

1 PathIntegrate: Multivariate modelling approaches for 2 pathway-based multi-omics data integration

3 Cecilia Wieder¹, Juliette Cooke², Clement Frainay², Nathalie Poupin², Russell Bowler³,
4 Fabien Jourdan⁴, Katerina J. Kechris⁵, Rachel PJ Lai⁶, Timothy Ebbels¹

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

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

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

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

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

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

14

15

Abstract

16 As terabytes of multi-omics data are being generated, there is an ever-increasing need for methods
17 facilitating the integration and interpretation of such data. Current multi-omics integration
18 methods typically output lists, clusters, or subnetworks of molecules related to an outcome. Even
19 with expert domain knowledge, discerning the biological processes involved is a time-consuming
20 activity. Here we propose PathIntegrate, a method for integrating multi-omics datasets based on
21 pathways, designed to exploit knowledge of biological systems and thus provide interpretable
22 models for such studies. PathIntegrate employs single-sample pathway analysis to transform multi-
23 omics datasets from the molecular to the pathway-level, and applies a predictive single-view or
24 multi-view model to integrate the data. Model outputs include multi-omics pathways ranked by
25 their contribution to the outcome prediction, the contribution of each omics layer, and the
26 importance of each molecule in a pathway. Using semi-synthetic data we demonstrate the benefit of
27 grouping molecules into pathways to detect signals in low signal-to-noise scenarios, as well as the
28 ability of PathIntegrate to precisely identify important pathways at low effect sizes. Finally, using
29 COPD and COVID-19 data we showcase how PathIntegrate enables convenient integration and
30 interpretation of complex high-dimensional multi-omics datasets. The PathIntegrate Python
31 package is available at <https://github.com/cwieder/PathIntegrate>.

32

33

Author summary

34 Omics data, which provides a readout of the levels of molecules such as genes, proteins, and
35 metabolites in a sample, is frequently generated to study biological processes and perturbations
36 within an organism. Combining multiple omics data types can provide a more comprehensive
37 understanding of the underlying biology, making it possible to piece together how different
38 molecules interact. There exist many software packages designed to integrate multi-omics data, but
39 interpreting the resulting outputs remains a challenge. Placing molecules into the context of
40 biological pathways enables us to better understand their collective functions and understand how
41 they may contribute to the condition under study. We have developed PathIntegrate, a pathway-
42 based multi-omics integration tool which helps integrate and interpret multi-omics data in a single
43 step using machine learning. By integrating data at the pathway rather than the molecular level, the
44 relationships between molecules in pathways can be strengthened and more readily identified.
45 PathIntegrate is demonstrated on Chronic Obstructive Pulmonary Disease and COVID-19
46 metabolomics, proteomics, and transcriptomics datasets, showcasing its ability to efficiently extract
47 perturbed multi-omics pathways from large-scale datasets.

48

49 Introduction

50 Multi-omics data integration is rapidly becoming a mainstream strategy used to elucidate
51 complex molecular mechanisms in biological systems. Data profiled using diverse
52 modalities, including genomics, epigenomics, transcriptomics, proteomics, and
53 metabolomics provides complementary insights into the regulation of diverse biomolecules
54 and their cellular functions¹. Multi-omics data integration can delineate the transition from
55 genotype to phenotype, while providing a more holistic view of a biological system. Despite
56 the promise that multi-omics integration holds, the field itself is relatively young and faces
57 numerous challenges¹⁻⁶. Among these is the question of which method to use, and how to
58 interpret the results. Several review papers categorise multi-omics integration methods
59 according to underlying concepts, models, or intended purposes⁷. The choice of method used
60 will depend highly on the desired outcome, which can be broadly split into outcome
61 prediction (e.g. sample stratification) or elucidating molecular mechanisms (but often a
62 combination of these). Studies focused on outcome prediction may leverage integration
63 methods based on kernels or deep learning to optimise predictive performance⁸⁻¹⁰, whereas
64 those where the goal is hypothesis generation may opt for more explainable models using
65 classical supervised^{11,12} or unsupervised learning approaches¹²⁻¹⁵, joint pathway analysis¹⁶⁻
66 ¹⁹, network models^{12,20}, or Bayesian statistics⁷. The latter ‘hypothesis generation’-based
67 analysis, regardless of the method used, will often output results in the form of lists of
68 molecules (i.e. genes, proteins, metabolites), typically ranked by their contribution to the
69 model. Depending on the parameters and outputs of the model, the end-user may have
70 multiple latent variables¹³, clusters^{21,22}, or networks²³ composed of many molecules (genes,
71 proteins, and metabolites) to analyse. Doing so is not only be time consuming but requires
72 expert domain knowledge to place biomolecules into a functional context.

73 Pathway analysis (PA) refers to computational methods that have been specifically
74 developed to alleviate the task of analysing long lists of molecules by placing them into a
75 functional context based on curated pathway collections²⁴. Generally, conventional PA
76 methods such as over-representation analysis or gene set enrichment analysis use statistical
77 tests to determine which pathways are associated with a phenotype of interest^{25,26}. The
78 output is typically a list of significantly enriched pathways and their associated test statistics
79 and *p*-values. PA methods are frequently used due to their convenient representation of
80 omics data in the form of pathway descriptors, providing a straightforward interpretation of
81 the biological processes that may contribute to disease phenotypes. Multi-omics pathway
82 analysis is a relatively new but promising area of research²⁷. Tools such as MultiGSEA¹⁹,
83 ActivePathways¹⁷, PaintOmics²⁸, and IMPaLA¹⁶ all leverage multiple layers of biological
84 information to compute enrichment of multi-omics pathways, associated statistical
85 significance levels, and visualisations as an end-result. While highly useful, these methods
86 lack certain desirable features, including the ability to predict outcomes, enabling model
87 performance evaluation, or obtaining a representation of the data in a lower-dimensional
88 space. These goals can be achieved by using pathway-based predictive models, which use
89 pathway rather than molecular-level features to model and predict new data, and infer
90 pathway enrichment through feature importance²⁹⁻³². We provide a detailed overview of
91 related methods in supplementary information, but to the best of our knowledge, we are

92 unaware of any one method which provides predictive, integrative modelling of multi-omics
93 data at the pathway-level.

94 In this work we introduce PathIntegrate, a modelling framework and corresponding Python
95 toolkit to facilitate pathway-based multi-omics integration. PathIntegrate employs single-
96 sample pathway analysis approaches (ssPA), which transform molecular-level abundance
97 data matrices into pathway-level matrices, by using summarisation approaches (e.g.
98 principal component analysis (PCA)) to condense molecular-level measurements into
99 pathway scores for each individual sample in a dataset³³⁻³⁷. By using pathway-transformed
100 multi-omics datasets as input to multivariate supervised models, multi-omics data can be
101 integrated at the pathway-level, providing the user with a range of outputs including i)
102 interpretation of multi-omics pathways associated with the outcome, ii) prediction of
103 outcomes, iii) contribution of each omics view to the model and prediction (in the case of
104 multi-view models), iv) projection of the multi-omics data to a lower dimensional space (in
105 the case of latent variable models). An inherent challenge in multi-omics integration is the
106 heterogeneity between omics datatypes, both in terms of the number of features profiled
107 and the range of numerical values. PathIntegrate addresses these within the pathway-
108 transformation step, where disparate omics datasets are brought to a common scale, i.e. in
109 terms of pathway ‘activity’. Compared to their molecular-level counterparts, pathway-based
110 multi-omics integration models can provide a more parsimonious model when there are
111 fewer input pathways than molecules, while also enabling the detection of multiple small,
112 correlated signals that may not be detected in the molecular-level data. Moreover, pathway-
113 based modelling could increase robustness to data noise by maximising biological variation
114 and simultaneously reducing technical variation²⁹.

115 PathIntegrate consists of two supervised learning frameworks for pathway-based multi-
116 omics integration: PathIntegrate Single-View, which produces a multi-omics pathway-
117 transformed dataset and applies a classification or regression model to the data, and
118 PathIntegrate Multi-View, which uses a multi-block partial least regression (MB-PLS) model
119 to model interactions between pathway-transformed omics datasets. Note that both
120 PathIntegrate Multi-View and Single-View are multi-omics integration methods, and here we
121 use the terms ‘Multi’ and ‘Single’ to refer to the type of predictive model applied (multi-view
122 or single-view³⁸). As both these frameworks rely on pathway transformation (ssPA) of the
123 input omics data, we first demonstrate the ability of univariate methods to detect pathway
124 signals at higher power than molecular-level signals in low signal-to-noise scenarios. We
125 then show that PathIntegrate models can precisely detect enriched pathways even at low
126 effect sizes, as well as use this information to accurately classify samples. PathIntegrate was
127 benchmarked against DIABLO¹¹, a popular multi-omics integration tool with a similar
128 predictive framework, but which does not use pathway transformation. Finally, we showcase
129 the benefits of using PathIntegrate to interpret complex data using case studies on Chronic
130 Obstructive Pulmonary Disease (COPD) and COVID-19 multi-omics datasets, illustrating the
131 ability of the method to identify important and relevant pathway signatures. The
132 PathIntegrate Python package is freely available at
133 <https://github.com/cwieder/PathIntegrate>, and is designed to be compatible with many
134 SciKitLearn³⁹ functions, enabling fast and efficient model optimisation and evaluation.

135 PathIntegrate models are fitted in minutes and can run on a laptop with standard hardware
136 (e.g. 8GB RAM, 1.4 GHz processor).

137 **Results**

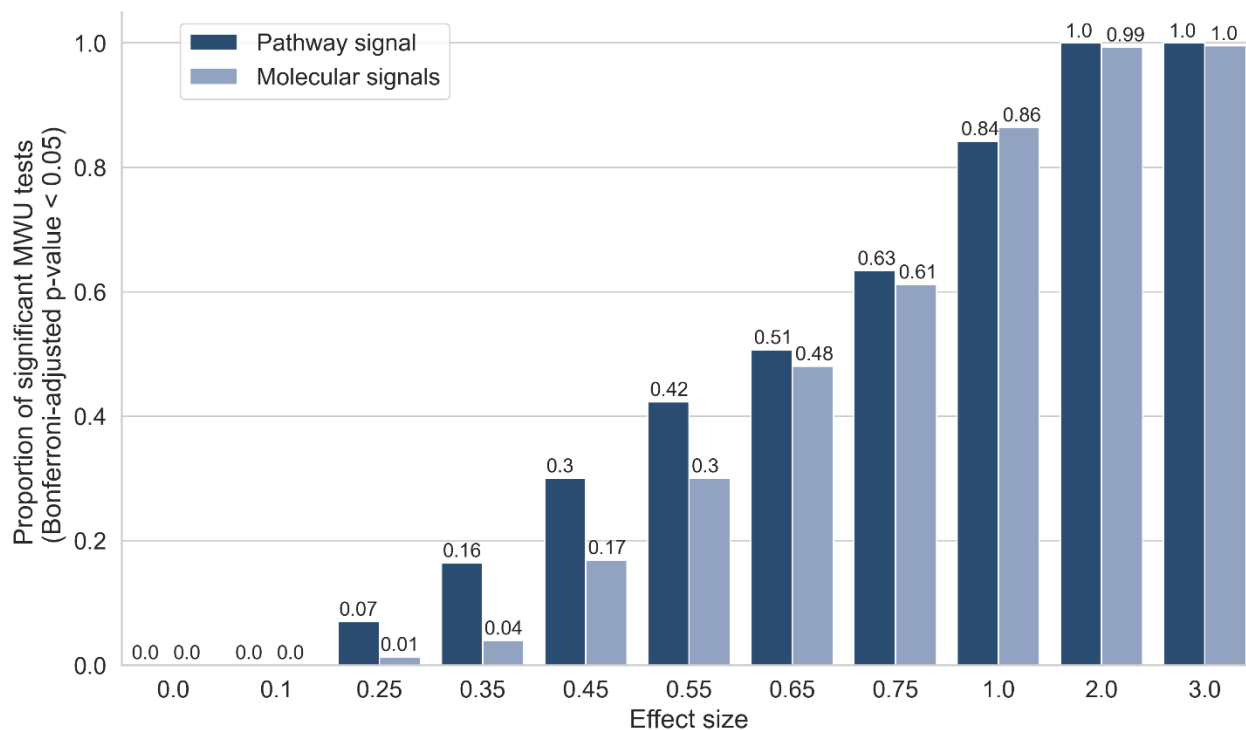
138 **Pathway transformation increases sensitivity to coordinated, low signal-to-noise** 139 **biological signals**

140 Aside from improvements in interpretability, we hypothesized that pathway-based
141 modelling or transformation of data can also provide increased sensitivity in detection of
142 pathway signals in the data, particularly in low signal-to-noise scenarios. By combining
143 abundance levels of correlated individual molecules within a pathway, we anticipate that
144 statistical methods will be able to detect the pathway signal with higher power than
145 individual molecular signals alone. Throughout this work, we refer to ‘molecular-level’
146 models as those with individual molecular entities (such as genes, proteins, and metabolites)
147 as input features, as opposed to ‘pathway-level’ models, which take ssPA pathway-
148 transformed data as input and hence features represent a combination of molecules in each
149 pathway. Briefly, ssPA methods require an $X_{N \times M}$ matrix of molecules as input and combine
150 the abundance values of molecules in a set of predefined pathways to provide an $A_{N \times P}$
151 pathway-level matrix, where features represent pathways and each sample has an ‘activity
152 score’ for each pathway (see Methods).

