

Strategy to remove batch effect between different omic data types



PRINCIPE FELIPE
CENTRO DE INVESTIGACION



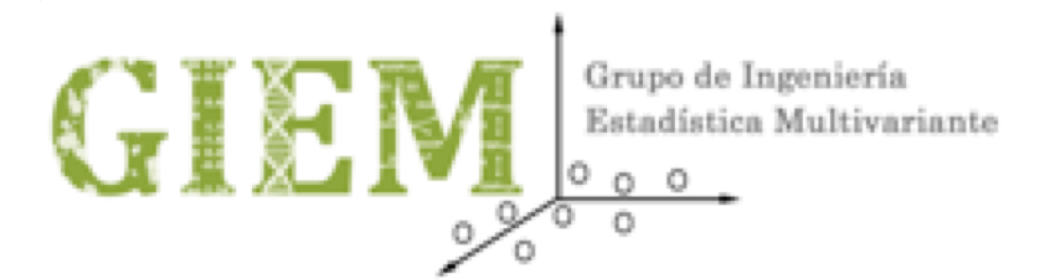
Manuel Ugidos^{1,*}, Sonia Tarazona^{1,2}, J.M. Prats-Montalbán²,
Alberto Ferrer², Ana Conesa^{1,3}

¹ Genomics of Gene Expression Laboratory, Centro de Investigación Príncipe Felipe, Valencia, Spain

² Department of Applied Statistics, Operations Research and Quality, Universitat Politècnica de València, Spain

³ Microbiology and Cell Science Dep, Institute for Food and Agricultural Sciences, University of Florida, Gainesville, USA

* Student information: **PhD Programme:** Statistics and Optimization; **UPV e-mail:** mauggue@doctor.upv.es

