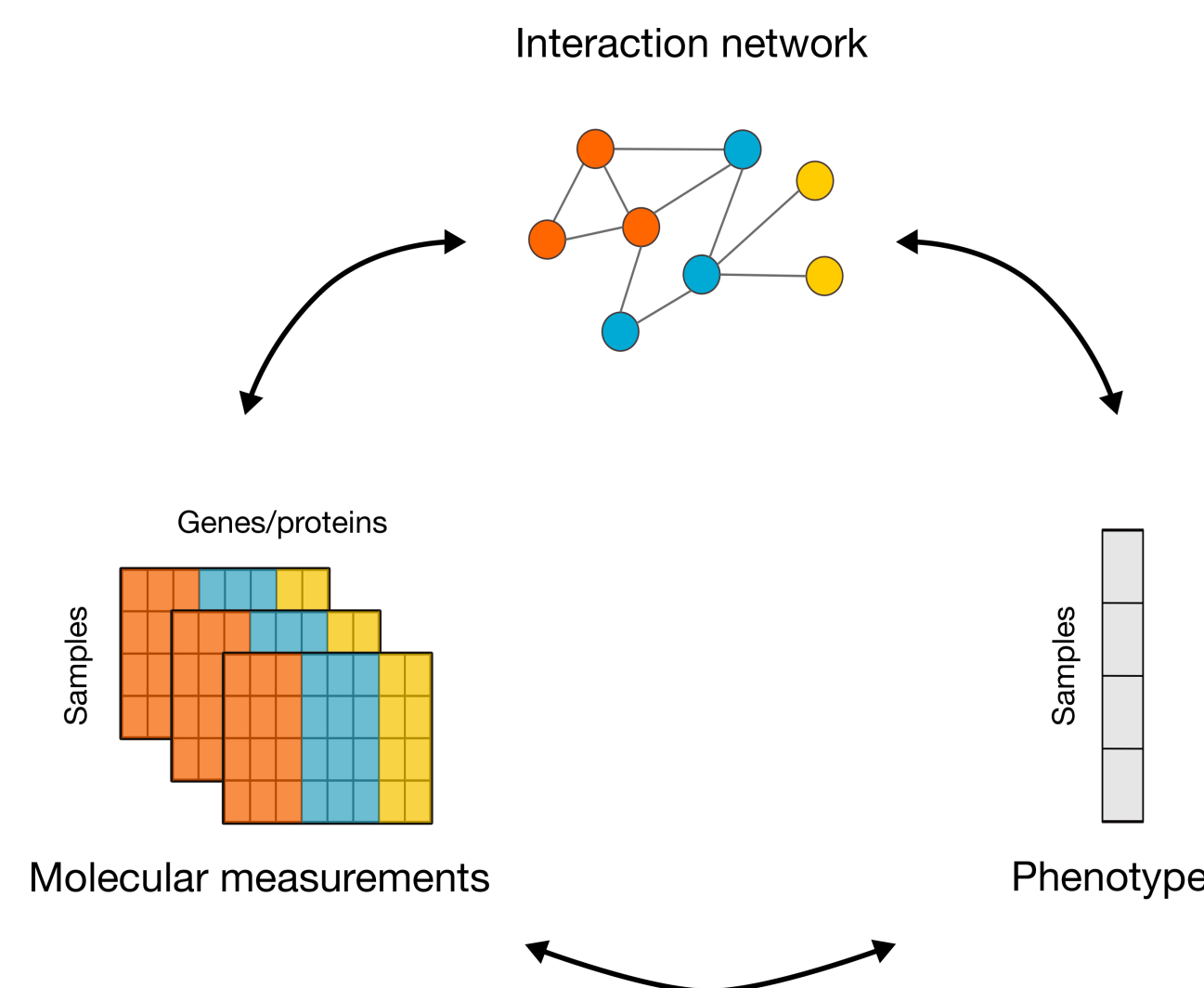


Motivation

What

- Accurate **phenotype classification**
- Biological **interpretability**
- **Robustness** to noise
- Handling **curse of dimensionality**