153 The use of ‘semi-synthetic’ data, in which artificial biological signals are inserted into
154 experimental multi-omics data, provides us with a ground truth we can use to benchmark
155 methods throughout this work³³. We used semi-synthetic multi-omics (metabolomics and
156 proteomics) data derived from COPD and COVID-19 studies (see Methods) to examine
157 whether pathway transformation of multi-omics data allowed pathway signals to be
158 detected by univariate analysis (Mann Whitney-U tests (MWU)) at higher power than
159 individual molecular signals (Fig. 1 and Supplementary Fig. 3). Each omics dataset was
160 transformed to the pathway level using ssPA, using the kPCA ssPA method³³ (see Methods).
161 At each realisation of the simulation, repeated for each Reactome pathway accessible in the
162 datasets, we enriched all the molecules in the pathway (metabolites and/or proteins) in the
163 simulated disease group for a range of effect sizes, corresponding to the range of log₂ fold
164 changes observed in the original datasets (Supplementary Fig. 1, Supplementary Fig. 2).

165 We applied MWU tests to detect differences between the simulated phenotype groups based
166 on the enrichment of each of the individual molecules in the molecular level data or ssPA
167 scores of the target pathway itself. For the molecular level simulation, we applied Fisher’s
168 method to combine p -values in the target pathway if at least 50% constituent molecules were
169 significant ($p \leq 0.05$), otherwise the combined p -value was set to 1. Encouragingly, at lower
170 effect sizes (i.e. 0.25-0.55), we observed a higher proportion of significant p -values in the
171 pathway-transformed data than in the molecular level data. The same trends were observed
172 irrespective of the dataset used to create the simulation (Fig. 1 and Supplementary Fig. 3).
173 This suggests that pathway-transformation approaches could improve the detection of low
174 signal-to-noise, correlated signals in multi-omics datasets, and motivates the use of

175 PathIntegrate models in the remainder of this work, which use ssPA pathway transformation
176 to enable pathway-based multi-omics integration.



177 **Figure 1: Pathway transformation enhances sensitivity to low signal-to-noise signals.** *y*
178 *axis shows proportion of MWU tests significant at Bonferroni $p \leq 0.05$, performed either on*
179 *the pathway-level data or the molecular level data, at varying effect sizes shown on x-axis.*
180 *Semi-synthetic data based on COVID-19 dataset.*

181 PathIntegrate: Supervised pathway-based multi-omics integration frameworks

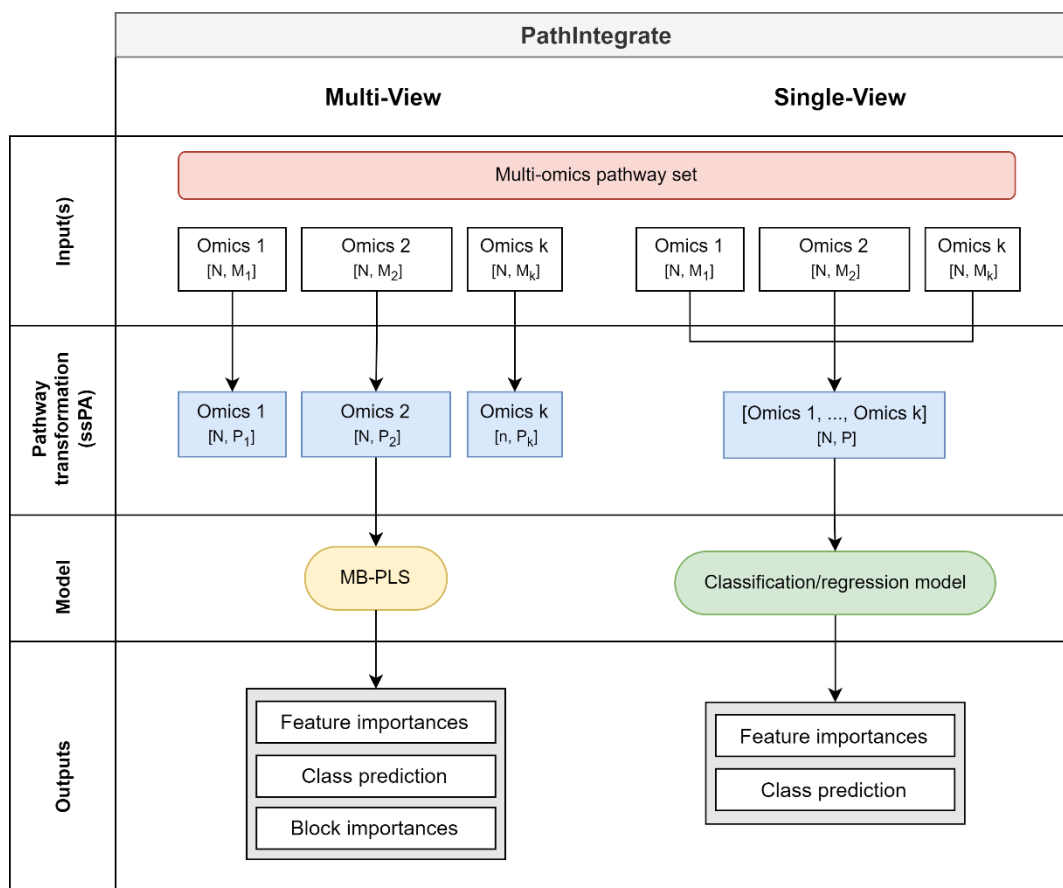
182 In this study we present and investigate the use of the PathIntegrate modelling frameworks
183 for multi-omics pathway-based integration (Fig. 2). PathIntegrate provides two supervised
184 models: Multi-View and Single-View. They are both designed to take two or more (k)
185 $X_{N \times M}$ sample-by-molecule omics abundance matrices as well as a labelled outcome vector y
186 as input and apply a single-sample pathway analysis transformation (facilitated by our
187 recently published ssPA Python package³³) before a predictive model is applied to the data.
188 PathIntegrate can model both continuous and binary outcomes using classification and
189 regression models, but for simplicity we have demonstrated it using binary (e.g. case-
190 control) outcomes throughout this work. Both frameworks achieve the same key outcomes:
191 i) using pathway scores to predict an outcome, and ii) ranking multi-omics pathways by
192 importance in the prediction. PathIntegrate Multi-View uses a multi-table integration model
193 and can therefore provide interpretable insights both within and between omics views,
194 whereas PathIntegrate Single-View provides more flexibility on the high-level predictive
195 model applied and can be better tuned towards prediction. Both models use a single set of
196 multi-omics pathways P , where each pathway has a unique identifier and description, and

197 contains a set of molecular identifiers which can either belong to different omics (i.e.
198 metabolites, proteins, and genes) or in some cases only one omics (i.e. only proteins). Using
199 these pathways, PathIntegrate Multi-View computes pathway scores on each omics view
200 separately, whereas Single-View computes them from multi-omics data.

201 PathIntegrate Multi-View uses a multi-block partial least squares (MB-PLS) latent variable
202 model to integrate ssPA-transformed multi-omics data. Each omics block is transformed to
203 the pathway level individually and the resulting $k A_{N \times P_i}$ blocks are used as input to the MB-
204 PLS model. This preserves the block structure of each omics view and importantly allows
205 users to compute how much each view contributes to the prediction of the outcome variable
206 y , as well as extract within- and between-omics level results such as pathway importances
207 and latent variable representations (scores and superscores⁴⁰⁻⁴²). Importantly, the latent
208 variable model used by Multi-View enables extraction of orthogonal biological effects,
209 similar to PCA, possibly capturing contrasting processes. Furthermore, such models are ideal
210 for pathway-level data, where there is expected to be a high degree of overlap and co-
211 linearity which is accounted for by the PLS framework.

212 PathIntegrate Single-View begins by computing multi-omics pathway scores by performing
213 ssPA transformation on molecular abundance or expression profiles obtained across
214 multiple omics data blocks (e.g. genes, proteins, and metabolites). A single $A_{N \times P}$ pathway-
215 level matrix is returned, in which each feature represents the ‘activity’ of each sample in a
216 multi-omics pathway. The resulting multi-omics pathway scores are used as input to a
217 predictive model (any SciKitLearn compatible model e.g., partial least squares discriminant
218 analysis (PLS-DA), logistic regression, support vector machine, random forest, etc). Pathway
219 importances can be obtained using variable selection approaches appropriate for the model
220 used (e.g., Gini impurity for random forests or the β coefficient for regression-based models).

221 By describing and evaluating the two PathIntegrate modelling frameworks we aim to help
222 users select the method best suited to their study design and research questions.



223

224 **Figure 2: PathIntegrate Multi-View (left) and Single-View (right) modelling**
 225 **frameworks for multi-omics pathway-based integration.** Frameworks are outlined in
 226 terms of their input data, pathway-transformation stage, statistical model, and outputs. Blue
 227 data blocks represent omics data which has been transformed from the molecular ($X_{N \times M}$)
 228 space to the pathway ($A_{N \times P}$) space using ssPA. Both Single-View and Multi-View make use of
 229 the same multi-omics pathway set.

230 PathIntegrate performance evaluation

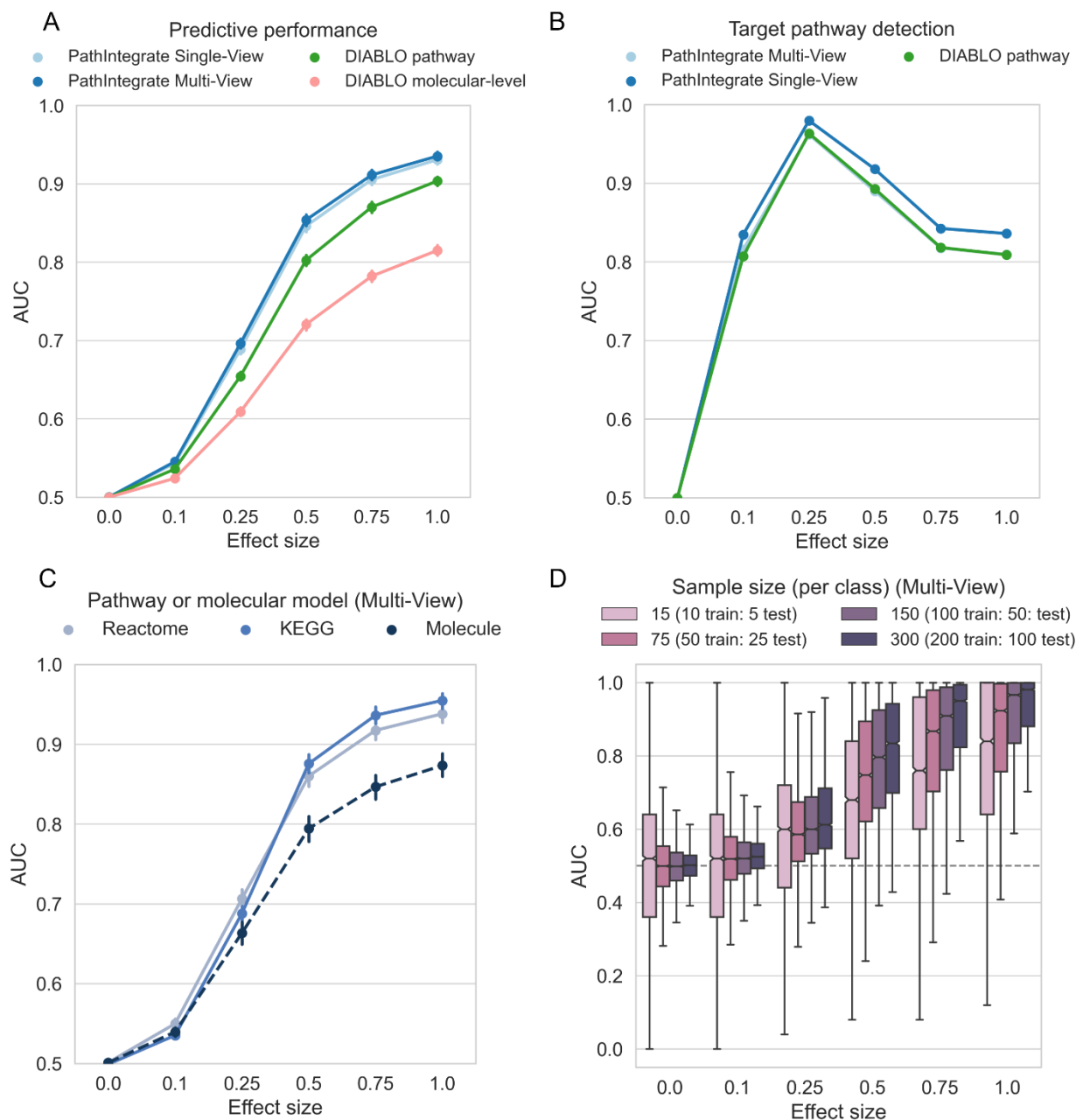
231 PathIntegrate Multi-View and Single-View were evaluated in a classification setting by a) the
 232 ability to discriminate between sample classes based on important pathways, and b) the
 233 ability to rank important pathways highly. Using semi-synthetic simulated metabolomics
 234 and proteomics data (see Methods) we enriched one target Reactome pathway containing
 235 metabolites and/or proteins at a time, at varying effect sizes, and repeated this for each
 236 pathway accessible in the datasets. For simplicity and consistency between datasets we
 237 integrated two omics throughout the performance evaluation section. Results based on
 238 COPDgene semi-synthetic data are shown in Fig. 3, and results based on COVID-19 semi-
 239 synthetic data are shown in Supplementary Fig. 8. Note that this simulation design is rather
 240 conservative, because only one pathway is enriched in each realisation (although its
 241 constituent molecules may overlap with other pathways), whereas in a real biological system
 242 we may expect multiple pathways to be enriched at once. PathIntegrate Multi-View used

243 multi-block PLS as the underlying predictive model, and for purposes of comparison,
244 PathIntegrate Single-View used standard PLS-DA.

