

1 **S1 Supporting Information**

2

3 **Authors:**

4 Cecilia Wieder¹, Juliette Cooke², Clement Frainay², Nathalie Poupin², Russell Bowler³,
5 Fabien Jourdan⁴, Katerina J. Kechris⁵, Rachel PJ Lai⁶, Timothy Ebbels*¹

6 **Affiliations:**

7 ¹Section of Bioinformatics, Division of Systems Medicine, Department of Metabolism, Digestion, and
8 Reproduction, Faculty of Medicine, Imperial College London, London, United Kingdom

9 ²Toxalim (Research Centre in Food Toxicology), Université de Toulouse, INRAE, ENVT, INP-Purpan, UPS,
10 Toulouse, France

11 ³National Jewish Health, 1400 Jackson Street, Denver, CO, 80206, USA

12 ⁴MetaboHUB-Metatoul, National Infrastructure of Metabolomics and Fluxomics, Toulouse, France

13 ⁵Department of Biostatistics and Informatics, Colorado School of Public Health, University of Colorado Anschutz
14 Medical Campus, Aurora, CO, United States of America

15 ⁶Department of Infectious Disease, Faculty of Medicine, Imperial College London, London, United Kingdom

16

17 Related work

18 DIABLO [1] is a supervised method for multi-omics data integration based on generalised
19 canonical correlation analysis (GCCA). It uses singular value decomposition to find a lower-
20 dimensional representation of multiple omics input matrices and selects correlated
21 variables which are associated with the phenotype of interest. It requires the user to specify
22 a design matrix, representing the expected correlation between omics datasets in the model.
23 The inputs for DIABLO are scaled N-by-M omics data matrices, rendering it also compatible
24 with pathway-transformed data matrices.

25 MOGSA [2] is an unsupervised method for multi-omics data integration, designed to output
26 a matrix of multi-omics single-sample pathway scores. It begins by integrating the data at
27 the molecular level using multiple-factor analysis, followed by projecting a binary matrix of
28 pathway-membership information onto the observations in the latent space, and finally
29 multiplying together the latent space matrices of samples and pathways to produce an N-by-
30 P pathway score matrix. The final pathway score matrix can be decomposed to investigate
31 the contribution of each omics dataset. MOGSA, unlike PathIntegrate and DIABLO, is not a
32 predictive model but rather a method for generating multi-omics pathway scores, which
33 could be used as input to predictive models like PathIntegrate.

34 Like MOGSA, Multi-Omics Pathway Analysis (MOPA) [3] generates pathway-score matrices
35 using non-negative tensor decomposition. It is designed for gene-based omics data such as
36 mRNA, methylation, and miRNA data. MOPA uses a two-step process for generating pathway
37 scores, firstly it employs a non-negative tensor decomposition to perform feature selection
38 to find genes significantly associated with a phenotype, and secondly computes pathway
39 scores using these genes with a method similar to Gene Set Variation Analysis [4]. Like
40 MOGSA, MOPA allows the calculation of an 'omics contribution rate', to understand how
41 different omics contribute to pathway score calculation.

42 Multi-Omics Factor Analysis (MOFA) [5] is an unsupervised latent-variable method for
43 multi-omics data integration. It uses group factor analysis to decompose multiple omics
44 matrices into loadings and score matrices, which can be sparse. MOFA could be used with
45 ssPA score matrices as input, to form an unsupervised pathway-based multi-omics
46 integration model. Similar to PathIntegrate, users can extract variable importances for each
47 latent factor and the contribution of each omics to each factor.

48 Lilikoi 2.0 [6] is a metabolomics-specific pathway-based deep learning model. It uses
49 Pathifier [7] to produce ssPA scores which are then input to a deep neural network, or other
50 classifiers such as random forest or logistic regression. It offers prognosis prediction using a
51 Cox proportional hazards model, as well as network-based pathway visualisation options for
52 downstream analysis.

53 PathwayPCA [8] is a toolkit offering multiple pathway-analysis based utilities: 1) testing
54 pathway association with an outcome (similar to conventional pathway analysis), 2)
55 extracting important genes within a pathway using sparse modelling, 3) compute pathway
56 scoring on important genes, which can be used as input for multi-omics analysis. The
57 pathway scores are computed using Adaptive, Elastic-net, Sparse PCA) or Supervised PCA

58 (SuperPCA), introduced by the same authors. Similar to Lilikoi, the pathway-transformed
59 output can be input to various downstream analysis such as survival analysis.

60 Integrative directed random walk-based method utilizing pathway information (iDRW)
61 [9,10] is a method for generating ssPA scores based on utilising gene-gene topological
62 interactions within pathways. Combining a gene-gene directed graph based on KEGG
63 pathways and a random walk algorithm, iDRW was used to integrate gene-expression and
64 copy number alteration data, resulting in a pathway score matrix. The authors demonstrated
65 using iDRW scores improved survival prediction compared to molecular-level data as well
66 as other ssPA scoring approaches.

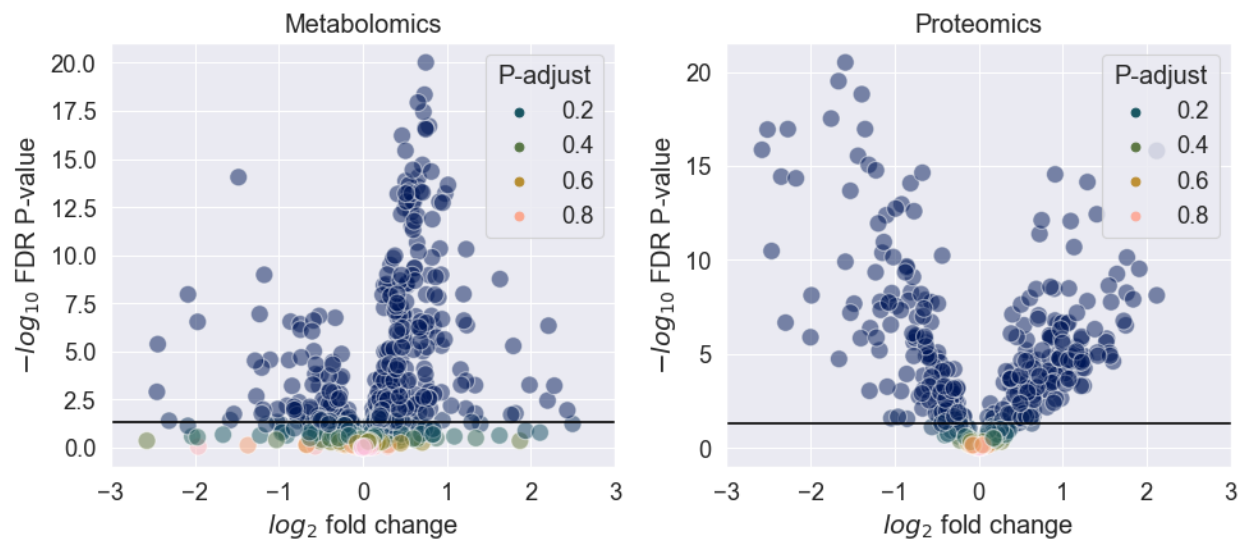
67 Finally, we refer the interested reader to a comprehensive review by Maghsoudi et al. [11]
68 which provides a systematic evaluation of 32 integrative pathway analysis methods. While
69 the aforementioned methods all provide useful functionality for either multi-omics
70 integration at the molecular level (DIABLO, MOFA), or the generation of pathway scores at
71 either the single omics level (Lilikoi, PathwayPCA), or the multi-omics level (MOGSA, MOPA,
72 iDRW), none of these provide a framework for pathway-based multi-omics data integration.
73 PathIntegrate seeks to fill this gap, providing a user-friendly Python implementation of the
74 Multi-View and Single-View frameworks which a) generate multi-omics pathway scores
75 (based on the user's choice of ssPA methods), and b) apply state-of-the-art predictive models
76 to identify perturbed pathways. Furthermore, the majority of methods for generating multi-
77 omics pathway scores are not designed to incorporate metabolomics data and are primarily
78 based on gene/protein identifiers. PathIntegrate is specifically designed for (but not limited
79 to) the integration of metabolomics data alongside other omics, providing multi-omics
80 pathways containing gene (ENSEMBL), protein (UniProt), and metabolite (ChEBI)
81 identifiers. Finally, to enhance ease-of-use and seamless integration with other pipelines,
82 PathIntegrate models are compatible SciKit-Learn estimators, enabling the use of various
83 predictive models and parameter optimisation functions available in the SciKit-Learn
84 library.