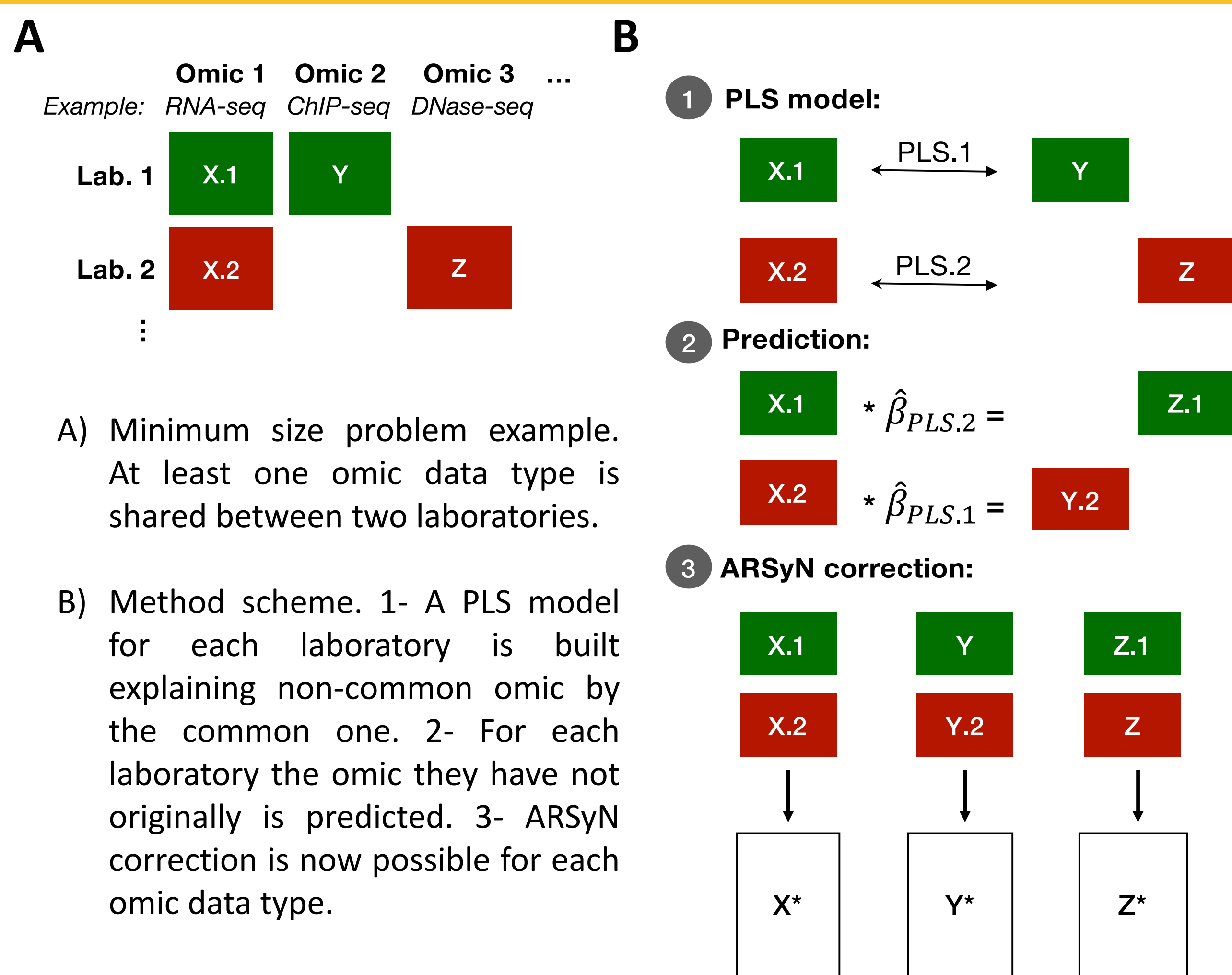


Introduction

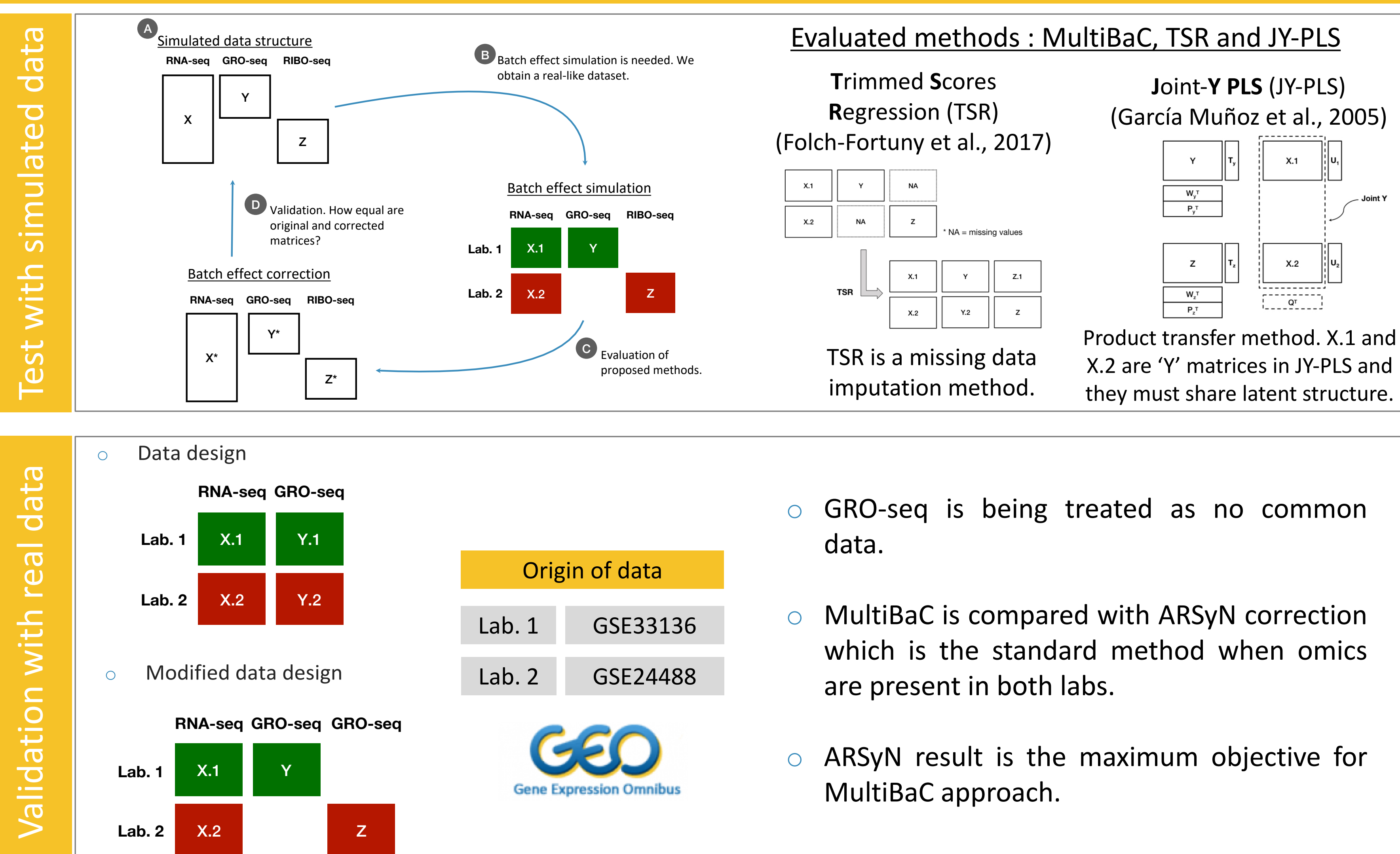
Omic technologies have expanded in diversity in the last years and the number of omics integration analysis possibilities has also increased. However, the costs of the different techniques are still high and most of research groups cannot afford research projects where many different omics techniques are analyzed. Nevertheless, as most research share their data in public repositories, there is a possibility of utilization of datasets from other laboratories to construct a multiomic study. An important issue when we want to integrate data from different studies is the batch effect. There are already several methods described which are able to correct batch effect on common omic data between different studies (e.g. ARSyN from M.J. Nueda et al., 2012) but they cannot be used to correct no common data (i.e. the omic data modality that has been analyzed at only one lab). We have developed **MultiBaC**, a strategy to correct batch effect on no common omic information which let us integrate different omic data types from different studies.