245 We compared PathIntegrate to the state-of-the-art multi-omics integration method DIABLO
246 from MixOmics^{11,43}. To the best of our knowledge, DIABLO is the most similar multi-omics
247 integration method developed to date which makes use of a multi-view framework. As
248 DIABLO is flexible as to the input data matrices, we compared standard DIABLO (using
249 molecular-level omics data, 'DIABLO molecular-level'), as well as a pathway-based DIABLO
250 ('DIABLO pathway') using the same ssPA-transformed omics matrices as input to
251 PathIntegrate Multi-View. Importantly, although we are comparing the performance of
252 PathIntegrate to DIABLO, we do not expect significant increases in predictivity or ability to
253 detect the target pathway, due to the similarity of the underlying generalised canonical
254 correlation analysis model to MB-PLS. Instead, we aim to highlight the flexibility of using
255 pathway scores as input to supervised integrative models, such as DIABLO, and that even
256 using different multivariate algorithms can yield predictive models capable of identifying
257 target pathways with high sensitivity and specificity, and thus generating more interpretable
258 results.

259

260



261

262 **Figure 3: Performance of PathIntegrate and DIABLO vs. effect size, based on semi-**
 263 **synthetic data measured by AUROC.** COPDgene metabolomics and proteomics data were
 264 integrated in each model. a. Ability to correctly predict sample outcomes (case vs. control).
 265 We compared PathIntegrate Multi-View and Single-View to DIABLO using both molecular and
 266 pathway-level multi-omics data. b. Ability to correctly recall target enriched pathway. We
 267 compared DIABLO RGCCA model loadings to the Multi-View MB-PLS VIP and Single-View PLS
 268 VIP statistics for pathway importance. c. Comparison of PathIntegrate Multi-View
 269 classification performance using KEGG and Reactome pathway databases as well
 270 molecular-level model. d. Effect of sample size on PathIntegrate Multi-View classification

271 *performance. For panels a-c error bars indicate 95% confidence intervals on the mean AUROC*
272 *(in some cases they appear smaller than point sizes).*

273 A fundamental question is whether modelling data using pathways can yield improvements
274 in predictive performance compared to using molecular level data. Fig. 3a shows the ability
275 of PathIntegrate Multi-View, PathIntegrate Single-View, and DIABLO to predict samples in
276 an unseen test set based on AUROC (Fig. 3a, Supplementary Fig. 8a). All methods began to
277 discriminate sample classes even at low effect sizes (0.1 - 0.25), concordant with findings
278 from the univariate simulation. The pathway-based models (PathIntegrate Multi-View,
279 Single-View and 'DIABLO pathway') exhibited improved performance compared to the
280 'DIABLO molecular-level' (standard) model across all effect sizes. As effect size increased
281 from 0.25-1.0 the PathIntegrate methods performed similarly to 'DIABLO pathway'. Overall,
282 these results suggest that using pathway-level models may yield improved predictive
283 performance compared to molecular-level models.

284 We also compared the predictive performance of PathIntegrate models using pathways from
285 two different databases, Reactome and KEGG, as well as the performance of MB-PLS and PLS
286 models using the molecular-level data (i.e. PathIntegrate without the pathway-
287 transformation step) (Fig. 3c/Supplementary Fig. 8c shows PathIntegrate Multi-View and
288 Supplementary Fig. 5/Supplementary Fig. 9 show PathIntegrate Single-View and DIABLO).
289 Results for the molecular level simulation can vary depending on the number of molecules
290 enriched at each realisation, which correspond to the size of the pathway in the equivalent
291 pathway-level simulation. Because Reactome and KEGG have differing distributions of
292 pathway sizes⁴⁴, we randomly sampled the number of molecules enriched in each realisation
293 based on the combined distribution of Reactome and KEGG pathway sizes, in order to reduce
294 dependence on database pathway size. At lower effect sizes (0.1 - 0.25), both the molecular
295 and pathway-level models performed similarly, whereas at moderate-to-high effects the
296 pathway-based models exhibited an increase in predictive performance concordant with
297 trends observed in Fig. 3a. Models based on KEGG pathways appear to perform marginally
298 better than Reactome pathways at larger effect sizes, which may be due to KEGG pathways
299 being larger on average (see Supplementary Fig. 4 and Supplementary Table 1 for pathway
300 database size statistics).

301 We next evaluated the ability of PathIntegrate and 'DIABLO pathway' to accurately detect
302 the target enriched pathway. (Fig. 3b, Supplementary Fig. 8b). For PathIntegrate Single-View
303 and Multi-View methods, variable importance in projection (VIP and multi-block-VIP) were
304 used to evaluate feature importances⁴¹. p -values for the significance of each pathway feature
305 VIP or MB-VIP value were computed empirically based on 10,000 sample permutations with
306 BH-FDR correction. For 'DIABLO pathway', the RGCCA loadings on component 1 were used
307 to infer feature importance, and p -values were subsequently computed using the same
308 permutation testing approach. A true positive enriched pathway was defined as being the
309 target enriched pathway and having an adjusted p -value of ≤ 0.05 (see Methods for full
310 description of the confusion matrix computation). Both PathIntegrate and DIABLO models
311 performed well in terms of target pathway detection, even being able to detect the target
312 pathway with high AUC (≥ 0.90) at low effect and high noise scenarios (effect size = 0.25).
313 PathIntegrate Multi-View performed almost identically to 'DIABLO pathway'. All methods
314 experience a decrease in AUC at higher effect sizes (0.5-1), which is expected due to

315 pathways overlapping with the target pathway reaching significance, and in-built
316 normalisation of the model weights/loadings causing the magnitude of the coefficient of the
317 target pathway to shrink slightly in comparison to those of highly overlapping pathways. For
318 simplicity, these overlapping pathways are treated as false positives, though they contain
319 truly differentially abundant molecules. Thus, this decrease does not point to a lower
320 performance of the method in identifying pathways relevant to prediction of the outcome.
321 Furthermore, while the primary emphasis of this work is not on contrasting regularized and
322 non-regularized models, it is worth noting that sparse models are widely used for feature
323 selection. We also compared the ability of the models to select the target pathway with a
324 sparse version of DIABLO (using the L1 norm, see Methods) (Supplementary Fig. 6,
325 Supplementary Fig. 8b). At low to moderate effect sizes, the sparse model identified the
326 target pathway at similar AUC to the PathIntegrate/non-regularised DIABLO model, but at
327 high effect sizes it showed slight improvements in target pathway identification as the
328 sparsity constraint prevented high numbers of overlapping pathways reaching significance.

329 Finally, we investigated the effect of sample size, which is well known to influence model
330 performance, on PathIntegrate models. We down-sampled each of the two classes in the
331 data, keeping a 1:1 ratio between classes, and evaluated the predictive ability of the models
332 at varying effect sizes (Fig. 3d/Supplementary Fig. 8d and Supplementary Fig. 7/
333 Supplementary Fig. 10 show results for Multi-View and Single-View respectively). As
334 expected, the lower the number of samples in the model, the more variability observed in
335 the predictions. Particularly at lower effect sizes, smaller sample numbers were more likely
336 to result in false positives and spurious results. While it is not possible to state the minimum
337 number of samples necessary to apply PathIntegrate models, it is important for users to test
338 the performance of the model using appropriate cross-validation approaches to be confident
339 that the conclusions are statistically robust.

340 While these results demonstrate the predictive ability of PathIntegrate models, it is
341 challenging to create a realistic simulation scenario which accurately reflects molecular
342 activities and their participation in pathways in a biological system. Hence, we have applied
343 PathIntegrate to the COPDgene and COVID-19 experimental datasets in the Application
344 section to further illustrate model performance and interpretation.

345

346 **PathIntegrate Multi-View applied to COPDgene data**

347 The COPDgene cohort consists of 10,198 smokers at baseline with and without chronic
348 obstructive pulmonary disease (COPD) ⁴⁵. We integrated metabolomics, proteomics, and
349 transcriptomics multi-omics data measured at Phase 2 (~5 years after baseline) profiled on
350 a subset of individuals with all three omics data (n=522) using PathIntegrate to identify
351 Reactome pathways associated with COPD pathology. The Multi-View model of
352 PathIntegrate allows users to gain rich insights into the underlying data, from high-level
353 interpretation of the global rankings of enriched pathways, to being able to investigate the
354 importance of pathways in each omics block and latent component individually. We applied
355 the kPCA ssPA method to produce pathway score matrices for each omics view and using 5-
356 fold cross validation, we found that four latent variables yielded an optimised MB-PLS model

357 (mean cross-validated AUC: 0.70) (Supplementary Fig. 11). The MBPLS superscores for each
358 of the four latent variables coloured by COPD status are shown in Fig. 4a, providing a visual
359 representation of the ability of multi-omics pathways to identify differences between COPD
360 and non-COPD groups, in which each of the four latent variables exhibit a visible difference
361 between groups.

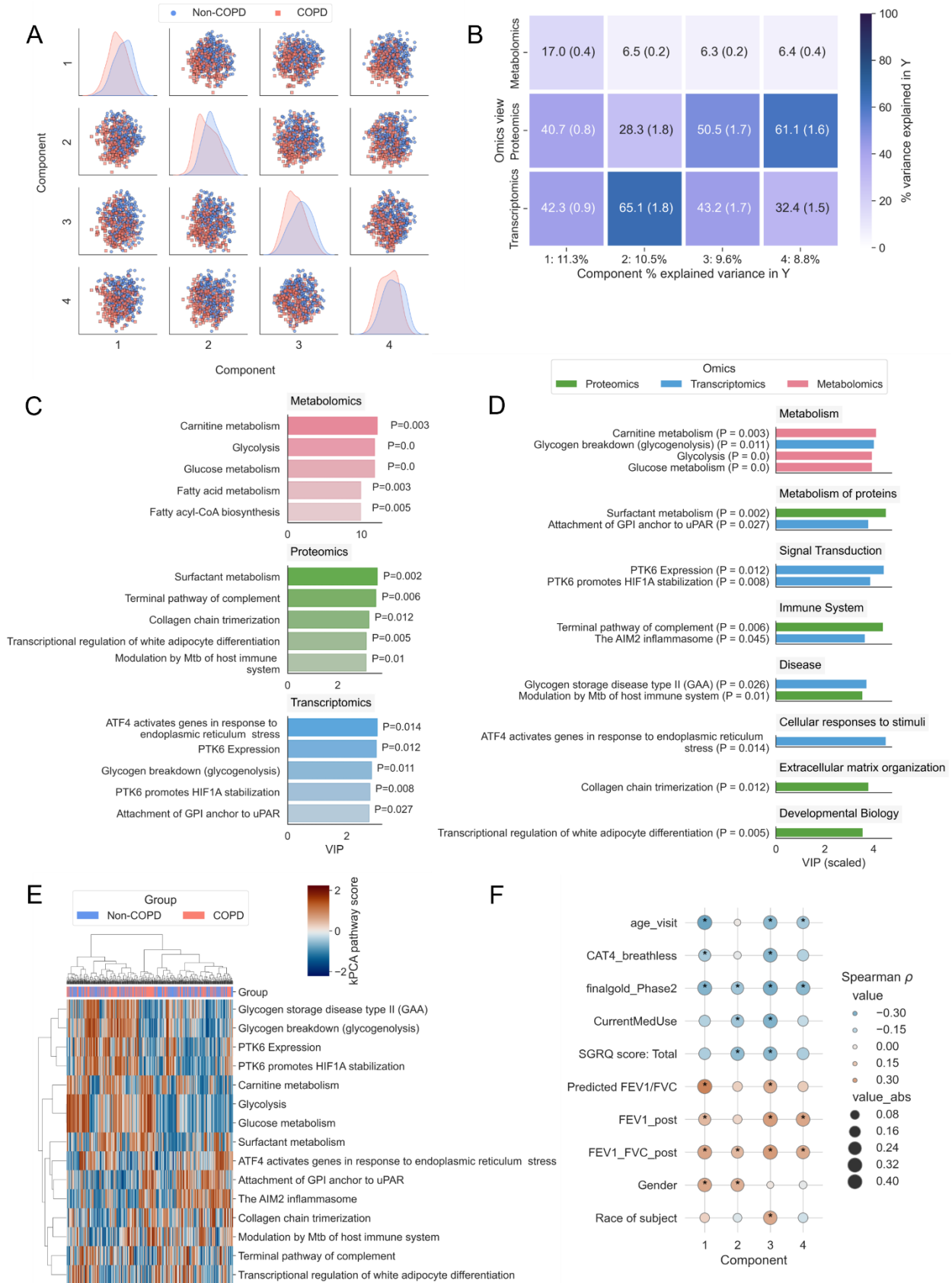
362 One of the primary insights obtained from the Multi-View model is the contribution of each
363 omics view to the variance explained in the outcome variable y (Fig. 4b). In the first latent
364 variable, all three omics accounted for a considerable proportion of the variance explained
365 in y , suggesting the pathway scores correlate well in the latent space. In the further three
366 latent variables, transcriptomics and proteomics views tend to contribute most to the
367 outcome prediction. Although metabolomics describes less of the variance in y than the
368 other omics, based on 100 bootstrap samples the mean variance explained across all latent
369 variables remained between 6 and 17 percent. The dominance of transcriptomics and
370 proteomics views may suggest that the COPD vs non-COPD distinction is best captured by
371 gene and protein-level signalling pathways as opposed to metabolic pathways, but it may
372 also be due to the lower metabolite coverage, and smaller set of pathways accessible using
373 these molecules (Table 3, Supplementary Table 1).