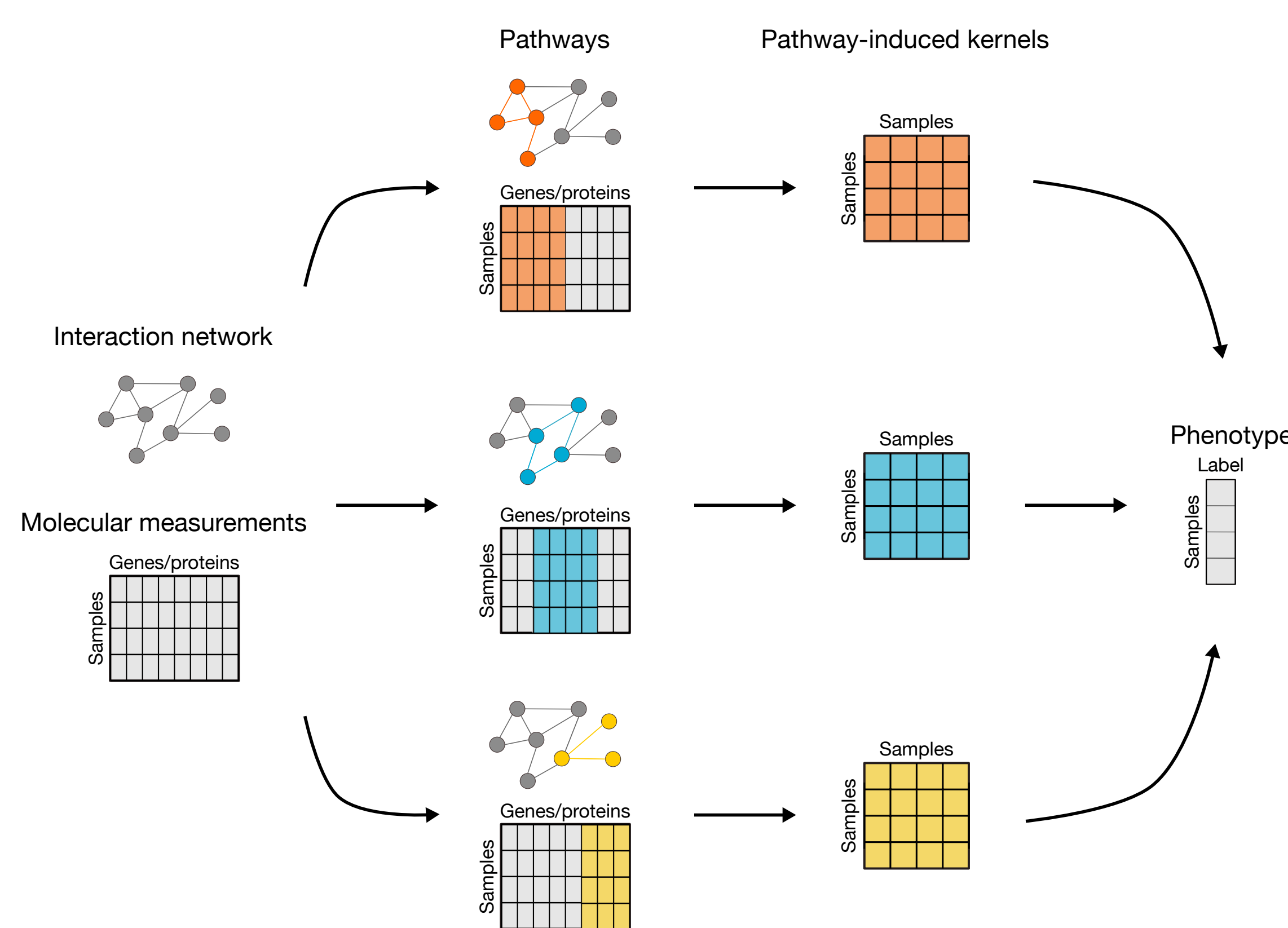


How

- Exploit **prior knowledge** from biological networks
- Apply **multiple kernel learning** for feature encoding
- Use **pathway annotations** to enable interpretability

PIMKL¹

Concept



Formulation

Given:
 $X \in \mathbb{R}^{n,m}$ molecular measurements
 $y \in \mathbb{R}^n$ phenotype of interest
 \mathcal{P} set of pathways