85

86

87 **Supplementary figures and tables**

Mild vs. severe COVID-19 volcano plots



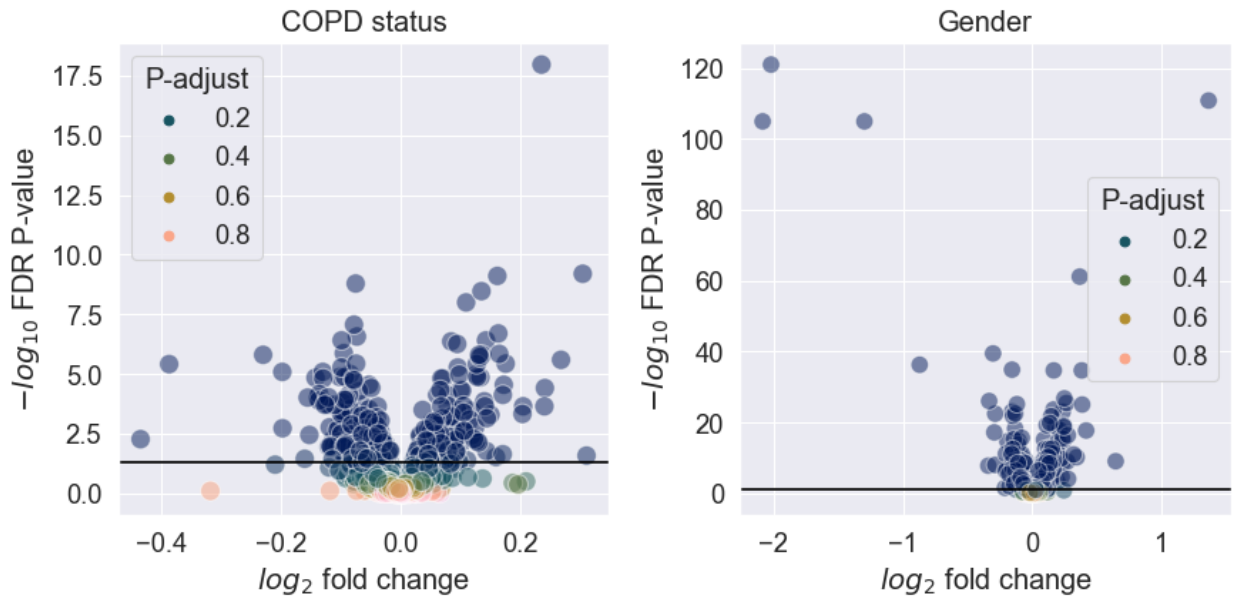
88

89 **Fig A in S1 Supporting Information:** Fold changes in COVID-19 multi-omics data based on
90 *outcome (mild vs. severe cases).*