Materials & Methods

MultiBaC (Multi-omic Batch Correction) method

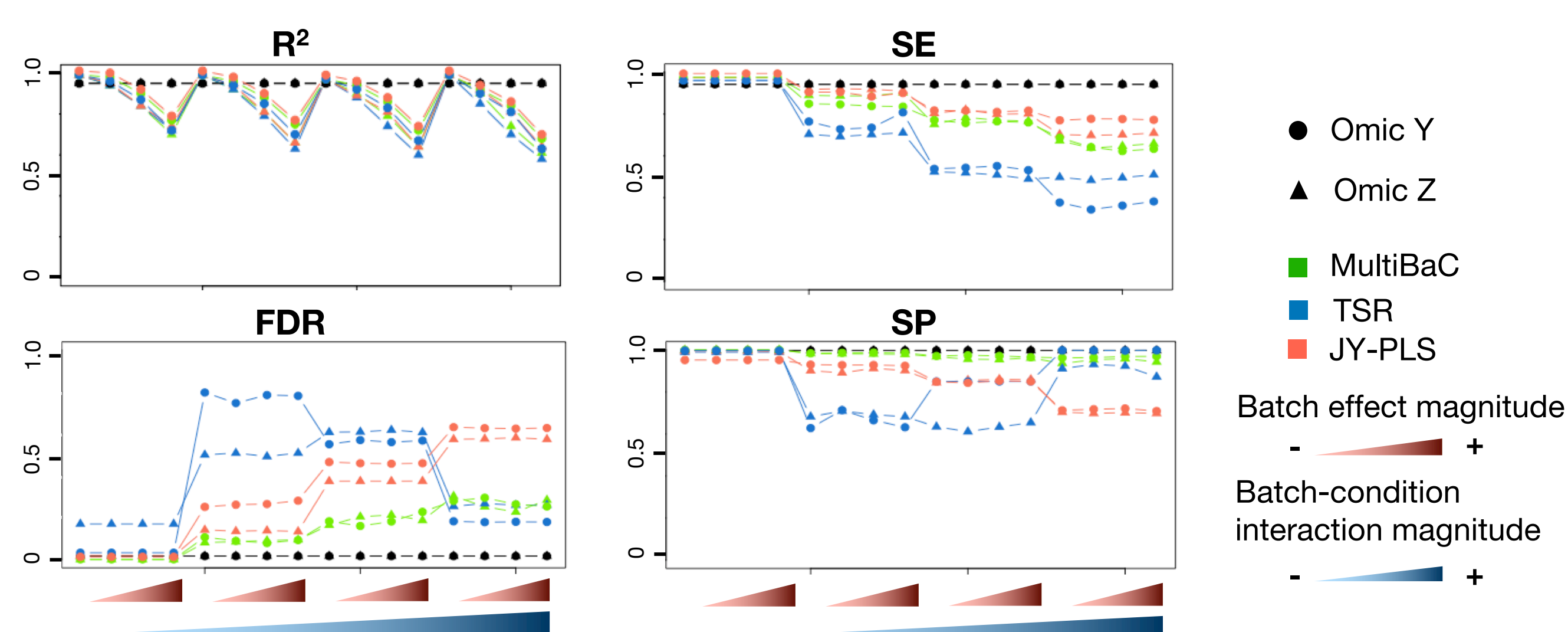


Test and Validation

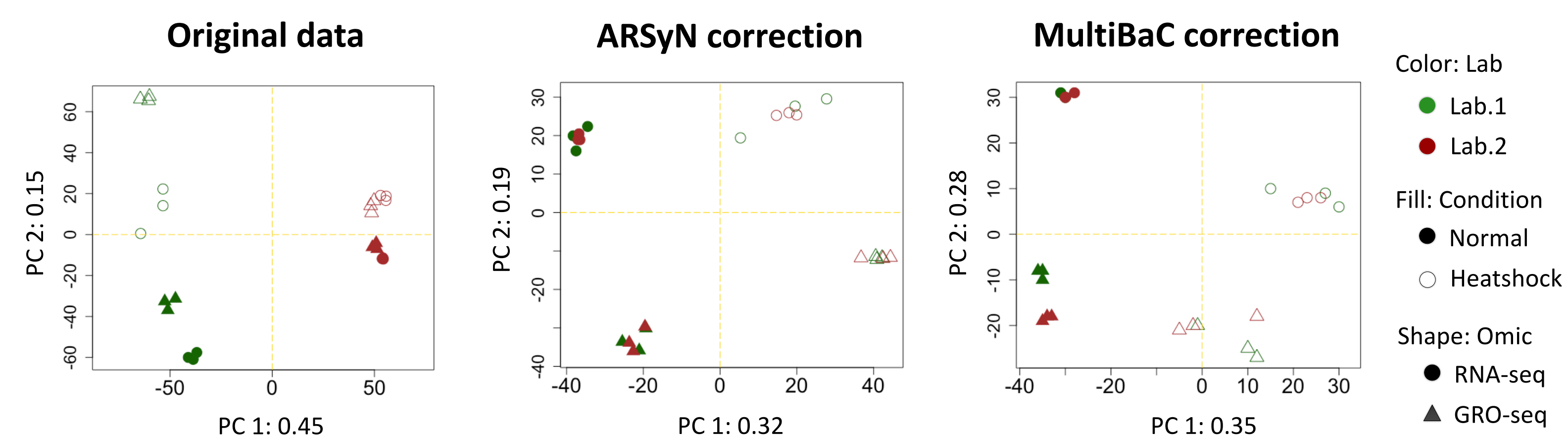


Results

With simulated data



With real data



- Measures of similarity between original simulation (no batch effect) and corrected matrices:
 - R²: Latent structure.
 - FDR, Sensitivity and Specificity: Differences in significant differentially expressed genes.
- MultiBaC reaches the best performance in comparison with TSR and JY-PLS.
- Real batch effects are not supposed to be as high as maximum magnitudes tested, so MultiBaC is suitable for real cases of batch effect.

- PCAs showing sources of variability.
- MultiBaC correction is not as perfect as ARSyN one but after MultiBaC correction the batch is not an important source of variability. ARSyN result is the maximum level of correction and MultiBaC reaches almost the same result.

Conclusions

- We have developed MultiBaC, a strategy to remove batch effect between different omic data types coming from different studies.
- MultiBaC approach reaches almost the same performance as ARSyN method which is the maximum level of correction, when both can be applied.
- When there is a no common omic MultiBaC is able to correct low and moderate batch effect magnitudes.
- MultiBaC does not work well with high batch effect and interaction magnitudes but these so high magnitudes have not been seen in real datasets.

Acknowledgments

This work is part of a research project that is totally funded by Generalitat Valenciana through PROMETEO grants programme for excellence research groups.

References

[1] Maria j. Nueda, Alberto Ferrer, Ana Conesa; ARSyN: a method for the identification and removal of systematic noise in multifactorial time course microarray experiments, *Bioinformatics*, 13 (2012) 553–566.

[2] A.Folch-Fortuny, F.Arteaga, A.Ferrer, Missing data imputation toolbox for MATLAB. *Chemometrics and Intelligent Laboratory Systems*. 154 (2016) 93-100.

[3] S. Garcia-Munoz, J. F. MacGregor, T. Kourti, Product transfer between sites using Joint-Y PLS. *Chemometrics and Intelligent Laboratory Systems*. 79(2005) 101-114.