Build *pathway-induced* kernels:

$$G_p = (V_p, E_p) \quad L_p = D_p - A_p$$

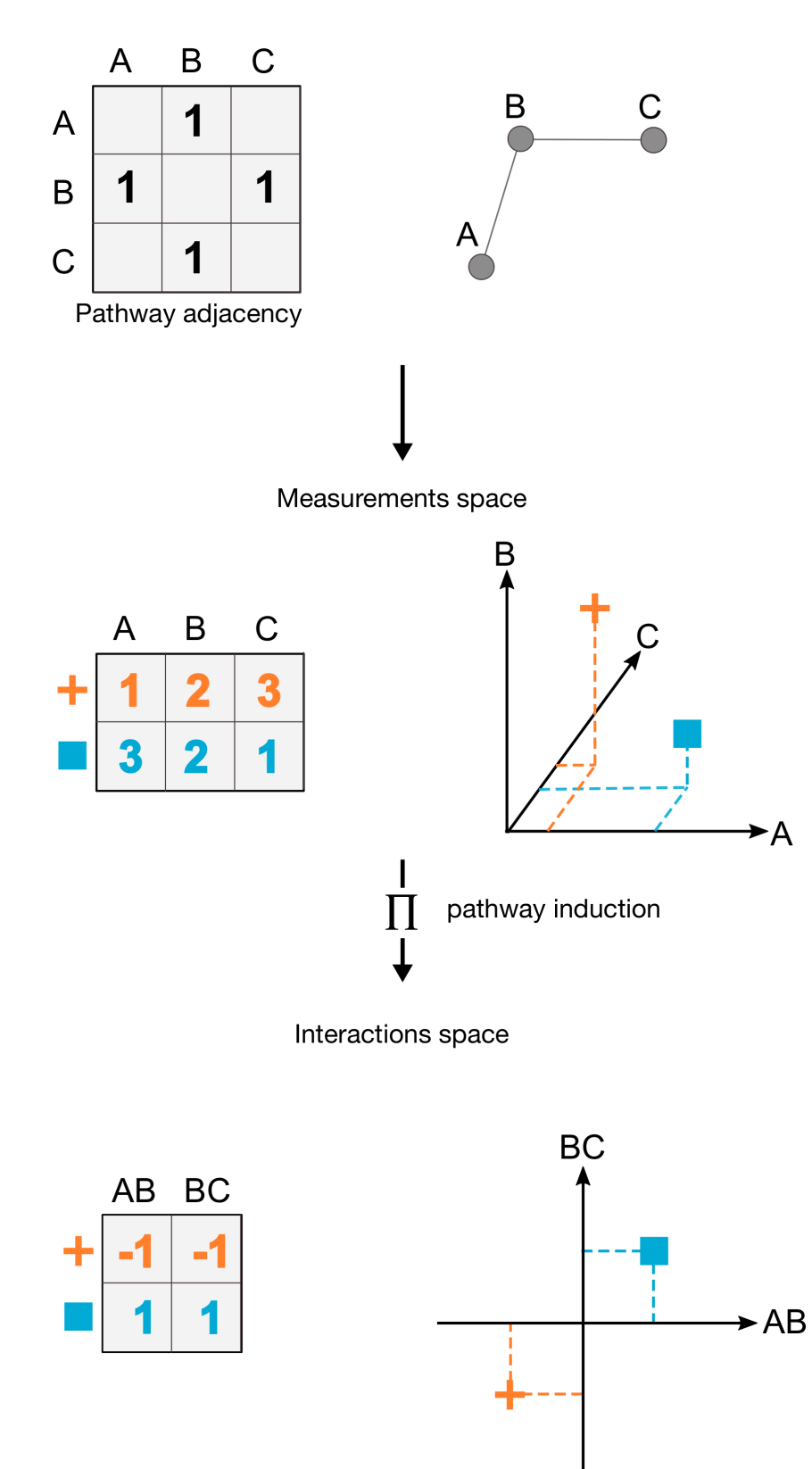
$$X_p \in \mathbb{R}^{n,m_p} \quad K_p = X_p L_p X_p^T$$

Consider the mixture of kernels:

$$K = \sum_{p \in \mathcal{P}} w_p K_p$$

Optimize the weights for phenotype classification (custom implementation of EasyMKL^{2,3})