91

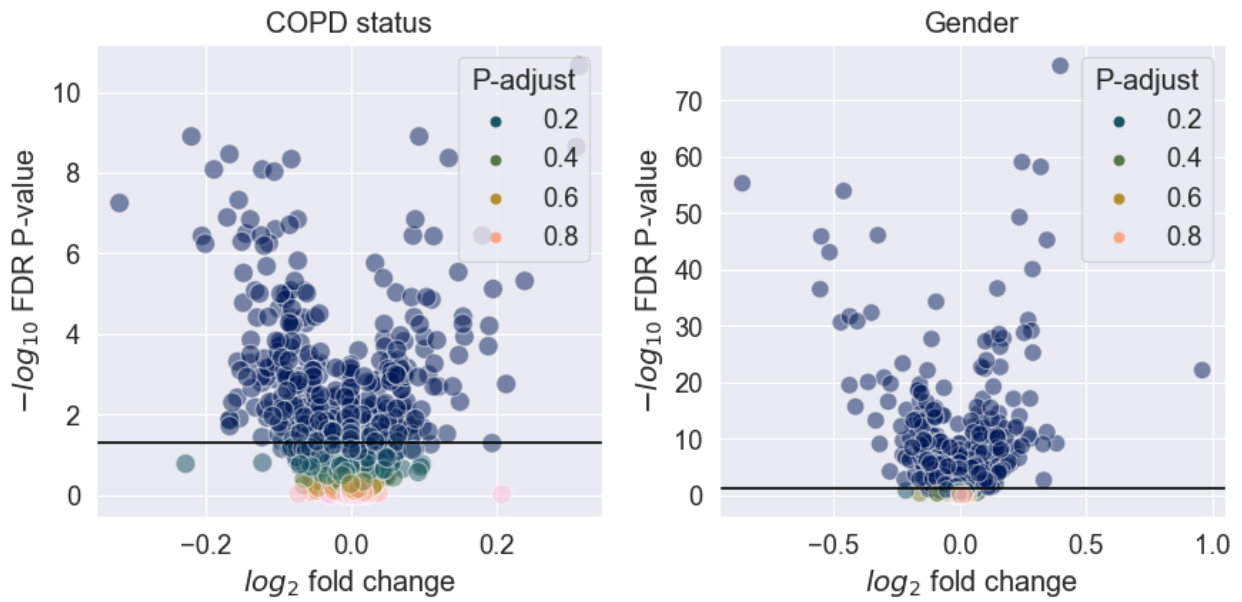
92

COPDgene proteomics volcano plot



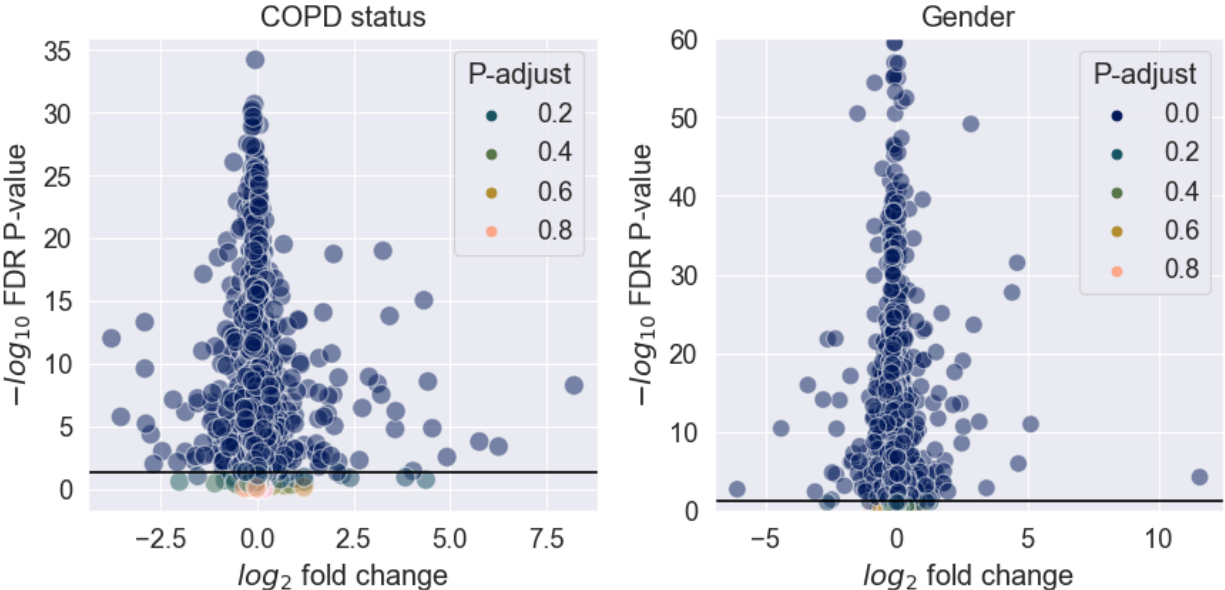
93

COPDgene metabolomics volcano plot

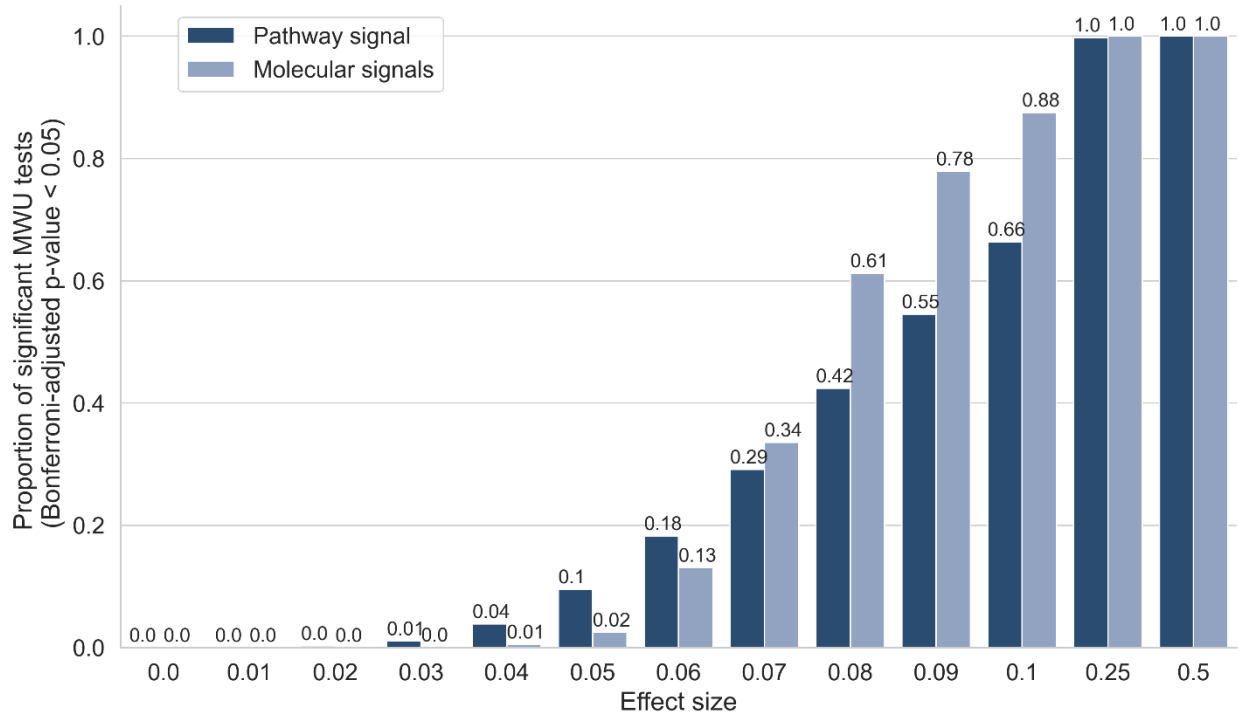


94

COPDgene transcriptomics volcano plot



95
96 **Fig B in S1 Supporting Information:** Fold changes in COPDgene multi-omics data based on
97 either COPD status or gender outcomes.
98



100

101 **Fig C in S1 Supporting Information: Pathway transformation enhances sensitivity to**
 102 **low signal-to-noise signals (COPDgene semi synthetic data).** Y axis shows proportion of
 103 MWU tests significant at Bonferroni $p \leq 0.05$, performed either on the pathway-level data or
 104 the molecular level data, at varying effect sizes shown on X-axis.

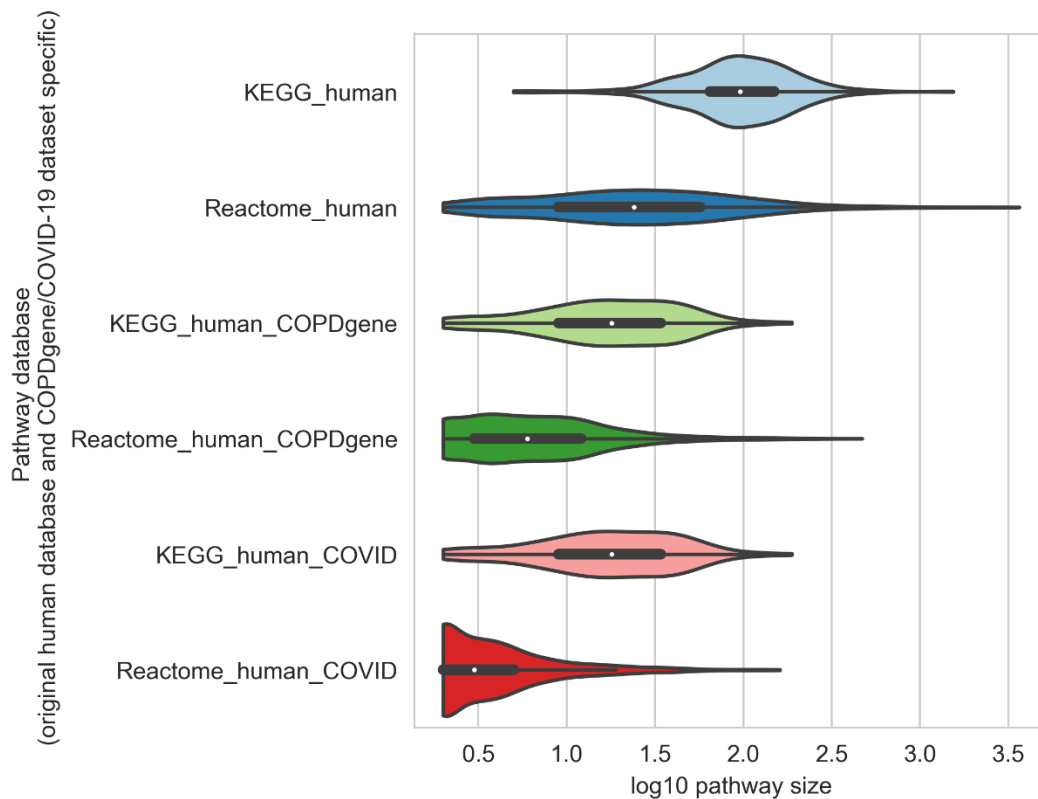
105

106 Pathway database influences model performance

107 The performance of pathway-based models is strongly dependent on the pathway
108 definitions used. The number and composition of pathways varies between databases, and
109 factors such as size, level of overlap, ratio of compounds to proteins/genes, etc. can all impact
110 the models. We investigated the size of pathways in the Reactome human versus the KEGG
111 human multi-omics pathway databases (those used in this work, where pathways can
112 contain a combination of metabolites, proteins, and genes), and found KEGG to contain on
113 average larger pathways (median size 96 molecules) than Reactome (median size 24
114 molecules). Reactome however contains more pathways (2,583) than KEGG (352).
115 Importantly, these pathway database statistics are influenced by the molecules profiled in
116 the dataset at hand, as only molecules that map to pathway identifiers will be included in the
117 modelling. We investigated the pathway size distribution in two datasets, COPDgene and
118 COVID-19 and found that the general trend was the same: KEGG pathways are generally
119 larger than Reactome pathways (*Fig D in S1 Supporting Information*).