374 We then investigated the pathways ranked highly by MB-VIP across all latent variables.
375 Pathway importances can be queried at an individual omics level (Fig. 4c), or at a multi-omics
376 level with VIP normalised across all views (Fig. 4d). The same is also possible at the
377 individual latent variable level, and as superscores are orthogonal, each latent variable
378 contains a different combination of pathways contributing to the prediction of y . p -values for
379 the MB-VIP statistic were computed empirically using permutation testing (see Methods). In
380 Fig. 4c we observe that the metabolic pathways implicated in COPD pathology relate broadly
381 to fatty acid metabolism, including carnitine metabolism, as well as central carbon
382 metabolism⁴⁶. The transcriptomics layer also highlighted the importance of glycogenolysis
383 (glycogen breakdown), which alongside alterations in lipid metabolism have been found to
384 be implicated in severe COPD, where there is an increased dependence on glucose for energy
385 production due to impaired lipolysis, and hence an increased rate of glycolysis⁴⁷. Carnitine
386 metabolism was one of the top ranked (metabolic) pathways overall, with Fig. 4d showing
387 its significance was driven by the Metabolomics layer ($p=0.003$). The ‘Carnitine metabolism’
388 pathway is composed of both metabolites and proteins, of which there was also sufficient
389 coverage in the transcriptomics data to produce ssPA scores for this pathway. In the
390 transcriptomics data however, this pathway was not significant ($p=0.55$); this demonstrates
391 the benefit of multi-omics modelling to gain a broader perspective of the molecular basis of
392 disease. Fatty acid metabolism has been shown to be part of a metabolic reprogramming that
393 occurs in respiratory disease including COPD^{48,49}. In COPD specifically, impairments in the
394 carnitine shuttle system in the mitochondria (preventing long-chain fatty acids from being
395 transported into the mitochondria) have been shown to result in lipotoxicity within the cell
396 cytosol^{50–52}. Conversely, ‘Surfactant metabolism’, which did not have sufficient coverage to
397 be included in the model in the metabolomics view, but was found relevant in the proteomics
398 data ($p=0.002$), is an important process by which phospholipid surfactants are produced by
399 the alveoli to ensure optimal lung function⁵³. The surfactant lipidome has been found to be
400 significantly different in COPD patients compared to healthy controls and is a potential

401 therapeutic target⁵³. Finally, several relevant proteomics and transcriptomics pathways
402 involved focus on innate immune processes, highly important in the chronic inflammatory
403 nature of COPD, such as inflammasome action ('The AIM2 Inflammasome') and the
404 complement system ('Terminal Pathway of Complement'). The AIM2 inflammasome has
405 recently been implicated in COPD pathogenesis, correlating with COPD severity and cigarette
406 smoke exposure⁵⁴. The full list of significant pathways is available in Supplementary File 1.

407 To demonstrate alternative visualisation strategies possible with PathIntegrate, we
408 extracted the top 15 pathways across all omics ranked by MB-VIP from the Multi-View model
409 and used the ssPA scores for these pathways to cluster the samples (Fig. 4e). Hierarchical
410 clustering showed two distinct clusters of pathways, one relating to metabolic processes
411 such as central carbon and fatty acid metabolism, as well as hypoxia-associated signalling
412 pathways ('PTK6 expression', 'PTK6 promotes HIF1A stabilization'), and the other consisting
413 of processes involved in the innate immune response ('The AIM2 inflammasome', 'Terminal
414 pathway of complement').

415 Further interpretation of the model can be gained by examining the correlation between the
416 superscores for each latent variable and clinical metadata, enabling investigation of the
417 relationship between clinical features and pathways (Fig. 4f). For example, we found
418 pathways in latent variables 1, 3, and 4 to be significantly associated with age, whereas
419 pathways in latent variable 3 were significantly associated with the race of subjects.



421 **Figure 4: PathIntegrate Multi-View applied to COPDgene multi-omics data.** A.
 422 *Superscores plot based on multi-omics (metabolomics, proteomics, and transcriptomics*
 423 *pathways) across four latent variables. B. Omics view importances across latent variables.*
 424 *Values represent mean and SEM across 100 bootstrap samples. C. Top five pathways per*
 425 *omics block. D. Top 15 pathways across omics blocks categorised by Reactome parent*
 426 *pathway. E. kPCA ssPA scores from top 15 pathways used to cluster samples using Euclidean*
 427 *distance and Ward linkage. F. Heatmap showing Spearman correlation between superscores*
 428 *across four latent variables and clinical metadata. Asterisks indicate Bonferroni p-value \leq*
 429 *0.05. Definitions of clinical variables are in Supplementary Table 2.*

430
 431 Finally, to check that the pathway-based modelling approach does not appreciably degrade
 432 prediction performance, we examined the performance of PathIntegrate Multi-View versus
 433 a molecular-level MB-PLS model using the COPDgene dataset (Table 1). In the case of
 434 predicting COPD using plasma multi-omics data (metabolomics, proteomics, and
 435 transcriptomics), for example, the pathway level model achieved an average AUC of 0.70
 436 (± 0.02), and the molecular level model also achieved an average AUC of 0.70 (± 0.02) when
 437 using all molecules available (inc. those not mapping to pathways), but required more latent
 438 variables to do so (4 vs. 6), resulting in a more complex model (Table 1).

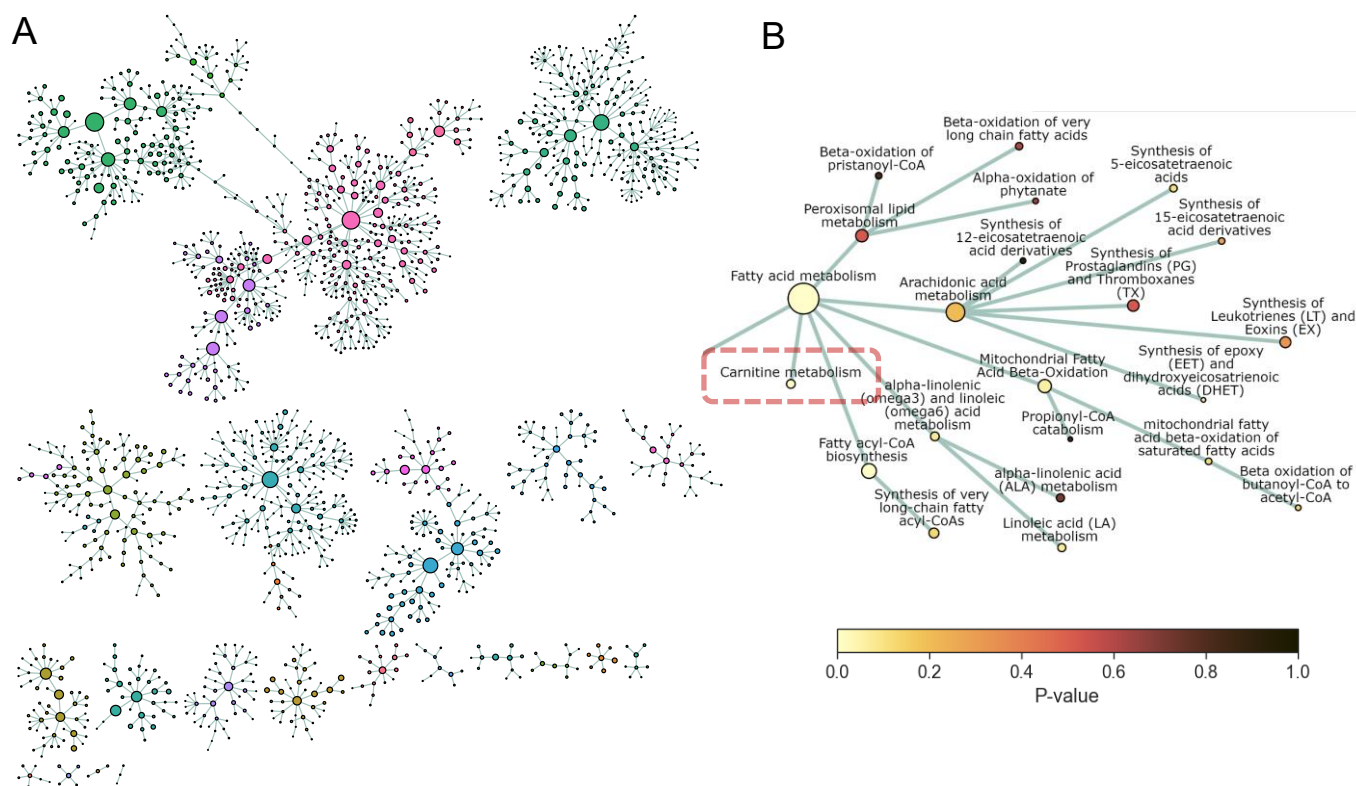
439 **Table 1: Performance comparison of PathIntegrate Multi-View using pathways versus**
 440 **using the molecular-level COPDgene dataset (mean AUC and 95% CI, as well as the number**
 441 **of latent variables (LV) used).** In both pathway and molecular-level scenarios the model was
 442 used to predict binary COPD status. The molecular-level model was fit both with all molecules
 443 available in the datasets, as well as only those mapping to pathways. AUC values are averaged
 444 across 5-times repeated 5-fold cross validation.

	All omics	Metabolomics and proteomics	Metabolomics and transcriptomics	Transcriptomics and proteomics	Metabolomics	Proteomics	Transcriptomics
AUC (pathway)	0.70 (0.67, 0.72) (4 LV)	0.67 (0.66, 0.69) (3 LV)	0.69 (0.67, 0.71) (3 LV)	0.68 (0.66, 0.70) (4 LV)	0.63 (0.61, 0.64) (1 LV)	0.67 (0.66, 0.68) (3 LV)	0.65 (0.63, 0.66) (3 LV)
AUC (molecular)	0.70 (0.69, 0.72) (6 LV)	0.71 (0.70, 0.72) (2 LV)	0.70 (0.68, 0.71) (6 LV)	0.71 (0.70, 0.73) (7 LV)	0.66 (0.65, 0.69) (2 LV)	0.72 (0.71, 0.74) (3 LV)	0.68 (0.66, 0.69) (5 LV)

AUC								
(molecular – only those mapping to pathways)	0.72 (0.70, 0.74) (7 LV)	0.72 (0.70, 0.74) (2 LV)	0.67 (0.66, 0.69) (6 LV)	0.70 (0.69, 0.72) (6 LV)	0.68 (0.67, 0.7) (2 LV)	0.71 (0.70, 0.73) (3 LV)	0.66 (0.64, 0.68) (7 LV)	

445

446 Visualisation of high-dimensional omics data in the context of many hundreds of pathways
447 remains a challenge. Alongside typical graphical outputs from the model, the PathIntegrate
448 package provides an interactive network explorer app designed to visualise the results of
449 PathIntegrate models on the Reactome pathway hierarchy graph (Supplementary Fig. 12).
450 Nodes in the network represent pathways and edges represent parent-child relationships
451 between them as part of a directed acyclic graph (DAG). Nodes can be coloured by feature
452 importance in the PathIntegrate model, so that users can intuitively visualise important
453 pathways and their relationships to other areas of the pathway network. Various
454 hierarchical and force-directed layouts are available, and images can be exported for further
455 annotation and customisation. Fig. 5a shows a global overview of the Reactome pathway
456 network based on coverage of the COPDgene dataset (full pathway hierarchy legend shown
457 in Supplementary Fig. 13). We coloured nodes by MB-VIP p -values in Fig. 5b to identify
458 important pathways linked to COPD, as well as other pathways which may be affected by
459 proximity in the network. Fig. 5b highlights the ‘Carnitine metabolism’ pathway ($p \leq 0.05$), as
460 well as other pathways which may not have reached statistical significance but may be of
461 interest such as ‘Arachidonic acid metabolism, or ‘Mitochondrial fatty acid beta
462 oxidation’^{48,55}. Encouragingly, related pathways in the close neighbourhood of ‘Carnitine
463 metabolism’ have lower p -values than those further from it.



464
 465 **Figure 5: Network visualisation with PathIntegrate interactive network explorer.**
 466 PathIntegrate Multi-View was applied to COPDgene multi-omics data. A. Multi-omics network
 467 view of global Reactome hierarchy DAG. Only pathways with sufficient coverage (≥ 2
 468 molecules per pathway) are shown as nodes. Edges represent parent-child relationships
 469 between pathways as defined by Reactome. Nodes are coloured by Reactome superpathway
 470 membership. Node size corresponds to pathway coverage. B. Network view of 'Carnitine
 471 metabolism' pathway (zoomed-in subset of (a)) and close neighbourhood within the
 472 Reactome pathway hierarchy. Nodes are coloured by p-values obtained from PathIntegrate
 473 Multi-View model.

474 Taken together, these results demonstrate how PathIntegrate Multi-View can be used to
 475 investigate various aspects of pathway regulation associated with a specific phenotype.
 476 COPD-associated pathways can be explored both within omics (individual views) and across
 477 omics (global view), and superscores of the latent variables can be used to identify
 478 correlations between pathways and other data, e.g. clinical measurements. The contribution
 479 of each omics to the prediction can be easily obtained from the Multi-View model, which
 480 obtains a lower-dimensional representation of the data that maximises covariances between
 481 omics view blocks and the y outcome, but also keeps data blocks separate in order to retain
 482 this level of granularity.