Pathway induction



Results

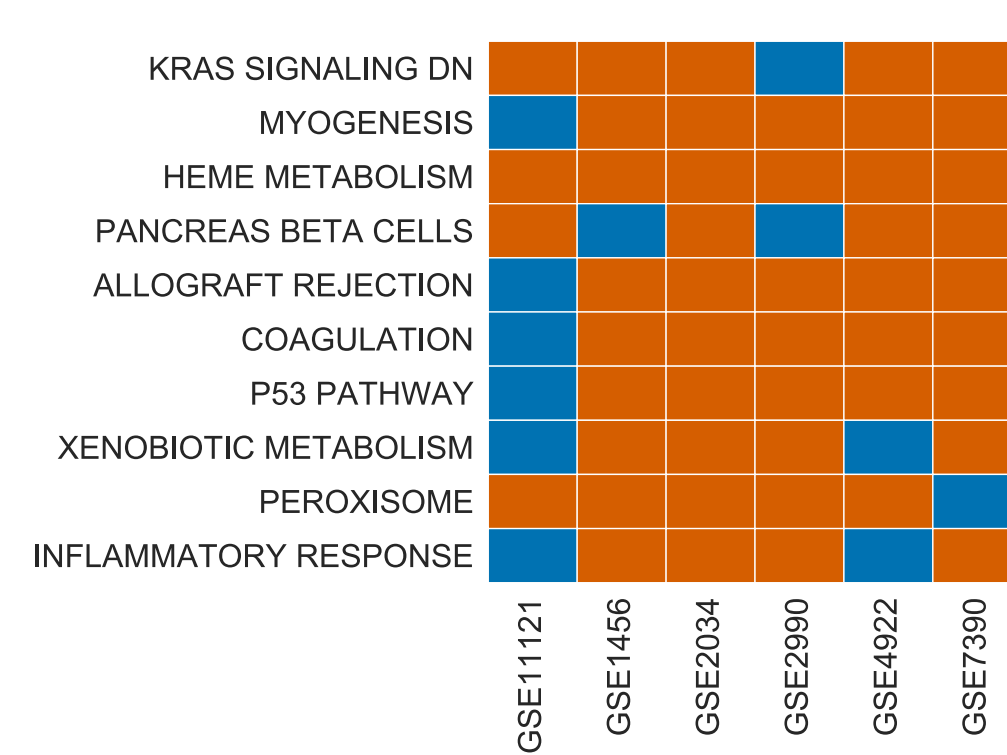
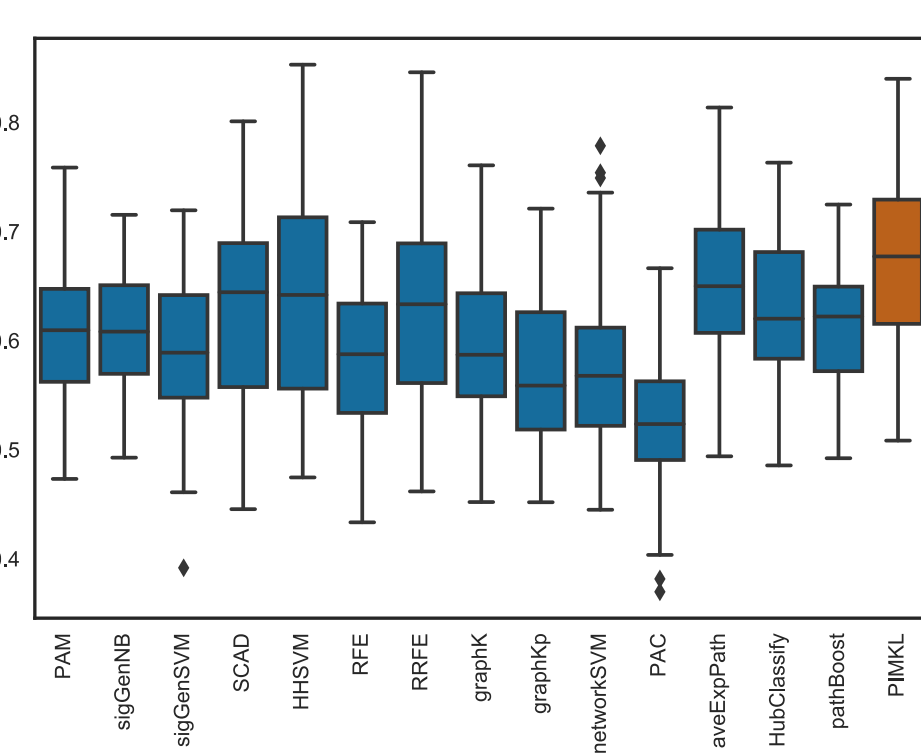
Data