120 We also investigated the pathway annotation levels of genes, proteins, and metabolites, i.e.
121 the percentage of molecules profiled in a dataset with a valid pathway database identifier
122 (ENSEMBL, Uniprot, or ChEBI) assigned to pathways. Although results are highly dataset and
123 assay-dependent, when considering the COPDgene and COVID-19 datasets and the Reactome
124 pathway database (*Table A in S1 Supporting Information*), we found proteomics data to have
125 the highest percentage of total molecules profiled mapping to pathways (>70% for both
126 datasets). Metabolomics data had the lowest percentage of molecules mapping to pathways
127 (16.9% for COPDgene and 23.9% for COVID-19). This is likely due to the specificity of the
128 ChEBI identifiers, particularly for chemical subclasses such as fatty acids, where molecules
129 i.e. lipids can be annotated to a very high level of specificity depending on side chain
130 composition etc, but these are not yet annotated to pathway databases at such a high level of
131 specificity. Bulk transcriptomics data was not available for the COVID-19 data, but in the
132 COPDgene dataset only 27% of ENSEMBL genes mapped to Reactome pathways,
133 demonstrating that the annotation issue is not specific only to metabolomics data, but can
134 also affect sequencing-based omics such as transcriptomics, where thousands of genes are
135 yet to be added to pathways.

136



137 **Fig D in S1 Supporting Information:** *Violin plots showing log10 pathway size for KEGG and*
 138 *Reactome human databases, both for the original databases as well as the database specific*
 139 *coverage (COPDgene and COVID-19). Pathways used are Reactome and KEGG human multi-*
 140 *omics pathways, containing both metabolites and proteins.*

141

142

143

144

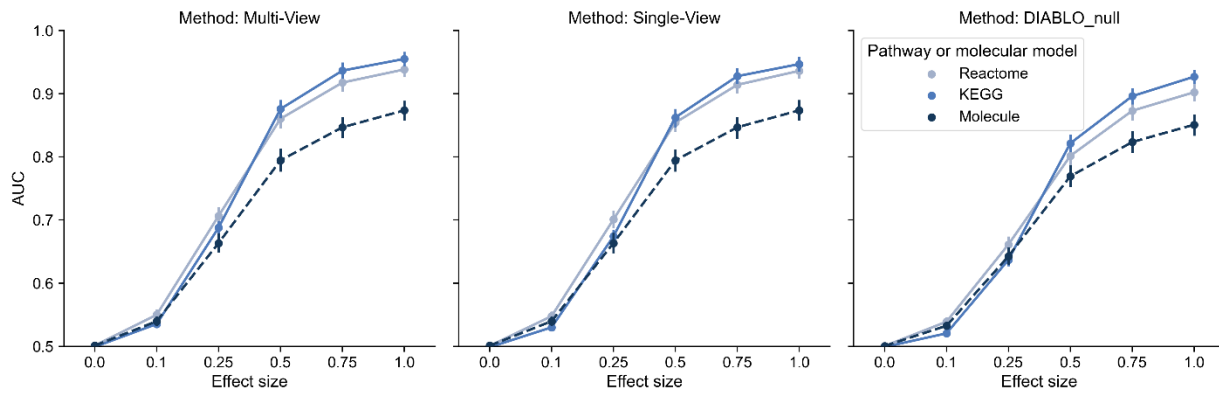
145

146 **Table A in S1 Supporting Information:** *Percentage of molecules with a valid identifier*
147 *(ChEBI, UniProt, or ENSEMBL) in single omics mapping to Reactome human pathways. A lower*
148 *percentage of molecules mapping to pathways means a greater percentage of molecules do not*
149 *yet map to pathways and are not incorporated into pathway-based analyses.*

	% of molecules with an identifier mapping to pathways		
Dataset	Metabolomics (ChEBI)	Proteomics (UniProt)	Transcriptomics (ENSEMBL)
COPDgene	16.9	81.5	27.5
COVID-19	23.9	77.9	NA (No transcriptomics data)

150

151



153 **Fig E in S1 Supporting Information.** Comparison of PathIntegrate methods classification
154 performance using KEGG and Reactome pathway databases as well as molecular-level model
155 based on semi-synthetic COPDgene data.