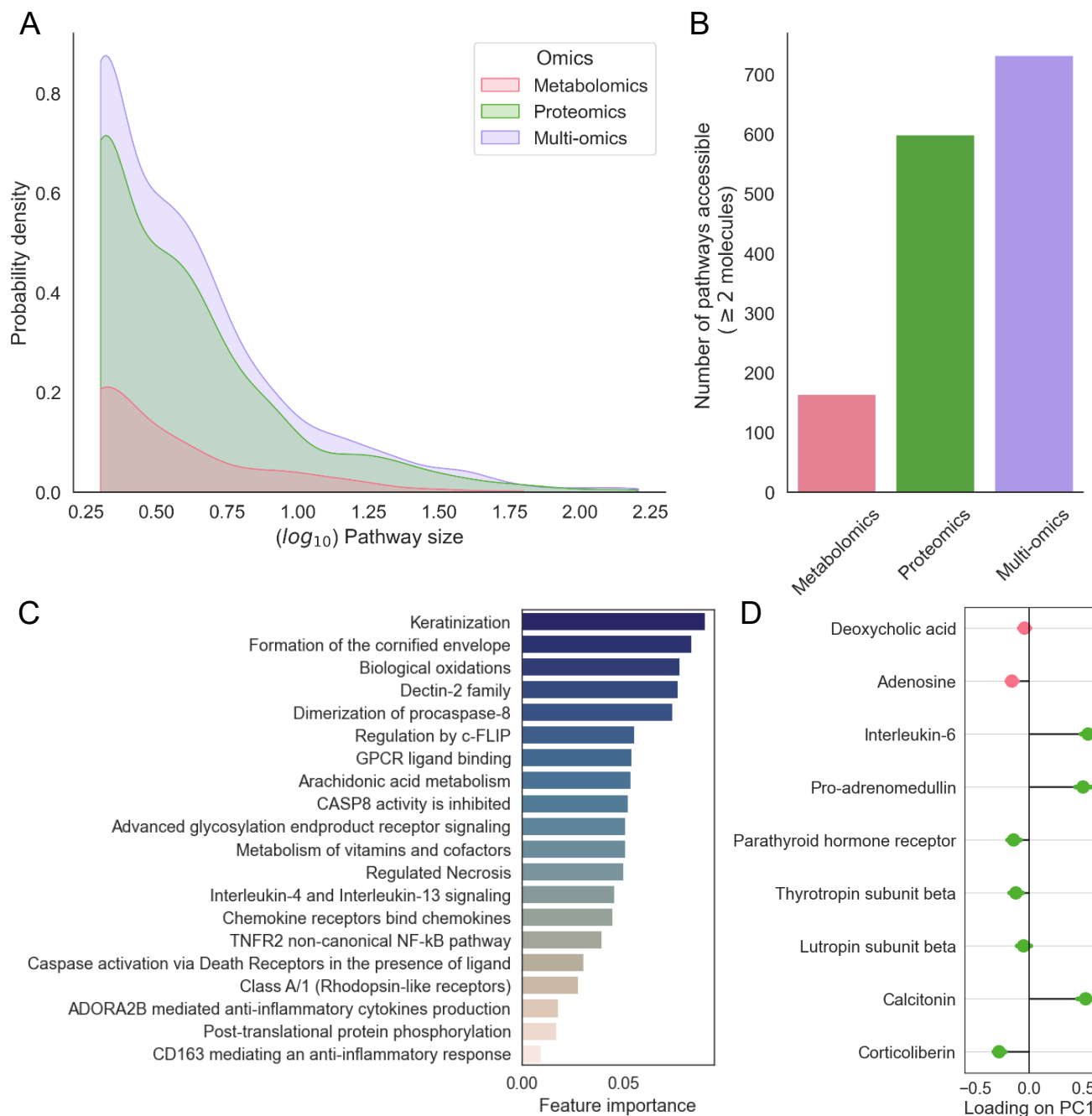
483 PathIntegrate Single-View applied to COVID-19 multi-omics data

484 We applied PathIntegrate Single-View to data from a multi-omics study of COVID-19 severity
485 ⁵⁶ to understand pathways driving the transition from mild to moderate/severe COVID-19
486 pathogenesis. Proteomics and LC-MS metabolomics data were integrated using
487 PathIntegrate Single-View, in which the concatenated omics data were transformed to multi-
488 omics ssPA scores using the SVD method ³⁷ and a random forest model was applied to the
489 resulting Reactome pathway score matrix.

490 An advantage of the Single-View model is that it computes each pathway score based on
491 multi-omics data, providing a broader coverage of pathways by doing so (Fig. 6a). Multi-
492 omics pathways had a greater mean pathway coverage (number of molecules in the data
493 mapping to each pathway, mean: 6.39 versus 6.21 and 4.86 for proteomics and metabolomics
494 separately). This enabled more pathways to be included as they contained enough molecules
495 to meet the minimum filtering threshold (732 pathways versus a maximum of 599 and 169
496 for proteomics and metabolomics separately, a total of 701 unique pathways); we used a
497 liberal threshold of ≥ 2 molecules per pathway) (Fig. 6b).

498 We found the PathIntegrate Single-View model to perform similarly in terms of classification
499 AUC on the unseen test set (AUC 0.95) compared to the concatenated molecular level omics
500 data (AUC: 0.98), suggesting that in this case pathway-level modelling can aid interpretation
501 without substantial loss of prediction performance. We next inspected the important multi-
502 omics pathway features using random forest recursive feature elimination, which identified
503 20 of the most informative pathways (Fig. 6c). Within this set, there are several immune-
504 related processes known to be implicated in COVID-19 severity such as 'Interleukin-5 and
505 interleukin-13 signalling'^{57,58} and 'Caspase activation via death receptors in the presence of
506 ligand'⁵⁸.

507 Finally, if certain ssPA methods are used (e.g. SVD³⁷), it is possible to obtain information on
508 how individual molecules contribute to the formation of the overall multi-omics pathway
509 score. As we used SVD scores in this model, we can use the loadings on principal component
510 1 as the importance of each molecule in the pathway score (Fig. 6d). In Fig. 6d, which shows
511 the molecular-level importance for the 'ADORA2B mediated anti-inflammatory cytokines
512 production' pathway as an example, we observe that metabolites deoxycholic acid and
513 adenosine are correlated with four proteins, all with negative loadings on PC1, while three
514 proteins: interleukin-6 (IL-6) and the hormones pro-adrenomedullin and calcitonin had
515 positive loadings with greater magnitudes. In chronic COVID-19, elevated levels of both IL-6
516 and adenosine have been observed, with IL-6 contributing to the proinflammatory 'cytokine
517 storm' and adenosine being considered as a potential therapeutic for severe cases due to its
518 anti-inflammatory effects⁵⁹. Such investigations can help researchers pinpoint the specific
519 molecules contributing most to pathway scores, reducing the number of molecules required
520 in developing biomarker assays, as well as providing understanding of how molecules from
521 different omics correlate in the latent space.



522

523 **Figure 6: PathIntegrate Single-View applied to COVID-19 multi-omics data. A. Kernel**
 524 **density distribution of \log_{10} pathway sizes in the COVID dataset per omics view. Pathway size**
 525 **refers to the number of molecules annotated to each pathway present in the COVID datasets.**
 526 **B. Number of pathways with sufficient coverage in the COVID dataset in each omics view. C.**
 527 **Multi-omics pathway features identified using recursive feature elimination from the**
 528 **PathIntegrate Single-View random forest model, ranked by Gini importance. D. Molecular**
 529 **level importances derived from the 'ADORA2B mediated anti-inflammatory cytokines**

530 *production' (R-HSA-9660821) SVD pathway scores. Datapoints represent mean and standard*
531 *deviation of loadings of each molecule on PC1 across 200 bootstrap samples.*

532 Discussion

533 This study contributes a new approach to the rapidly growing body of multi-omics
534 integration methods^{3,5,6,27}, specifically by providing insights into the use of pathways as a
535 basis for interpretable predictive modelling of multi-omics data, and by introducing the
536 PathIntegrate framework for doing so. The use of pathways for modelling omics data is a
537 promising avenue of research, with several studies highlighting its potential in recent years
538 ^{27,60}. However, there is limited research available on the use of pathways for multi-omics
539 integration, or evaluation of the performance of pathway-based versus molecular-level
540 integration models. Here, we have introduced the PathIntegrate Multi-View and Single-View
541 modelling frameworks for multi-omics pathway-based integration and evaluated their
542 performance using semi-synthetic and experimental data.

543 To demonstrate the ability of pathway transformation to increase statistical power by
544 combining correlated molecular signals we applied a series of univariate tests to evaluate
545 the ability to detect pathway or molecular level enrichment across various effect sizes. At
546 lower effect sizes, we found that the univariate tests could recover more pathway-level
547 signals than molecular signals, demonstrating the benefit of pathway transformation of
548 multi-omics data, which often have low effect sizes, particularly in heterogeneous clinical
549 studies, and especially those where a phenotype is not well defined. Additionally, pathway
550 transformation naturally reduces the number of tests required, thereby reducing the
551 multiple testing correction burden. This motivated our development of PathIntegrate, which
552 uses ssPA pathway transformation as a basis for pathway-based multi-omics integration.

553 We compared PathIntegrate to DIABLO¹¹, a highly-cited multi-omics integration tool, which
554 uses a similar underlying multi-view model as PathIntegrate Multi-View. We found
555 PathIntegrate methods to perform similarly to DIABLO (when using pathway score matrices
556 as input). Overall, however, we wish to emphasise the benefit of using pathway-transformed
557 data as input to multivariate models and show that even using a different predictive model
558 (DIABLO RGCCA vs PathIntegrate Multi-View MB-PLS) similar results can be obtained. We
559 compared PathIntegrate Multi-View to a molecular level MB-PLS model and demonstrated
560 the ability of the pathway-based model to classify samples with improved AUC across effect
561 sizes. A full comparison of PathIntegrate, a pathway-based predictive model, to conventional
562 pathway analysis approaches, such as ORA²⁵, GSEA²⁶, or integrated pathway analysis e.g.
563 MultiGSEA¹⁹ and IMPaLA¹⁶ is beyond the scope of the present work. This is because pathway-
564 based predictive models leverage multivariate modelling to identify pathways most
565 associated with an outcome, whereas conventional pathway analysis methods typically test
566 pathways in a univariate and non-predictive manner. Although the question of 'which
567 pathways are perturbed in a phenotype?' is similar in both approaches, the way results are
568 derived and the differences in outputs would render a direct comparison challenging, yet an
569 interesting avenue for future research.

570 We applied PathIntegrate to two datasets: COPDgene and COVID19 multi-omics. Both case
571 studies highlighted the benefits of pathway-based modelling for integration, interpretation,
572 and visualisation of multi-omics data. In terms of predictive performance, in both case
573 studies, as expected, PathIntegrate performed similarly to the molecular level counterpart.
574 Pathway coverage, the proportion of molecules in a pathway which can be observed in the
575 data, or pathway annotation, the proportion of known documented biomolecules annotated
576 to pathways in databases are both inherent bottlenecks of pathway-based analyses. These
577 issues particularly affect certain datatypes such as metabolomics, where even multiple
578 assays are not enough to provide high coverage of the metabolic pathway network⁴⁴. Despite
579 this, in the COVID-19 case study where 314 metabolites were annotated to ChEBI identifiers,
580 and 456 proteins to UniProt identifiers, the PathIntegrate Single-View model based on 732
581 multi-omics pathway scores was still able to achieve an AUC of 0.95 in predicting COVID-19
582 severity based on the pathway coverage provided by these molecules. Another important
583 consideration is pathway database choice, as pathway definitions can differ greatly between
584 databases, as well as the level of overlap between pathways and possible hierarchical
585 structure^{44,61-63}. As expected, we found PathIntegrate to exhibit minor changes in predictive
586 performance based on the database used.

587 Although PathIntegrate Multi-View uses an MB-PLS model and Single-View uses any
588 SciKitLearn-compatible predictive model (e.g., random forest), we endeavour to provide
589 readers with a general framework for pathway-based multi-omics integration which they
590 can build upon to complement their experimental design or analysis goals. For example, if
591 prediction of a phenotype with high accuracy is a desired outcome, a deep feed forward
592 neural network could be applied within the Single-View framework, to classify samples
593 based on pathways. Model interpretability can also be further enhanced by customising the
594 model inputs, such as using bespoke pathway sets or ontologies to generate the pathway
595 score input layer. For example, in PathIntegrate Multi-View, an additional omics block could
596 be added composed of lipidomics data, and pathway scores could be computed using the
597 LipidMaps⁶⁴ classification system to reflect enrichment patterns of lipid subclasses. Note
598 that in this work we focused on supervised pathway-based integration models; however
599 similar frameworks using unsupervised methods is also feasible and may be explored
600 further. We decided to focus on supervised methods as firstly an outcome is directly
601 modelled and there is less risk of confounding variation obscuring the interpretation, and
602 secondly, users can evaluate model performance in a straightforward manner by examining
603 prediction accuracy.

604 Both PathIntegrate Single-View and Multi-View are designed to handle multiple omics views.
605 In this work we have demonstrated the use of two or three omics views, however both
606 models can accommodate further (3+) omics views as long as they contain continuous
607 measurements (rather than binary e.g., genomics data) and the features can be mapped to
608 pathway identifiers, enabling the pathway-transformation stage to be performed. Data
609 blocks from the same omics type e.g. metabolomics but profiled on different biofluids or
610 tissues can also be integrated using PathIntegrate, to understand how pathways in different
611 biological matrices contribute to the phenotype. Although the focus of this work was on
612 pathway-based models, both PathIntegrate models can be made hybrid in the sense that both

613 pathway-transformed omics data and other data e.g., clinical metadata, genomics data,
614 metagenomic data, etc., can be integrated alongside one another.

615 PathIntegrate is unique in its specific support for metabolomics in multi-omics studies,
616 which is often omitted by other integration methods. Metabolomics is becoming frequently
617 profiled alongside gene-based omics, providing researchers with an essential snapshot on
618 the biochemical activities of small-molecules^{1,65}. Metabolomics data differs considerably
619 from gene-based omics in several ways including the molecular identifiers used, assay
620 coverage of the metabolome, and annotation uncertainties. PathIntegrate users can
621 download the latest release of Reactome pathways via the `sspa` Python package and obtain
622 a merged multi-omics pathway database object composed of protein (UniProt), gene
623 (ENSEMBL), and metabolite identifiers (ChEBI) to enable integration of these distinct omics
624 in a straightforward manner.

625 Our study shares several limitations with other pathway-analysis and multi-omics
626 integration studies, a key drawback being the lack of appropriate benchmarking data.
627 Ideally, a benchmarking dataset would contain two or more high-quality omics views, a large
628 sample size ($n \geq 1000$), and known biological signals at the molecular and pathway level
629 validated by laboratory experiments. Without access to such data, we employed the semi-
630 synthetic simulation strategy to artificially introduce known molecular and pathway-level
631 signals into a real experimental dataset. As described in our previous work³³, this approach
632 allows the simulation to retain important characteristics of real data such as the underlying
633 statistical distributions, correlations, and covariances between molecules and pathways. It
634 also enabled us to vary the effect size of pathway signals, which we based on the effect sizes
635 (\log_2 fold changes) detected in the experimental datasets used. Despite these efforts, it
636 remains a challenge to compare molecular vs. pathway-level models, as it is unknown how
637 many molecules in a pathway are differentially abundant at any one time, and pathway
638 definitions and sizes vary between databases^{44,63,66,67}.