- Affymetrix Human Genome U133A Array from 6 breast cancer microarray cohorts⁴
- A collection of 50 hallmark gene sets from Molecular Signatures Database (MSigDB) version 5.2⁵
- Molecular networks from KEGG^{6,7,8} and Pathway Commons⁹
- Illumina Human v3 microarray data and Affymetrix SNP 6.0 copy number data from METABRIC¹⁰

Benchmark

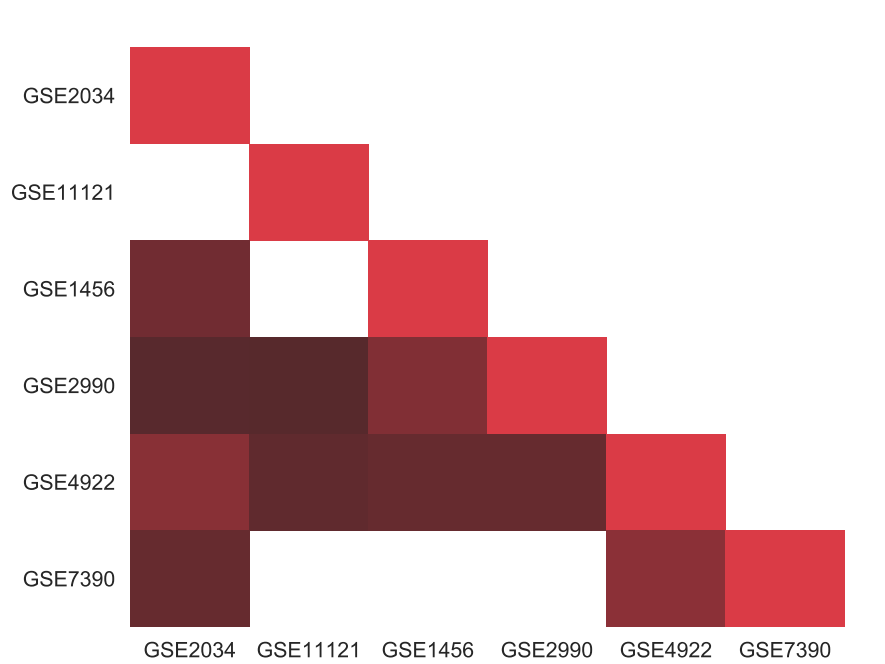
Comparison on a benchmark⁴ against 14 methods exploiting prior knowledge

Weights learned in different cohorts highlight relevant pathways for breast cancer