156

157

158

159

160

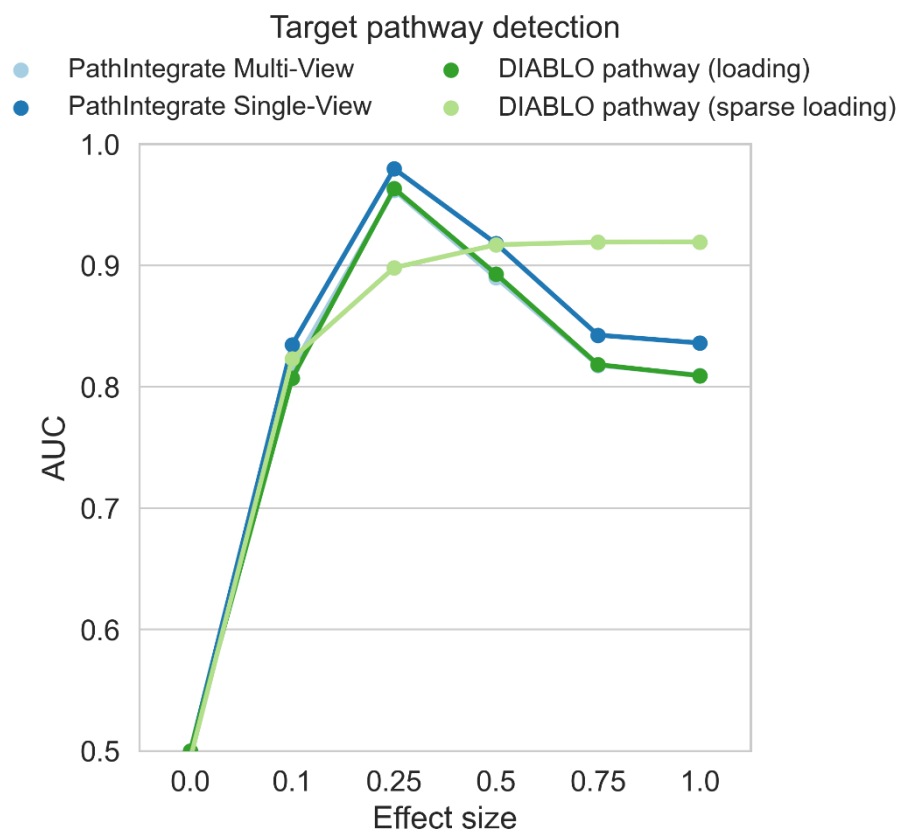
161

162

163

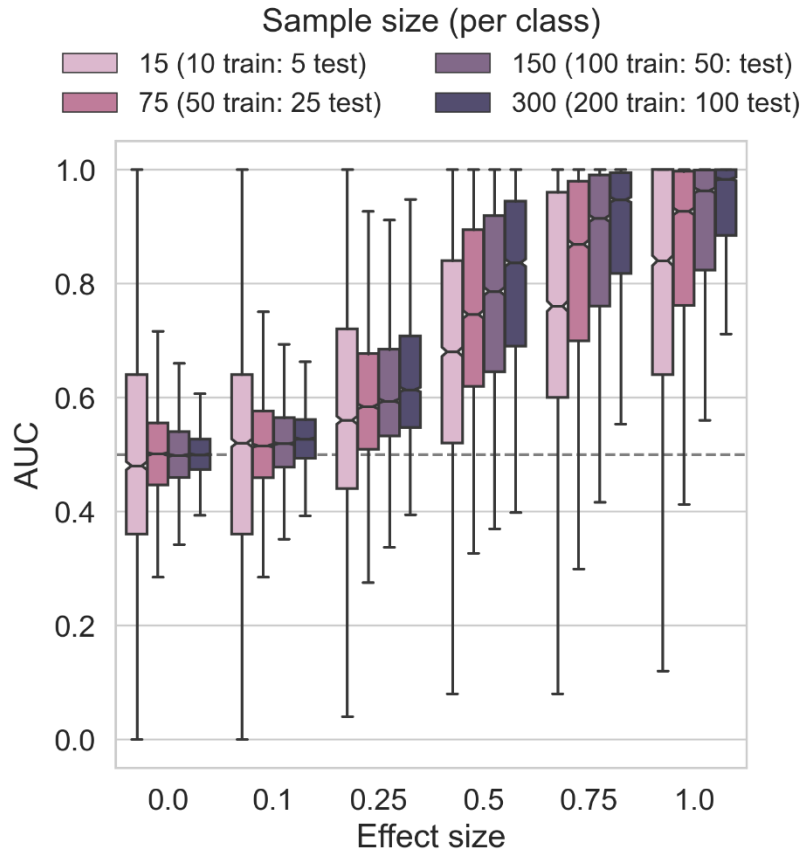
164

165



166 **Fig F in S1 Supporting Information.** Comparison of PathIntegrate and DIABLO full/sparse
 167 models ability to correctly recall target enriched pathway based on semi-synthetic COPDGene
 168 data. 'DIABLO pathway (loading)' uses an RGCCA model with no regularisation, whereas
 169 'DIABLO pathway (sparse loading)' uses an RGCCA model with L1 penalty.

170



171

172 **Fig G in S1 Supporting Information:** *Investigation of effect of sample size in PathIntegrate*
 173 *Single-View (PLS) classification performance on COPDgene data.*

174

175

176

177

178

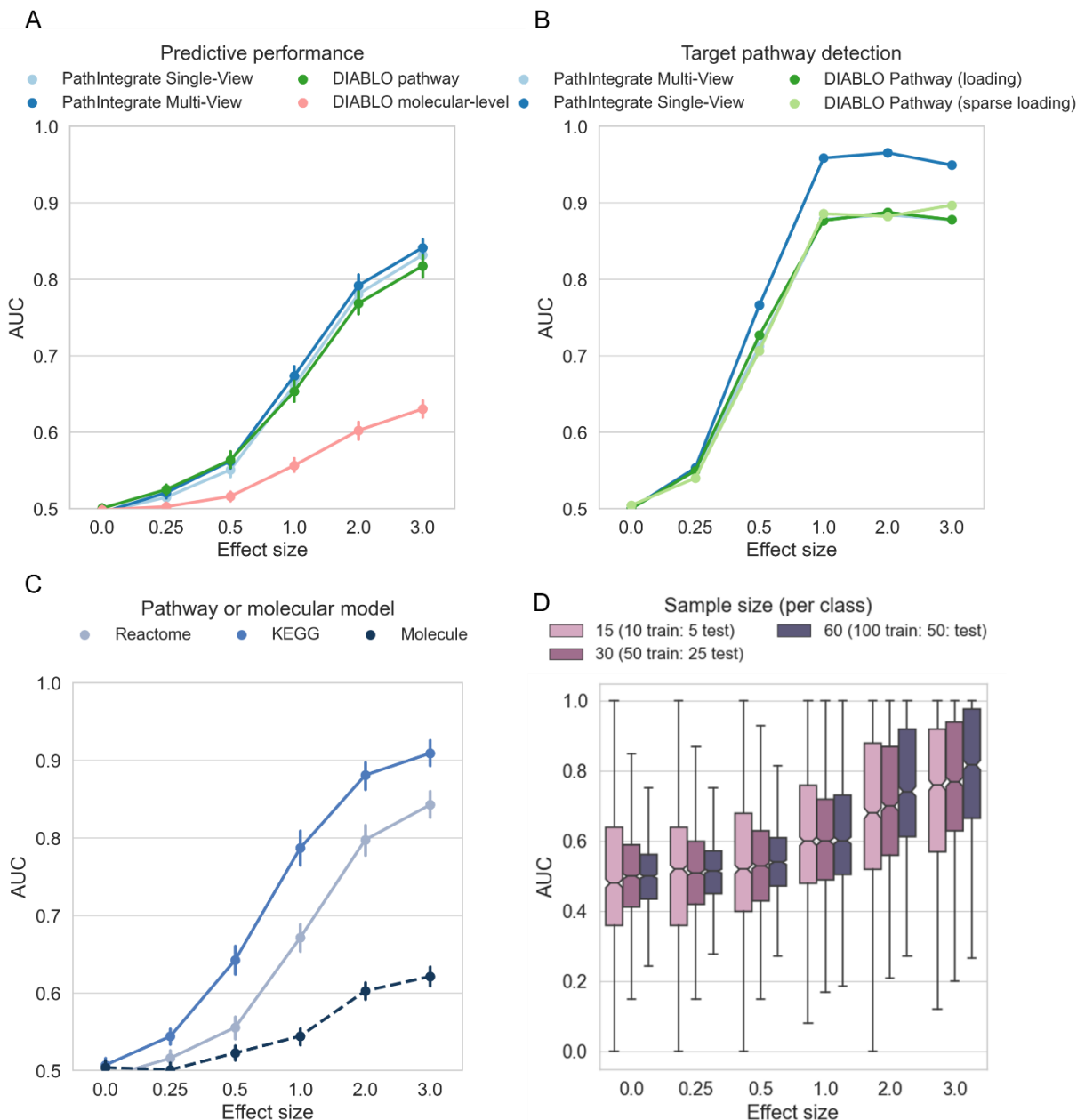
179

180

181

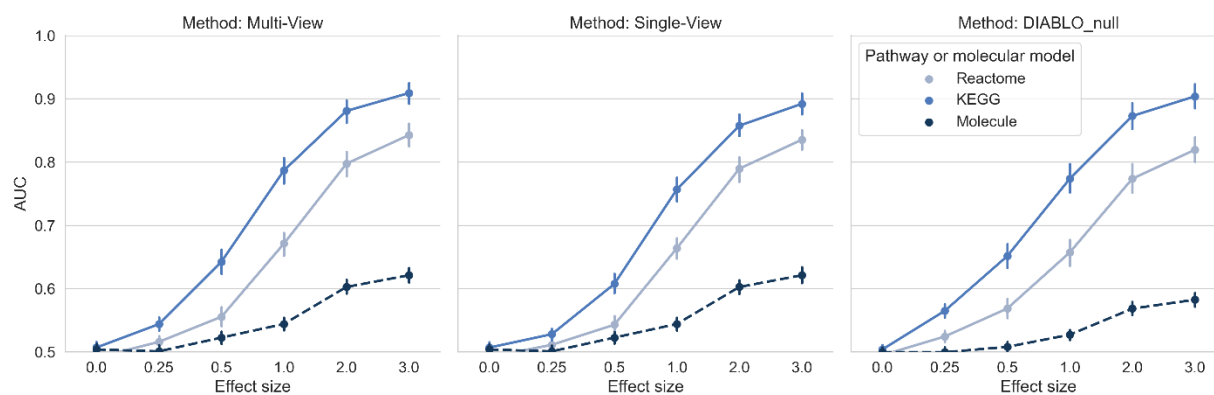
182

183



184 **Fig H in S1 Supporting Information: Performance of PathIntegrate and DIABLO vs.**
 185 **effect size, based on semi-synthetic data measured by AUROC. COVID-19 metabolomics**
 186 **and proteomics data were integrated in each model. A. Ability to correctly predict sample**
 187 **outcomes (case vs. control). We compared PathIntegrate Multi-View and Single-View to**
 188 **DIABLO using both molecular and pathway-level multi-omics data. B. Ability to correctly**
 189 **recall target enriched pathway. For 'DIABLO pathway' we compared the full RGCCA model**
 190 **loadings to the sparse model loadings for feature importance. C. Comparison of PathIntegrate**
 191 **Multi-View using KEGG and Reactome pathway databases as well as molecular-level model. D.**
 192 **Effect of sample size on PathIntegrate Multi-View classification performance. For panels A-C**
 193 **error bars indicate 95% confidence intervals on the mean AUROC (in some cases they appear**
 194 **smaller than point sizes).**