639 In common with many other statistical integration approaches, PathIntegrate requires all
640 input omics to be measured on the same individuals. This means samples from individuals
641 without data on all omics will have to be discarded, as PathIntegrate currently does not
642 support entire rows of missing data. Some models such as MB-PLS can handle sparse data
643 (using NIPALS algorithm)⁶⁸, however future work is required to determine how robust this
644 could be for high rates of missingness. Further work is needed to develop multi-omics
645 integration methods that can handle missing samples⁶⁹ and using pathway-based models
646 may aid in the robust imputation of data, by helping to capture biological rather than
647 technical variation.

648 **Conclusion**

649 As knowledge of biological pathways continues to evolve and pathway databases develop
650 alongside this, we anticipate that pathway-based models such as PathIntegrate will become
651 a valuable way of interpreting complex multi-omics datasets. This work contributes to our
652 understanding of such models, by evaluating the effectiveness of using pathways for multi-
653 omics integration, as well as introducing the PathIntegrate modelling framework.
654 PathIntegrate provides a novel solution to the challenge of integrating heterogeneous omics

655 datasets, by using pathway-transformation to bring omics to a common basis, followed by
656 state-of-the-art supervised modelling. The PathIntegrate framework presented here and
657 accompanying Python package will provide a useful resource to the research community,
658 streamlining the analysis of multi-omics data with the aim of providing an interpretable,
659 integrated set of results at the pathway level.

660

661 **Methods**

662 **Datasets**

663 **COPDgene data**

664 We integrated COPDgene Phase 2 (~5 years after baseline) plasma metabolomics
665 (Metabolon UHPLC-MS/MS), plasma proteomics (SOMAscan 1.3k assay), and bulk whole
666 blood transcriptomics data (Illumina HiSeq2000) from 522 samples which had data for all
667 three omics. As detailed in Regan et al., 2010⁴⁵: COPD was defined using spirometric
668 evidence of airflow obstruction [post-bronchodilator forced expiratory volume at one
669 second (FEV1)/forced vital capacity (FVC) ≥ 0.70], as well as a GOLD score of 1-4. The sub-
670 cohort comprised 273 COPD samples (GOLD 1-4) and 249 non-COPD samples (GOLD 0) from
671 smokers. Full details of the multi-omics datasets and pre-processing are available in the
672 original article²¹. We also obtained clinical data for samples, including COPD phenotypes and
673 demographic variables. Clinical data was filtered to include 260 variables measured in all
674 522 samples of the sub-cohort.

675 **COVID-19 data**

676 The publicly available COVID-19 multi-omics dataset was obtained from Su et al. 2020⁵⁶. Full
677 details of the multi-omics datasets and pre-processing are available in the original article⁵⁶.
678 We integrated plasma metabolomics (Metabolon UHPLC-MS/MS) and proteomics (Olink)
679 datasets with matched samples, of which 45 samples had 'mild' COVID (WHO status 1-2), and
680 82 had 'moderate-severe' COVID19 (WHO status 3-7), totalling 127 samples.

681 **Multi-omics data pre-processing and quality control**

682 All multi-omics datasets were subject to quality control and pre-processing as detailed in the
683 original articles^{45,56}. Metabolomics, proteomics, and transcriptomics abundances were \log_2
684 transformed followed by unit-variance scaling. Missing values were imputed using the
685 singular-value decomposition approach implemented in the fancyimpute Python package.
686 In the transcriptomics data, low-variance genes (below 25th percentile) were filtered out.
687 Table 2 shows the number of molecules in each omics remaining after identifier mapping
688 and quality control which were used in all analyses.

689 *Identifier mapping*

690 Identifier harmonisation of both the COPDgene and COVID metabolite datasets was
691 performed via the sspa package identifier conversion utility via the MetoboAnalyst⁷⁰ API
692 (<https://www.metaboanalyst.ca/docs/APIs.xhtml>.) HMDB metabolite identifiers provided
693 with the dataset were converted to ChEBI (for Reactome)/KEGG compound (for KEGG)
694 identifiers.

695 COPDgene and COVID-19 proteomic data was provided with UniProt identifiers which
696 directly map to Reactome pathways. KEGG gene IDs were obtained using the UniProt ID
697 matching tool (<https://www.uniprot.org/id-mapping>). COPDgene transcriptomics data was
698 provided with ENSEMBL IDs which directly map to Reactome pathways.

699

700 *Table 2: Number of molecules in each omics in COPDgene and COVID-19 datasets after*
701 *processing and identifier mapping.*

Dataset	Total number of samples	Number of metabolite features (mapping to ChEBI)	Number of protein features (mapping to UniProt)	Number of transcript features (mapping to ENSEMBL)
COPDgene	522	513	1305	14441
COVID-19	127	314	456	NA

702 Pathways

703 PathIntegrate Single-View and Multi-View models both make use of a single, merged set of
704 multi-omics pathways as input. Each pathway contains either a set of molecules from
705 different omics (metabolites (ChEBI), proteins (Uniprot), and genes (ENSEMBL)), or only
706 molecules from a single omics, depending on the pathway definition. The PathIntegrate
707 package enables download of multi-omics pathway sets (via ssPA) from Reactome, providing
708 a text file of the latest version for various supported organisms in standard GMT file format.
709 PathIntegrate is also flexible to the input pathway set and is not restricted to those provided
710 via the package. Any pathway set in GMT file format can be used as input, where each row
711 represents a pathway, and each pathway set is described by a name, a description, and its
712 constituent molecules (see Broad Institute website for further details on GMT format:
713 [#https://software.broadinstitute.org/cancer/software/gsea/wiki/index.php/Data_formats](https://software.broadinstitute.org/cancer/software/gsea/wiki/index.php/Data_formats)
714 [#GMT:_Gene_Matrix_Transposed_file_format_28.2A.gmt.29](#)).

715 In this work, Reactome human version 83 and KEGG human version 105 were used. Table 3
716 shows the number of pathways from each omics in the COPDgene/COVID-19 datasets
717 accessible using the molecules profiled in each dataset (≥ 2 per pathway).

718 *Table 3: Number of Reactome/KEGG pathways accessible in COPDgene and COVID-19 multi-*
719 *omics datasets*

Dataset	Number of Reactome pathways accessible (≥ 2 molecules mapping)	Number of KEGG pathways accessible (≥ 2 molecules mapping)
COPDgene metabolomics	202	125
COPDgene proteomics	1396	291
COPDgene transcriptomics	1902	341
COVID-19 metabolomics	169	122
COVID-19 proteomics	599	217

720 Semi-synthetic multi-omics data generation

721 To benchmark our methods, we applied the semi-synthetic simulation approach detailed in
722 Wieder et al., 2022³³ to insert artificial biological signals into existing multi-omics data. This
723 approach involves creating simulated datasets based on experimental data, with the
724 assumption that doing so will preserve the complex biological signals and statistical
725 distributions within the data, and more accurately reflect a real scenario as opposed to
726 approaches based on sampling from parametric distributions. Various experimental designs
727 can be simulated using this approach, but here we opt for a simple case-control design in
728 which we add the artificial signal only to molecules in the ‘case’ group. By adding the same
729 effect size to the abundances of all molecules within a pathway (detailed below), this
730 approach emphasises realism (by preserving the covariance structure of the original omics
731 data) without being overly complex.

732 The input data is a series of \log_2 transformed abundance matrices for the k omics types $X_k =$
733 $[x_1, x_2, \dots, x_{M_k}]$, each of size $(N \times M_k)$, and a set of N outcome labels $y_i, i = 1, \dots, N$. The
734 approach is as follows for each realisation of the semi-synthetic data:

- 735 1. Randomly shuffle outcome labels y_i . This results in a new ‘control’ group C and a new
736 ‘case’ group D of the same class sizes as the original dataset. The shuffling ensures any
737 biological effects correlated to the outcome are removed but preserves existing
738 covariances between molecules.
- 739 2. Add a constant α corresponding to desired effect size (e.g. $\log_2 \text{FC}=0.5$) to specified
740 target molecules only in samples in the new ‘case’ group $i \in D$, simulating increased
741 abundance of those molecules associated with the outcome (Equation 1).

$$742 \quad X_{i,j} \rightarrow X_{i,j}, \quad i \in C$$

$$743 \quad X_{i,j} \rightarrow X_{i,j} + \alpha, \quad i \in D$$

744 *Equation 1*

745 In this work we increase the abundance of all molecules in a single target pathway at
746 each realisation, at the same effect size. By adding a constant to \log_2 scale data this
747 simulates a multiplicative fold change in the original data.

748 In this work, we enriched all molecules in one randomly selected (Reactome/KEGG) ‘target’
749 pathway p_i at a time, at varying effect sizes. Here, effect size refers to the \log_2 fold change of
750 a molecule. We enriched the known target pathway by effect sizes of 0-1 in the COPDgene
751 dataset and 0-3 in the COVID-19 dataset, based on fold changes observed in the original data
752 (Fig S1, S2). For performance evaluation purposes, we performed the semi-synthetic
753 simulation approach using COPDgene and COVID-19 metabolomics and proteomics datasets.
754 We performed the semi-synthetic simulation once for each target pathway in the
755 Reactome/KEGG database that contained at least 3 molecules mapping to the input data
756 (1290 and 298 realisations for Reactome and KEGG respectively for COPDgene data; 456 and
757 256 for COVID-19 data). For each target pathway we used a different random shuffling of
758 outcome labels.

759 **Single-sample pathway analysis**

760 Reactome human pathways (R83) and KEGG human pathways (R105) were downloaded
761 using the `sspa` Python package v0.2.4 (<https://github.com/cwieder/py-ssPA>). The `sspa`
762 package creates multi-omics pathways by merging proteins/genes and metabolites
763 participating in the same pathway into a single multi-omics pathway.

764 Single-sample pathway analysis (ssPA) is an unsupervised method used to transform omics
765 data matrices into pathway score matrices, where columns represent pathways rather than
766 individual molecules. Importantly, all omics data input to ssPA must be standardised.
767 Throughout this work and in the `sspa` Python package, unit variance scaling is used, where
768 the mean of each feature is set to 0 and the standard deviation is set to 1. ssPA begins by
769 using the P pathways $P = \{p_1, p_2, \dots, p_P\}$ passing minimum coverage criteria for the dataset
770 (an integer defined by the user, default 2 molecules per pathway). The i 'th pathway p_i is
771 composed of L_i molecules (e.g. proteins), $p_i = \{m_1, m_2, \dots, m_{L_i}\}$. ssPA is performed to provide
772 pathway 'activity scores' for each sample, reflecting an estimate of the enrichment of each
773 pathway in each individual sample.

774 One of the most popular categories of ssPA methods is that based on dimensionality
775 reduction, specifically PCA. In the original PLAGÉ (referred throughout this work as 'SVD')
776 method by Tomfohr et al.³⁷, singular value decomposition is performed on the omics
777 abundance matrix retaining only the L_i columns (molecules) present in the i 'th pathway. For
778 each pathway, column vectors of abundance profiles belonging to molecules in pathway p_i
779 are concatenated to form a matrix Z_i (Equation 2).

$$780 \quad Z_i = [x_{m_1}, x_{m_2}, \dots, x_{L_i}]$$

781 *Equation 2*

782 Then, the first right singular vector (first principal component score) is used to represent the
783 pathway 'activity' scores a_i (size $N \times 1$) for the i 'th pathway. Pathway score vectors for each
784 pathway are combined to produce a sample-by-pathway matrix $A = [a_1, a_2, \dots, a_P]$. The
785 kPCA method we proposed in³³ uses a very similar approach, instead applying kernel PCA
786 with a radial basis function kernel and using the scores for principal component 1 to reflect
787 pathway activities. Full details of how ssPA is performed are available in^{33,36,37}. In this work
788 we used the kPCA method³³ in the benchmarking section and COPDgene application, and the
789 SVD method (PLAGE)³⁷ in the COVID-19 application section. The `sspa` package functions
790 `sspa_kPCA` and `sspa_SVD` were used to generate pathway score matrices used in both
791 PathIntegrate Multi-View and Single-View.

792 **Supervised modelling frameworks**

793 **PathIntegrate Single-View**

794 PathIntegrate Single-View is a predictive model applied to a single data matrix of multi-
795 omics ssPA scores (Fig. 2). Conceptually it is simpler than PathIntegrate Multi-View due to
796 the input being a single pathway-level matrix rather than multiple pathway-level matrices.

797 Note that both models integrate multi-omics data; the “Single-View” and “Multi-View” refer
798 to the machine learning framework used to effect this integration.

799 The first step of PathIntegrate Single-View involves computing ssPA scores at the multi-
800 omics level, using multi-omics pathway sets (i.e. pathways $p_i = \{m_1, m_2, \dots, m_{M_i}\}$ where the
801 m_i represent genes, metabolites, and proteins present in the omics data). All input omics data
802 matrices are unit-variance scaled. ssPA is performed on the multi-omics abundance matrices
803 Z_i using any sspa algorithm implemented in the sspa Python package to form the pathway
804 scores matrix A of size $(N \times P)$.³³

805 The second step of PathIntegrate Single-View applies a predictive model to the multi-omics
806 ssPA score matrix A to predict an outcome variable \hat{y} (Equation 3).

$$807 \quad \hat{y} = f(A; \theta)$$

808 *Equation 3*

809 where θ represents the parameters of the predictive model f . There is a single predictor
810 matrix, hence the term ‘Single-View’. The user can apply a variety of models (any of those
811 available in SciKitLearn are compatible with the PathIntegrate python package), including
812 random forest, PLS regression, support vector machine, etc. Important pathways are
813 determined using feature importance metrics specific to the predictive model used (e.g. Gini
814 impurity for random forests or VIP for PLS regression). In this work to demonstrate
815 PathIntegrate Single-View, we applied a PLS-DA model in the performance evaluation
816 section, and a Random Forest model in the COVID-19 case study.