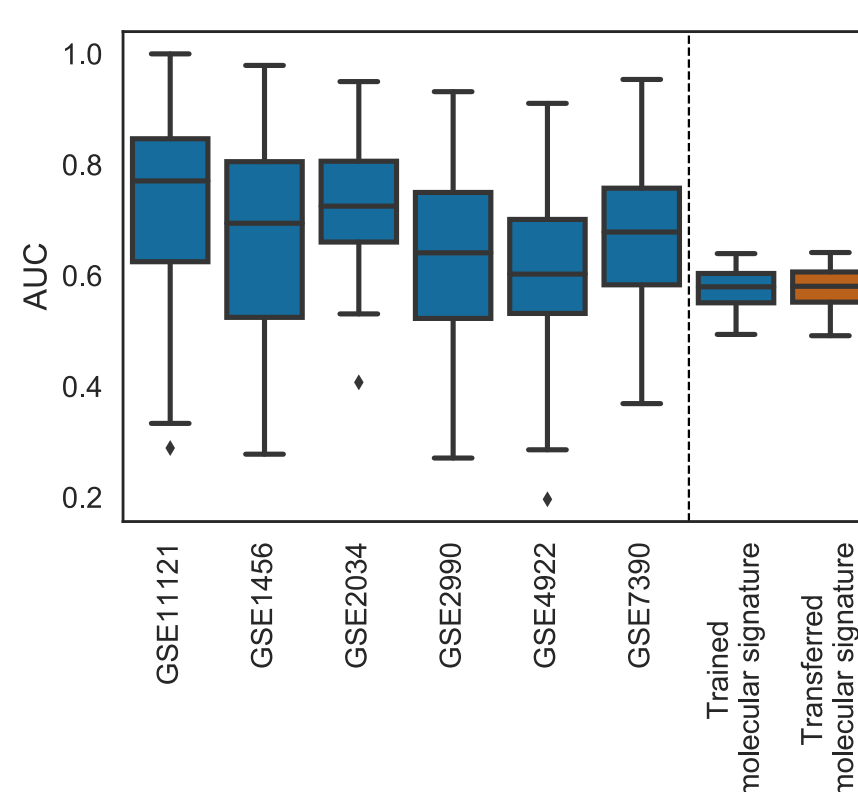


Transfer learning

Pathway signatures are significantly correlated across cohorts

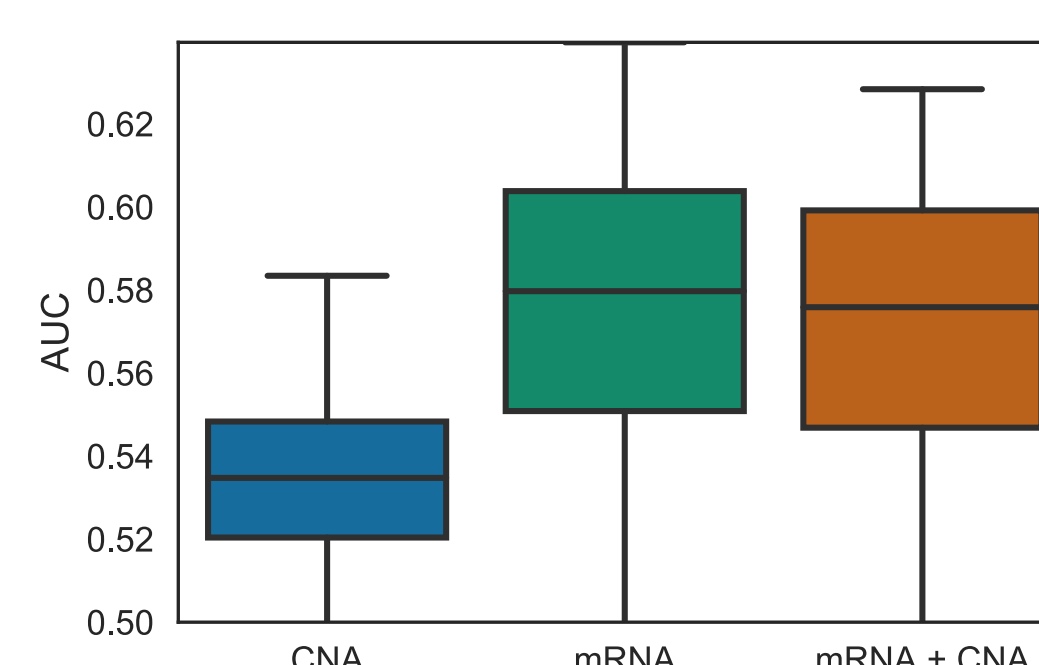


Signature learned on other cohorts successfully transferred to METABRIC



Multi-omics

In METABRIC case noisy kernels from CNA are integrated without significant performance loss



Availability

PIMKL is available as an open access web service on IBM Cloud



Find out more on the project website



<https://ibm.biz/ibmpimkl>

- 1) Manica M., et al., PIMKL: Pathway Induced Multiple Kernel Learning, npj Systems Biology and Applications, 2019
- 2) Aiolfi F., et al., A kernel method for the optimization of the margin distribution, Lecture Notes in Computer Science, 2008
- 3) Aiolfi F., et al., EasyMKL: A scalable multiple kernel learning algorithm, Neurocomputing, 2015
- 4) Cun Y., et al., Prognostic gene signatures for patient stratification in breast cancer-accuracy, stability and interpretability of gene selection approaches using prior knowledge, BMC bioinformatics, 2012
- 5) Liberzon A., et al., The Molecular Signatures Database Hallmark Gene Set Collection, Cell Systems, 2015
- 6) Kanehisa M., et al., KEGG: Kyoto encyclopedia of genes and genomes, Nucleic Acids Research, 2000
- 7) Kanehisa M., et al., KEGG as a reference resource for gene and protein annotation, Nucleic Acids Research, 2016
- 8) Kanehisa M., et al., KEGG: new perspectives on genomes, pathways, diseases and drugs, Nucleic Acids Research, 2017
- 9) Cerami E. G., et al., Pathway Commons, a web resource for biological pathway data, Nucleic Acids Research, 2011
- 10) Curtis C., et al., The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups, Nature, 2012