195

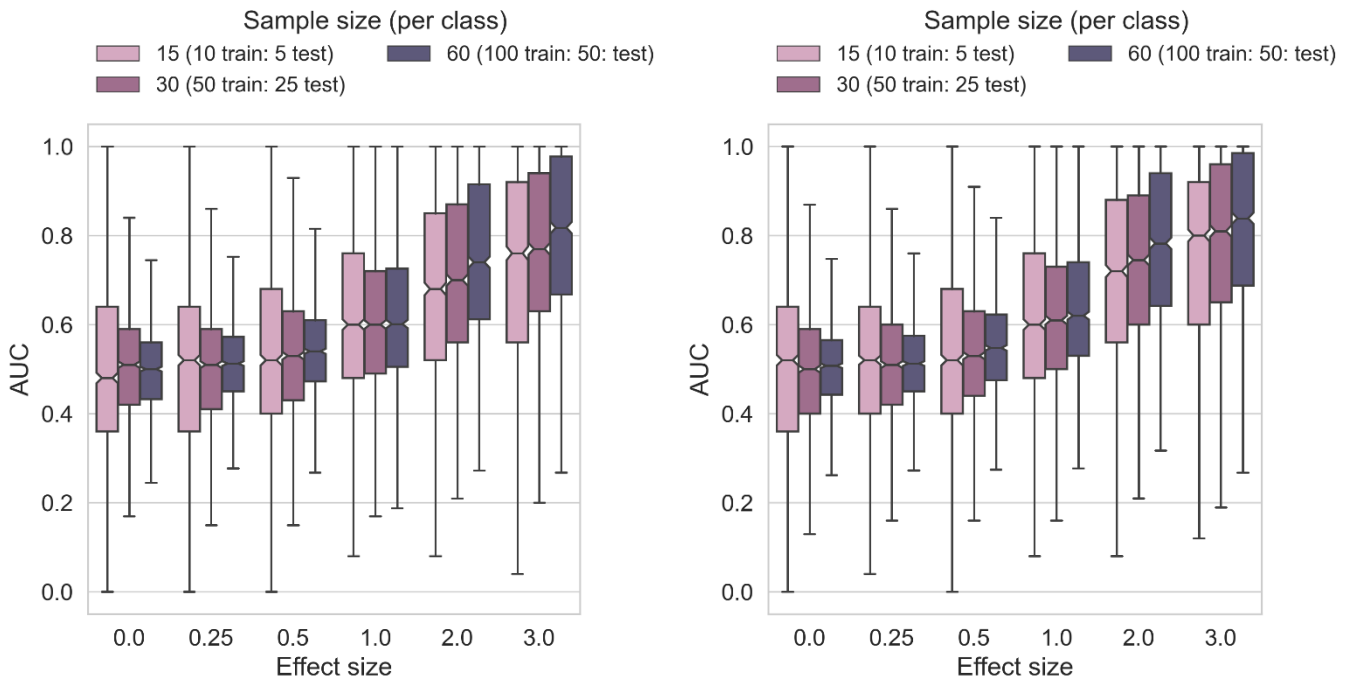


196

197 **Fig I in S1 Supporting Information.** Comparison of PathIntegrate classification performance
198 using KEGG and Reactome pathway databases as well as molecular-level model based on semi-
199 synthetic COVID-19 data.

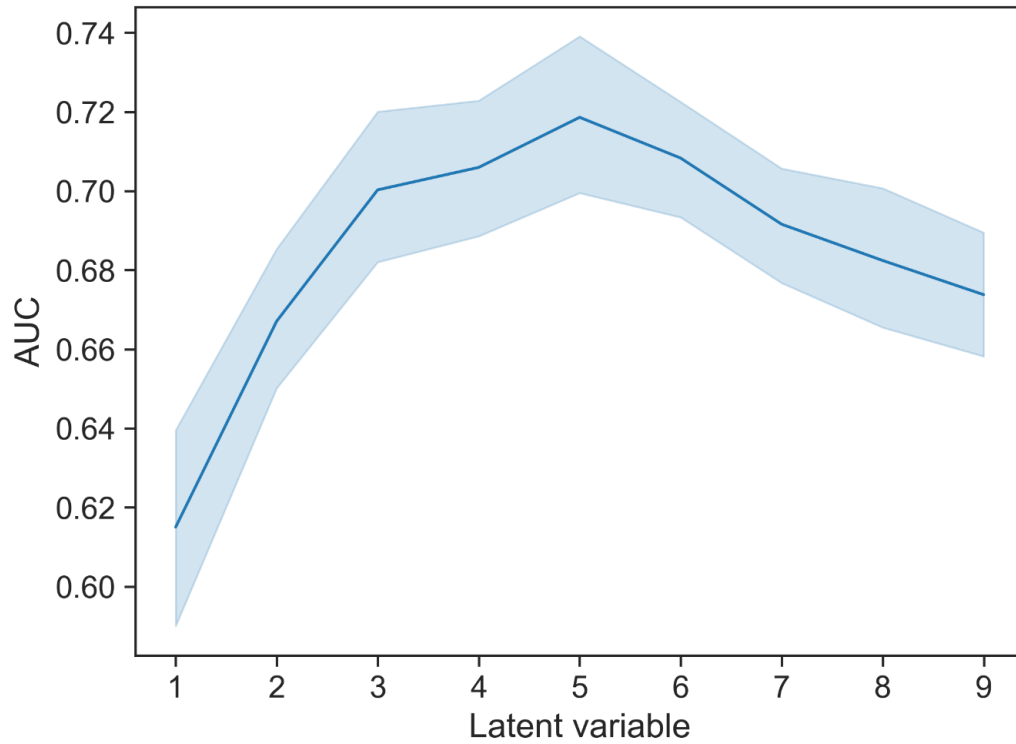
200

201



202 **Fig J in S1 Supporting Information:** Investigation of effects of sample size in PathIntegrate
203 Multi-View (left) and Single-View (PLS) (right) classification performance based on semi-
204 synthetic COVID-19 data.

205



206

207 **Fig K in S1 Supporting Information:** 5-times repeated nested 5-fold cross-validated results
208 for number of latent variables parameter tuning in PathIntegrate Multi-View for COPDgene
209 case study integrating metabolomics, proteomics, and transcriptomics data. X axis shows mean
210 AUC across inner folds. Error bars represent standard deviation.