817 **PathIntegrate Multi-View**

818 PathIntegrate Multi-View leverages multi-table integration approaches to build a predictive
819 model based on multiple, separate ssPA score matrices from each omics view (Fig. 2). There
820 are several ($k > 1$) predictor matrices here, hence the term ‘Multi-View’. In this work we used
821 a multi-block partial least squares (MB-PLS) model due to its ability to model multiple data
822 blocks (omics views) in relation to a response variable y . However, any multi-view
823 supervised machine learning technique could be used within the same framework. The MB-
824 PLS model was implemented using the mbpls Python package⁴⁰ using the NIPALS algorithm.
825 Again, all input omics data matrices are unit-variance scaled. As with PathIntegrate Single-
826 View, users can apply any ssPA algorithm implemented in the sspa package to perform the
827 first step of Multi-View, transforming each omics abundance matrix X_k of size $(N \times M_k)$ into a
828 pathway score matrix A_k of size $(N \times P_k)$. Then each pathway score matrix A_k is modelled by
829 MB-PLS, to predict an outcome variable. Important pathways are identified using the multi-
830 block variable importance in projection (MB-VIP) statistic, detailed below (Equation 14). In
831 this section we follow standard practice in describing how MB-PLS models an outcome Y
832 (which can be univariate or multivariate) using a several predictor matrices X_k , that, for
833 PathIntegrate, correspond to the pathway scores matrices A_k .

834 (Single block) partial least squares (PLS) regression⁷¹ is a supervised regression method
835 designed to work well on high-dimensional and highly co-linear datasets due to its latent
836 variable decomposition of both the predictor and response variables⁴¹. PLS performs a
837 simultaneous projection of the unit variance scaled predictor matrix X , of size $(N \times J)$, and a Y

838 response matrix, of size $(N \times H)$, into a lower dimensional space (defined by latent variables,
839 LVs) to maximise the covariance between the two projections (X scores, T and Y scores, U)
840 (Equation 4). The low dimensional representation of the X data can be used to predict Y
841 (Equation 6).

842 The PLS model as defined by Wold et al, 2001⁷¹ is as follows:

843 The X and Y matrices are decomposed into scores and loadings such that:

$$844 \quad X = TV^T + E$$

$$845 \quad Y = UC^T + F$$

846 *Equation 4*

847 Here, T and U represent X and Y scores respectively, each of size $(N \times R)$, for a model with R
848 latent variables. V, size $(J \times R)$, represents X loadings, C, size $(H \times R)$ represents Y weights, and
849 E and F refer to residual matrices, sizes $(N \times J)$ and $(N \times H)$ respectively, of independent and
850 identically distributed (iid) noise. Matrix transpose is denoted by ^T.

851 The X scores, T are linear combinations of the original X variables multiplied by the X weights
852 (coefficients):

$$853 \quad T = XW^*$$

854 *Equation 5*

855 Where W^* , size $(J \times R)$ denotes the weights matrix relating to the original variables, as opposed
856 to W , size $(J \times R)$, which denotes the weights matrix computed from the deflated matrices (see
857 Eqn. 8 below).

858 The X scores and Y weights are used to predict Y:

$$859 \quad \hat{Y} = TC^T + G$$

860 Equation 6, where G is a further residual matrix.

861 PLS is performed sequentially, obtaining scores, loadings, and weights for each of R latent
862 variables. Importantly, the first pair of latent vectors t and u are selected such that the
863 covariance between them is maximal:

$$864 \quad (t, u) = \underset{(t, u)}{\operatorname{argmax}} (\operatorname{cov}(t, u))$$

865 *Equation 7*

866 At each step, the model estimates, corresponding to the product of scores and loadings are
867 subtracted from the current X and Y matrices (this step is termed deflation) so that the next
868 set of latent vectors $r + 1$ can be computed from a new X_{r+1} and Y_{r+1} :

$$869 \quad X_{r+1} = X_r - t_r v_r^T$$

870
$$Y_{r+1} = Y_r - t_r c_r^T$$

871 *Equation 8*, with $X_1 = X$ and $Y_1 = Y$.

872 The optimal number of latent vectors is typically chosen using cross-validation approaches.

873 Using Equation 5, the prediction of Y can be re-written as:

874
$$\hat{Y} = XW^*C^T + G$$

875 *Equation 9*

876 (note the $*$ does not denote multiplication) and thus the regression coefficients for each X
877 variable are obtained using:

878
$$\beta = W^*C^T$$

879 *Equation 10*

880 The prediction of Y can finally be expressed in the form of a regression equation:

881
$$\hat{Y} = X\beta + G$$

882 *Equation 11*

883 Once the model is fit the scores, loadings, and weight matrices can be interpreted. Variable
884 selection approaches for PLS methods include inspection of β coefficients, as well as variable
885 importance in projection (VIP) ⁷². VIP is based on the PLS weights W weighted by the
886 proportion of Y explained in each latent variable (sum of squares) normalised by the total
887 sum of squares across all LVs, and explains the influence of each X feature on the model.

888 VIP for the j^{th} variable is given by⁷³:

889
$$\text{VIP}_j = \sqrt{\frac{J \cdot \sum_{r=1}^R (w_{rj}^2 \cdot \text{SSY}_r)}{\text{SSY}_{\text{cum}}}}$$

890 *Equation 12*

891 Here J represents the number of features in X , R is the number of latent variables (LVs), w_{rj}
892 is the weight of the j^{th} feature in the r^{th} LV, SSY_r is the sum of squares of Y explained by the
893 r^{th} LV, and SSY_{cum} is the cumulative sum of squares.

894 Often, variables with $\text{VIP} < 1$ are discarded, as the average of sum of squares of VIP scores is
895 equal to 1. However, a more reliable approach is to compute significance of the VIP values
896 using empirical p -value computation, described below in section 'Feature importance'.

897 Multi-block PLS is an extension of PLS that allows multiple data blocks $\{X_1, \dots, X_k\}$ as
898 predictors ⁴¹. The k^{th} X predictor block and Y response matrix can be decomposed as:

899
$$X_k = T_k V_k^T + E_k$$

900
$$Y = T_S C^T + F$$

901 *Equation 13*

902 *where T_S represents the X superscores.*

903 In the multi-block PLS case, block scores for each X block are combined to form superscores
904 $T_S = [T_1, T_2, \dots, T_K]$. The superscores are used to predict the response scores U , and also to
905 deflate the X_k blocks (if using the method proposed by Westerhuis and Coenegracht 1997),
906 rendering the superscores orthogonal.

907 VIP can be computed for MB-PLS models by using the superscores T_S across all blocks. In
908 Equation 14, SSY represents the proportion of Y explained across all X blocks, using the
909 superscores T_S rather than the scores T as in Eqn. 12 for single-block VIP.

910 MB-VIP for the j^{th} variable present in the k^{th} block is given by:

911

912
$$\text{MB-VIP}_j = \sqrt{\frac{f \cdot \sum_{r=1}^R (w_{krj}^2 \cdot SSY_r)}{SSY_{cum}}}$$

913 *Equation 14*

914 *where f is the number of features across all blocks.*

915 Similar to the original VIP definition, this MB-VIP metric satisfies the condition that the mean
916 of the sum of squares of VIP scores per X block equals 1.

917
$$\frac{SS(\text{MB-VIP})}{f \cdot k} = 1$$

918 *Equation 15*

919 *where $SS(\text{MB-VIP})$ represents the total sum of squares of the multi-block VIP values.*

920 **Univariate detection of pathway versus molecular-level signals**

921 Applying the semi-synthetic data generation approach detailed above, we generated semi-
922 synthetic data for each pathway accessible in the COPDgene and COVID-19 metabolomics
923 and proteomics datasets (1290 and 298 realisations for Reactome and KEGG respectively for
924 COPDgene data; 456 and 256 for COVID-19 data) at a range of different effect sizes.

925 For the pathway-level simulation, we used the ssPA kPCA method to generate ssPA scores
926 for each simulation. We then performed Mann Whitney U (WMU) tests to determine whether
927 there was a significant difference in the pathway scores of the target enriched pathway in
928 the simulated control and case groups. Bonferroni correction was used to obtain adjusted p -
929 values.

930 For the molecular-level simulation, we performed MWU tests to determine whether there
931 was a significant difference in each of the molecules in the target enriched pathway in the

932 simulated control and case groups. Bonferroni correction was used to obtain adjusted p -
933 values. To facilitate comparison with the pathway-level simulation, we used the Fisher
934 method to combine p -values from each molecule in the target pathway. If at least 50% of
935 molecules in the target pathway had a significant MWU test adjusted p -value (≤ 0.05), we
936 combined them using Fisher's method to obtain the final p -value. If less than 50% of the
937 molecules in the target pathway had an adjusted p -value of ≤ 0.05 , the combined p -value was
938 set to 1.

939 Performance evaluation

940 Unit-variance scaling, imputation, and ssPA transformation were performed separately on
941 the test-train splits in order to avoid data leakage when evaluating the results of multivariate
942 methods. Specifically, for ssPA, for each pathway the ssPA (PCA/kPCA) model is fit on the
943 training data only and ssPA scores for the test data are derived from the fitted model.
944 Hyperparameter tuning for the number of latent variables in the MBPLS/PLS models was
945 performed using 5-fold nested cross-validation, and for all semi-synthetic datasets the
946 optimal number of latent variables was 1 (as expected). Predictive performance was
947 computed using 5 times repeated 5-fold cross-validation, and evaluated using the area under
948 the Receiver Operator Characteristic (ROC) curve (AUC).

949 DIABLO

950 DIABLO requires tuning of a hyperparameter representing the design matrix, which
951 regulates the strength of correlation maximised between each omics block. In this work we
952 used DIABLO with a 'null' design (no correlation constraint) as in the original DIABLO
953 paper¹¹, as our simulation setup was not designed to incorporate correlations between
954 omics blocks.

955 Detection of target pathway simulation

956 For the target pathway simulation, we also used AUC to determine how well each method
957 was able to detect the artificially enriched target pathway in each simulation realisation. To
958 compute the AUC, the confusion matrix of true positives (TP), false positives (FP), true
959 negatives (TN) and false negatives (FN) was defined as follows:

- 960 • TP: The target enriched pathway with $p_{adj} \leq 0.05$
- 961 • FP: A non-target pathway with $p_{adj} \leq 0.05$
- 962 • TN: A non-target pathway with $p_{adj} > 0.05$
- 963 • FN: The target enriched pathway with $p_{adj} > 0.05$

964 p -values for each pathway's feature importance (e.g. VIP/MB-VIP/DIABLO loading) were
965 computed using permutation testing, see 'Feature importance' below.

966 When evaluating the ability of DIABLO to detect the enriched target pathway, we used two
967 methods referred to as 'DIABLO pathway (loading)' and 'DIABLO pathway (sparse)

968 loading)'. 'DIABLO pathway (loading)' involved using the loadings in a non-penalised single
969 component GCCA DIABLO model as the feature importances and calculating empirical p -
970 values for these loadings as described below. 'DIABLO pathway (sparse loading)' involves
971 using a sparse DIABLO rGCCA model with L1 penalty, where 5-fold, 5-times repeated cross-
972 validation is used to select the number of important features. Then, 25 bootstrap subsets of
973 the data are obtained (each containing 400 samples in the COPDgene data or 60 in the
974 COVID-19 data per class) and a sparse DIABLO model is fitted on each of these subsets. The
975 test statistic for feature importance is defined as the proportion of the 25 bootstraps in
976 which the pathway has a non-zero (sparse) loading. Intuitively, the target enriched
977 pathway should be of high importance to the sparse model and therefore often appear in
978 the significant features with a non-zero loading. Empirical p -values are also computed from
979 the 'DIABLO pathway (sparse loading)' test statistic as described below.

980 **Feature importance**

981 p -values for the significance of each feature (pathway) in the PathIntegrate models were
982 computed empirically using a standard permutation test. We permuted class labels (Y)
983 10,000 times to obtain p -values with a resolution of 0.0001. p -values for each feature were
984 calculated by counting the number of trials with test statistic (in this case VIP, MB-VIP,
985 DIABLO loading, or non-zero proportion for DIABLO sparse) greater than or equal to the
986 observed test statistic, and dividing this by 10,000. Multiple testing correction using the
987 Benjamini Hochberg FDR method was then applied.

988 **PathIntegrate network explorer app**

989 Plotly Dash Cytoscape v0.3.0 (<https://github.com/plotly/dash-cytoscape>) was used to
990 create the PathIntegrate network explorer app within the PathIntegrate python package. The
991 app can be launched from within the Python package and runs on a local host. NetworkX was
992 used to create the base network based on the Reactome pathway hierarchy, which was
993 downloaded from <https://reactome.org/download/> (ReactomePathwaysRelation.txt).
994 Nodes represent pathways and edges represent a parent-child relationship between them.
995 The app takes as input a PathIntegrate Multi-View or Single-View model object and uses
996 attributes such as feature importance to colour nodes.

997 **COPDgene case study**

998 A PathIntegrate Multi-View model was fitted to COPDgene metabolomics, proteomics, and
999 transcriptomics data, using multi-omics ssPA scores generated using the kPCA³¹ method.
1000 The optimal number of latent variables (4) used in the MBPLS model was identified using
1001 nested 5-fold cross-validation.

1002 The superscores were correlated to 260 clinical metadata variables using Spearman
1003 correlation, and p -values were corrected for using Bonferroni correction. Absolute
1004 correlations ≥ 0.3 and adjusted p -values ≤ 0.05 were used to filter for significantly correlated
1005 metadata variables.

1006 **COVID-19 case study**

1007 A PathIntegrate Single-View model was fitted to COVID-19 metabolomics and proteomics
1008 data, using multi-omics ssPA scores generated using the SVD (PLAGE³⁵) method, and
1009 employing a random forest for outcome prediction. The optimal hyperparameters for the
1010 SciKit-Learn RandomForestClassifier model selected via 5-fold cross-validation were:
1011 n_estimators=200, min_samples_split=2, min_samples_leaf=4, max_features='sqrt',
1012 max_depth=10, bootstrap=True, oob_score=True.

