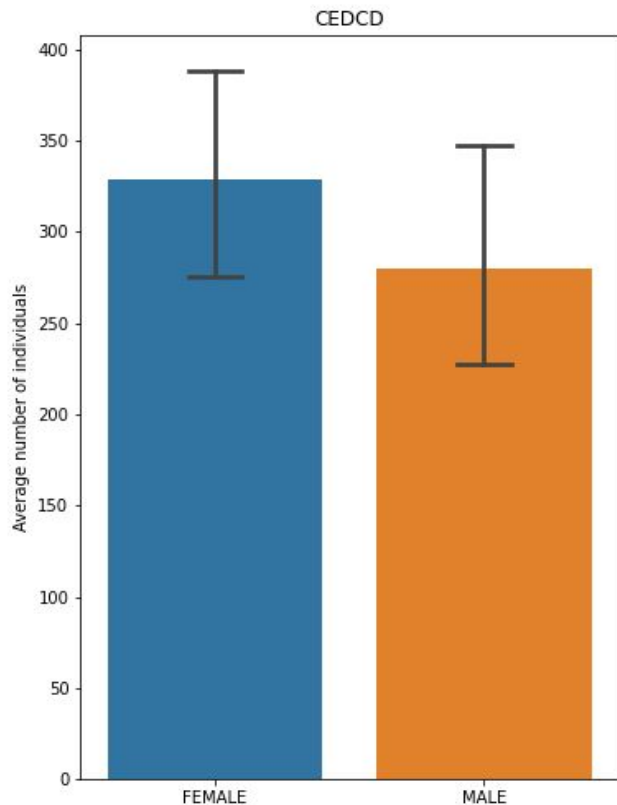
Explainable synthetic data generation for paediatric cancer



Sample size and label availability



Sample size and label availability



Conclusions

- **Multilayer networks** represent a powerful tool for heterogeneous data integration in **rare diseases** such as medulloblastoma.
- The study of multilayer **community structure at different scales** enables to detect strong associations between bio-entities.
- The study of **multilayer community trajectories** allows to accurately performing tasks such as dimensionality reduction and molecular interpretation.
- Explainable **synthetic data generation** enables to both augment the data and to identify genes that are relevant to data synthesis.

Alfonso Valencia
Iker Núñez-Carpintero
Alejandro Tejada Lapuerta
Mar Batlle Pérez
Hanns Lochmüller
Maria Rigau
Mattia Bosio
Emily O'Connor
Salvador Capella
Steve Laurie
Sergi Beltran
Yoshiteru Azuma
Ana Topf
Rachel Thompson
Peter-Bram van Hoen
Ivailo Tournev
Velina Guerguelcheva
Carolina Armengol
Marianyela Petrizzelli
Andrei Zinovyev
Anaïs Baudot
Alberto Valdeolivas
Léo Pio-Lopez



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Thank you



davide.cirillo@bsc.es