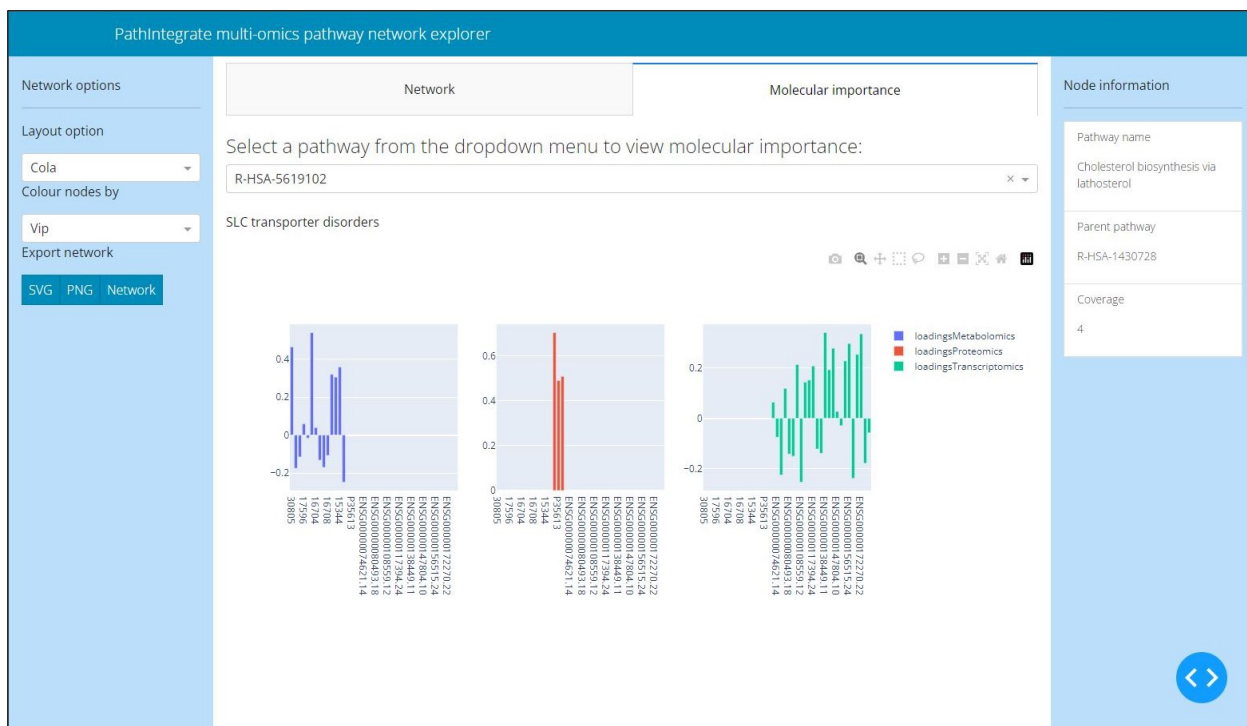
211

212

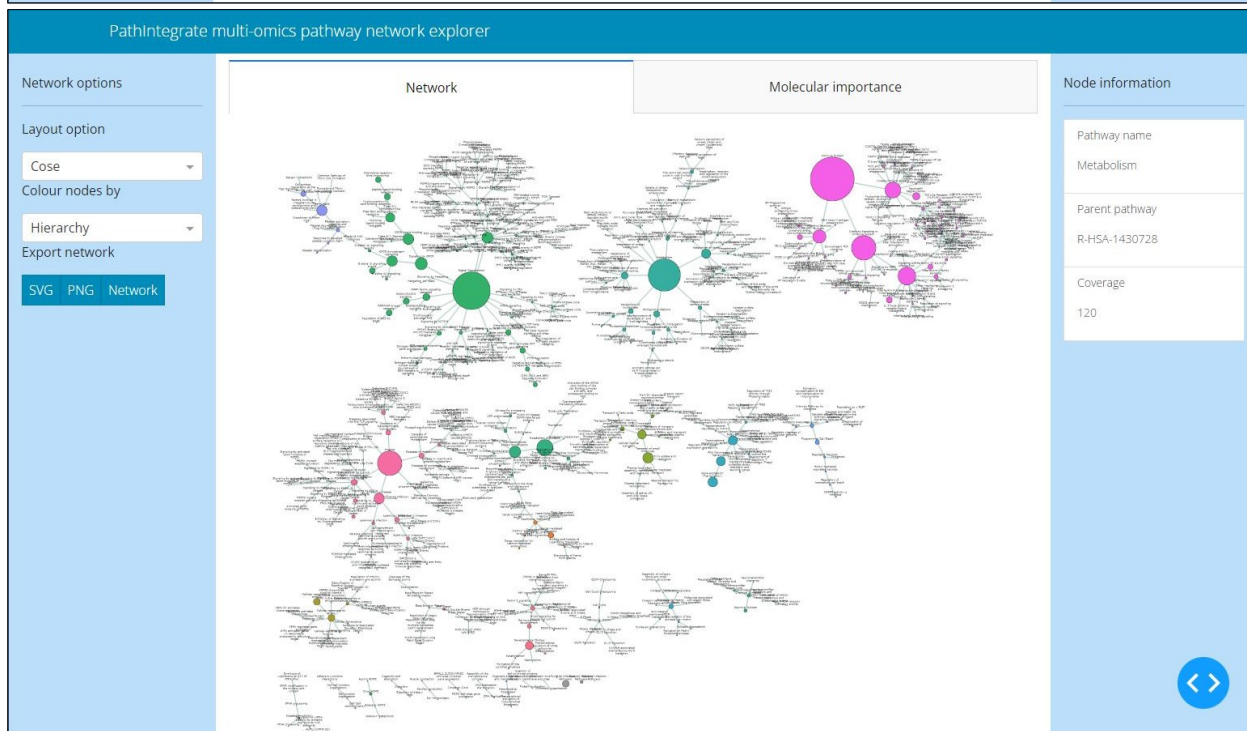
213 **Table B in S1 Supporting Information:** Clinical data definitions for significantly correlated
 214 clinical variables from COPDgene study shown in Fig 4F.

Variable	Definition
AGE_VISIT	Age in years
CAT4_Breathless	Cat questionnaire breathlessness
Finalgoldphase 2	GOLD stage at Phase 2
CurrentMedUse	Currently do you use medications to treat breathing problem
SGRQ score total	St George's Respiratory Questionnaire total score (1-100)
Predicted FEV1_FVC	Predicted ratio of the forced expiratory volume in the first one second to the forced vital capacity of the lungs
FEV1_post	Post-bronchodilator forced expiratory volume in one second
FEV1_FVC post	Post-bronchodilator forced expiratory volume in one second to the forced vital capacity of the lungs
Gender	Gender
Race of subject	Race of subject

215

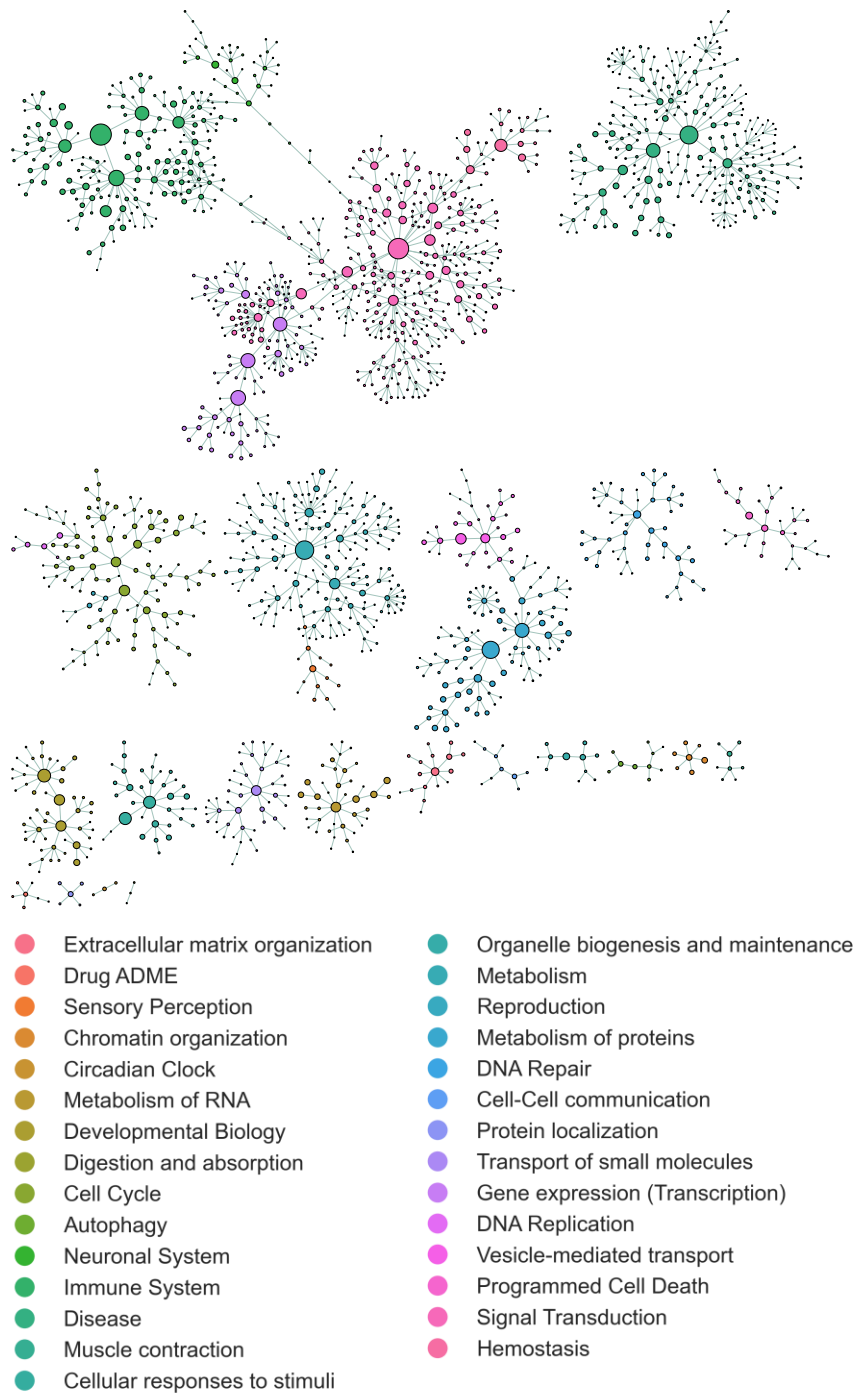


216



217

218 **Fig L in S1 Supporting Information:** Preview of PathIntegrate network explorer app (running
 219 on a local host server) showing an example of a multi-omics dataset being analysed. Interactive
 220 visualisations are facilitated by the open-source Plotly Dash framework (MIT license). Nodes in
 221 the network represent pathways and edges represent parent-child relationships between them.
 222 Users can zoom in and hover over nodes to see more information about the pathway.



224 **Fig M in S1 Supporting Information:** Reactome hierarchy network (based on coverage in
 225 COPDgene multi-omics data) coloured by root pathway membership with full legend. In the
 226 interactive app users can hover over nodes to see detailed information about pathway name,
 227 root pathway, and coverage in a dataset.

228

Variable	Dimension	Definition
X	$[N, M]$	Molecular level matrix of N samples by M molecular features
A	$[N, P]$	Pathway level matrix of N samples by P pathway features
N		Number of samples profiled
M		Number of molecular features profiled
P		Number of pathways accessible in an omics dataset based on minimum coverage threshold
L		Number of molecules present in a given pathway p
p_P		A pathway member of the total pathway set P set consisting of a set of molecules
m_L		A molecule member of the pathway p_P
Z	$[N, L]$	Sub matrix of X containing only the L columns (molecules) present in the i^{th} pathway
Y	$[N, H]$	Outcome variable
H		Number of columns in outcome variable (1 in univariate case)
\hat{Y}	$[N, H]$	Predicted outcome variable
β	$[M, 1]$	Set of regression coefficient of each variable in a regression model
VIP		Variable importance in projection statistic of PLS model
MB-VIP		Multi-block variable importance in projection statistic of MB-PLS model
α		Constant added to semi-synthetic data corresponding to magnitude of enrichment
C		Set of samples present in the semi-synthetic simulated control group
D		Set of samples present in the semi-synthetic simulated case group
θ		Model hyperparameters
T	$[N, R]$	PLS X score matrix
T_s	$[N, R]$	MB-PLS X super score matrix
V	$[M, R]$	PLS X loadings matrix (note usually denoted by P , but here we use P for pathways)
U	$[N, R]$	PLS Y scores matrix

C	$[M, R]$	PLS Y weights matrix
E, F, G	$[N, M]$	PLS model residual matrices
W	$[M, R]$	PLS X weights matrix
W^*	$[M, R]$	
R		Number of latent variables in PLS/MB-PLS model
k		Number of omics data matrices (predictor blocks)
J		Total number of features in a single X block
f		Total number of features across all k predictor blocks

230

231

232 **Reference**

- 233 1. Singh A, Shannon CP, Gautier B, Rohart F, Vacher M, Tebbutt SJ, et al. DIABLO: An
234 integrative approach for identifying key molecular drivers from multi-omics assays.
235 *Bioinformatics*. 2019;35: 3055–3062. doi:10.1093/bioinformatics/bty1054
- 236 2. Meng C, Basunia A, Peters B, Gholami AM, Kuster B, Culhane AC. MOGSA: Integrative
237 single sample gene-set analysis of multiple omics data. *Molecular and Cellular*
238 *Proteomics*. 2019;18: S153–S168. doi:10.1074/mcp.TIR118.001251
- 239 3. Jeon J, Han EY, Jung I. MOPA: An integrative multi-omics pathway analysis method
240 for measuring omics activity. *PLoS One*. 2023;18: e0278272.
241 doi:10.1371/JOURNAL.PONE.0278272
- 242 4. Hänzelmann S, Castelo R, Guinney J. GSVA: Gene set variation analysis for microarray
243 and RNA-Seq data. *BMC Bioinformatics*. 2013;14: 7. doi:10.1186/1471-2105-14-7
- 244 5. Argelaguet R, Velten B, Arnol D, Dietrich S, Zenz T, Marioni JC, et al. Multi-Omics
245 Factor Analysis—a framework for unsupervised integration of multi-omics data sets.
246 *Mol Syst Biol*. 2018;14: e8124. doi:10.15252/msb.20178124
- 247 6. Al-Akwaa FM, Yunits B, Huang S, Alhajaji H, Garmire LX. Lilikoi: an R package for
248 personalized pathway-based classification modeling using metabolomics data.
249 *Gigascience*. 2018;7: 1. doi:10.1093/gigascience/giy136
- 250 7. Drier Y, Sheffer M, Domany E. Pathway-based personalized analysis of cancer. *Proc*
251 *Natl Acad Sci U S A*. 2013;110: 6388–6393. doi:10.1073/pnas.1219651110
- 252 8. Odom GJ, Ban Y, Colaprico A, Liu L, Silva TC, Sun X, et al. PathwayPCA: an
253 R/Bioconductor Package for Pathway Based Integrative Analysis of Multi-Omics
254 Data. *Proteomics*. 2020;20: 1900409. doi:10.1002/PMIC.201900409
- 255 9. Kim SY, Jeong HH, Kim J, Moon JH, Sohn KA. Robust pathway-based multi-omics data
256 integration using directed random walks for survival prediction in multiple cancer
257 studies. *Biol Direct*. 2019;14: 1–13. doi:10.1186/S13062-019-0239-8/FIGURES/5
- 258 10. Kim SY, Choe EK, Shivakumar M, Kim D, Sohn K-A. Multi-layered network-based
259 pathway activity inference using directed random walks: application to predicting
260 clinical outcomes in urologic cancer. *Bioinformatics*. 2021;37: 2405–2413.
261 doi:10.1093/BIOINFORMATICS/BTAB086
- 262 11. Maghsoudi Z, Nguyen H, Tavakkoli A, Nguyen T. A comprehensive survey of the
263 approaches for pathway analysis using multi-omics data integration. *Brief Bioinform*.
264 2022 [cited 18 Oct 2022]. doi:10.1093/BIB/BBAC435

265