1013 **Identifying important pathways using PathIntegrate Single-View**

1014 Random forest recursive feature elimination with 5-fold cross validation was used to identify
1015 the optimal number of pathway features (20) for the Single-View model, implemented using
1016 the sklearn RFECV function.

1017 **Identifying important molecules within a pathway**

1018 For a pathway of interest, loadings on principal component 1 were used to represent the
1019 contribution of each molecule to the pathway scores across samples.

1020

1021 **Data and code availability**

1022 The COVID dataset is publicly available from Mendeley data
1023 (<https://data.mendeley.com/datasets/tzydswhhb5/5>)⁵⁶.

1024 The COPDgene multi-omics data can be found at the following sources: Clinical Data and
1025 SOMAScan data are available through COPDGene (<https://www.ncbi.nlm.nih.gov/gap/>, ID:
1026 phs000179.v6.p2). RNA-Seq data is available through dbGaP
1027 (<https://www.ncbi.nlm.nih.gov/gap/>, ID: phs000765.v3.p2). Metabolon data is available at
1028 Metabolomics Workbench (<https://www.metabolomicsworkbench.org/> ID: PR000907).

1029 PathIntegrate is available via the open-source PathIntegrate Python package
1030 (www.github.com/cwieder/PathIntegrate). Tutorials and documentation for PathIntegrate
1031 can be found at <https://cwieder.github.io/pathintegrate>. Source code for benchmarking and
1032 applications can be found at https://github.com/cwieder/PathIntegrate_scripts.

1033

1034 References

- 1035 1. Krassowski, M., Das, V., Sahu, S. K. & Misra, B. B. State of the Field in
1036 Multi-Omics Research: From Computational Needs to Data Mining and
1037 Sharing. *Front Genet* **11**, 1598 (2020).
- 1038 2. Subramanian, I., Verma, S., Kumar, S., Jere, A. & Anamika, K. Multi-omics
1039 Data Integration, Interpretation, and Its Application. *Bioinformatics and
1040 Biology Insights* vol. 14 Preprint at
1041 <https://doi.org/10.1177/1177932219899051> (2020).
- 1042 3. Eicher, T. *et al.* Metabolomics and multi-omics integration: A survey of
1043 computational methods and resources. *Metabolites* vol. 10 Preprint at
1044 <https://doi.org/10.3390/metabo10050202> (2020).
- 1045 4. Canzler, S. *et al.* Prospects and challenges of multi-omics data integration
1046 in toxicology. *Arch Toxicol* **94**, 371–388 (2020).
- 1047 5. Bersanelli, M. *et al.* Methods for the integration of multi-omics data:
1048 Mathematical aspects. *BMC Bioinformatics* **17**, 15 (2016).
- 1049 6. Huang, S., Chaudhary, K. & Garmire, L. X. More is better: Recent progress
1050 in multi-omics data integration methods. *Frontiers in Genetics* vol. 8 84
1051 Preprint at <https://doi.org/10.3389/fgene.2017.00084> (2017).
- 1052 7. Reel, P. S., Reel, S., Pearson, E., Trucco, E. & Jefferson, E. Using machine
1053 learning approaches for multi-omics data analysis: A review. *Biotechnol
1054 Adv* **49**, 107739 (2021).
- 1055 8. Zhang, L. *et al.* Deep learning-based multi-omics data integration reveals
1056 two prognostic subtypes in high-risk neuroblastoma. *Front Genet* **9**, 477
1057 (2018).
- 1058 9. Wang, T. *et al.* MOGONET integrates multi-omics data using graph
1059 convolutional networks allowing patient classification and biomarker
1060 identification. *Nature Communications 2021 12:1* **12**, 1–13 (2021).
- 1061 10. Yan, K. K., Zhao, H. & Pang, H. A comparison of graph- and kernel-based -
1062 omics data integration algorithms for classifying complex traits. *BMC
1063 Bioinformatics* **18**, 539 (2017).
- 1064 11. Singh, A. *et al.* DIABLO: An integrative approach for identifying key
1065 molecular drivers from multi-omics assays. *Bioinformatics* **35**, 3055–
1066 3062 (2019).
- 1067 12. Zhou, G., Ewald, J. & Xia, J. OmicsAnalyst: a comprehensive web-based
1068 platform for visual analytics of multi-omics data. *Nucleic Acids Res* **49**,
1069 W476–W482 (2021).

- 1070
1071
1072
13. Argelaguet, R. *et al.* Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets. *Mol Syst Biol* **14**, e8124 (2018).
- 1073
1074
14. Vahabi, N. & Michailidis, G. Unsupervised Multi-Omics Data Integration Methods: A Comprehensive Review. *Front Genet* **13**, 854752 (2022).
- 1075
1076
1077
15. Min, E. J. & Long, Q. Sparse multiple co-Inertia analysis with application to integrative analysis of multi-Omics data. *BMC Bioinformatics* **21**, 1–12 (2020).
- 1078
1079
1080
16. Kamburov, A., Cavill, R., Ebbels, T. M. D., Herwig, R. & Keun, H. C. Integrated pathway-level analysis of transcriptomics and metabolomics data with IMPaLA. *Bioinformatics* **27**, 2917–2918 (2011).
- 1081
1082
17. Paczkowska, M. *et al.* Integrative pathway enrichment analysis of multivariate omics data. *Nat Commun* **11**, 1–16 (2020).
- 1083
1084
1085
1086
18. Odom, G. J., Colaprico, A., Silva, T. C., Chen, X. S. & Wang, L. PathwayMultiomics: An R Package for Efficient Integrative Analysis of Multi-Omics Datasets With Matched or Un-matched Samples. *Front Genet* **12**, 783713 (2021).
- 1087
1088
1089
19. Canzler, S. & Hackermüller, J. multiGSEA: a GSEA-based pathway enrichment analysis for multi-omics data. *BMC Bioinformatics* **21**, 561 (2020).
- 1090
1091
1092
20. Rodríguez-Mier, P., Poupin, N., de Blasio, C., Le Cam, L. & Jourdan, F. DEXOM: Diversity-based enumeration of optimal context-specific metabolic networks. *PLoS Comput Biol* **17**, (2021).
- 1093
1094
21. Gillenwater, L. A. *et al.* Multi-omics subtyping pipeline for chronic obstructive pulmonary disease. *PLoS One* **16**, e0255337 (2021).
- 1095
1096
22. Mastej, E. *et al.* Identifying protein–metabolite networks associated with COPD phenotypes. *Metabolites* **10**, 124 (2020).
- 1097
1098
1099
23. Zhou, G., Pang, Z., Lu, Y., Ewald, J. & Xia, J. OmicsNet 2.0: a web-based platform for multi-omics integration and network visual analytics. *Nucleic Acids Res* **1**, 13–14 (2013).
- 1100
1101
1102
1103
24. Khatri, P., Sirota, M. & Butte, A. J. Ten years of pathway analysis: Current approaches and outstanding challenges. *PLoS Computational Biology* vol. 8 e1002375 Preprint at <https://doi.org/10.1371/journal.pcbi.1002375> (2012).
- 1104
1105
1106
25. Tavazoie, S., Hughes, J. D., Campbell, M. J., Cho, R. J. & Church, G. M. Systematic determination of genetic network architecture. *Nat Genet* **22**, 281–285 (1999).

- 1107 26. Subramanian, A. *et al.* Gene set enrichment analysis: A knowledge-based
1108 approach for interpreting genome-wide expression profiles. *Proc Natl*
1109 *Acad Sci U S A* **102**, 15545–15550 (2005).
- 1110 27. Maghsoudi, Z., Nguyen, H., Tavakkoli, A. & Nguyen, T. A comprehensive
1111 survey of the approaches for pathway analysis using multi-omics data
1112 integration. *Brief Bioinform* (2022) doi:10.1093/BIB/BBAC435.
- 1113 28. Liu, T. *et al.* PaintOmics 4: new tools for the integrative analysis of multi-
1114 omics datasets supported by multiple pathway databases. *Nucleic Acids*
1115 *Res* **50**, W551–W559 (2022).
- 1116 29. Segura-Lepe, M. P., Keun, H. C. & Ebbels, T. M. D. Predictive modelling
1117 using pathway scores: Robustness and significance of pathway
1118 collections. *BMC Bioinformatics* **20**, 543 (2019).
- 1119 30. Wu, S. *et al.* Integrated Machine Learning and Single-Sample Gene Set
1120 Enrichment Analysis Identifies a TGF-Beta Signaling Pathway Derived
1121 Score in Headneck Squamous Cell Carcinoma. *J Oncol* **2022**, (2022).
- 1122 31. Al-Akwaa, F. M., Yunits, B., Huang, S., Alhajaji, H. & Garmire, L. X. Lilikoi:
1123 an R package for personalized pathway-based classification modeling
1124 using metabolomics data. *Gigascience* **7**, 1 (2018).
- 1125 32. Fang, X. *et al.* Lilikoi V2.0: a deep learning-enabled, personalized
1126 pathway-based R package for diagnosis and prognosis predictions using
1127 metabolomics data. *Gigascience* **10**, 1–11 (2021).
- 1128 33. Wieder, C., Lai, R. P. J. & Ebbels, T. M. D. Single sample pathway analysis
1129 in metabolomics: performance evaluation and application. *BMC*
1130 *Bioinformatics* **23**, 481 (2022).
- 1131 34. Meng, C., Kuster, B., Culhane, A. C. & Gholami, A. M. A multivariate
1132 approach to the integration of multi-omics datasets. *BMC Bioinformatics*
1133 **15**, 162 (2014).
- 1134 35. Hänzelmann, S., Castelo, R. & Guinney, J. GSVA: Gene set variation
1135 analysis for microarray and RNA-Seq data. *BMC Bioinformatics* **14**, 7
1136 (2013).
- 1137 36. Lee, E., Chuang, H. Y., Kim, J. W., Ideker, T. & Lee, D. Inferring pathway
1138 activity toward precise disease classification. *PLoS Comput Biol* **4**,
1139 e1000217 (2008).
- 1140 37. Tomfohr, J., Lu, J. & Kepler, T. B. Pathway level analysis of gene
1141 expression using singular value decomposition. *BMC Bioinformatics* **6**,
1142 225 (2005).

- 1143 38. Li, Y., Wu, F. X. & Ngom, A. A review on machine learning principles for
1144 multi-view biological data integration. *Brief Bioinform* **19**, 325–340
1145 (2018).
- 1146 39. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of*
1147 *Machine Learning Research* **12**, 2825–2830 (2011).
- 1148 40. Baum, A. & Vermue, L. Multiblock PLS: Block dependent prediction
1149 modeling for Python. *J Open Source Softw* **4**, 1190 (2019).
- 1150 41. Westerhuis, J., ... T. K.-J. of C. & 1998, undefined. Analysis of multiblock
1151 and hierarchical PCA and PLS models. *Wiley Online Library* (1998)
1152 doi:10.1002/(SICI)1099-128X(199809/10)12:5<301::AID-
1153 CEM515>3.0.CO;2-S.
- 1154 42. Wangen, L. E. & Kowalski, B. R. A multiblock partial least squares
1155 algorithm for investigating complex chemical systems. *J Chemom* **3**, 3–20
1156 (1989).
- 1157 43. Rohart, F., Gautier, B., Singh, A. & Lê Cao, K.-A. mixOmics: An R package
1158 for 'omics feature selection and multiple data integration. *PLoS Comput*
1159 *Biol* **13**, e1005752 (2017).
- 1160 44. Wieder, C. *et al.* Pathway analysis in metabolomics: Recommendations
1161 for the use of over-representation analysis. *PLoS Comput Biol* **17**,
1162 e1009105 (2021).
- 1163 45. Regan, E. A. *et al.* Genetic Epidemiology of COPD (COPDGene) Study
1164 Design. <https://doi.org/10.3109/15412550903499522> **7**, 32–43 (2010).
- 1165 46. Schols, A. M. W. J. Nutritional and metabolic modulation in chronic
1166 obstructive pulmonary disease management. *European Respiratory*
1167 *Journal* **22**, 81s–86s (2003).
- 1168 47. Kao, C. C. *et al.* Glucose and pyruvate metabolism in severe chronic
1169 obstructive pulmonary disease. *J Appl Physiol* **112**, 42 (2012).
- 1170 48. Xuan, L. *et al.* Association between chronic obstructive pulmonary
1171 disease and serum lipid levels: a meta-analysis. *Lipids Health Dis* **17**,
1172 (2018).
- 1173 49. Gong, J. *et al.* Cigarette smoke reduces fatty acid catabolism, leading to
1174 apoptosis in lung endothelial cells: Implication for pathogenesis of
1175 COPD. *Front Pharmacol* **10**, 469190 (2019).
- 1176 50. Zhao, H., Dennery, P. A. & Yao, H. Metabolic reprogramming in the
1177 pathogenesis of chronic lung diseases, including BPD, COPD, and
1178 pulmonary fibrosis. *Am J Physiol Lung Cell Mol Physiol* **314**, L544–L554
1179 (2018).

- 1180
1181
1182
51. Suleman, M., Attia, A. & Elsammak, M. Carnitine deficiency in chronic obstructive pulmonary disease patients. *European Respiratory Journal* **42**, (2013).
- 1183
1184
1185
52. Conlon, T. M. *et al.* Metabolomics screening identifies reduced L-carnitine to be associated with progressive emphysema. *Clin Sci* **130**, 273–287 (2016).
- 1186
1187
1188
53. Agudelo, C. W. *et al.* Decreased surfactant lipids correlate with lung function in chronic obstructive pulmonary disease (COPD). *PLoS One* **15**, (2020).
- 1189
1190
1191
54. Tran, H. B. *et al.* AIM2 nuclear exit and inflammasome activation in chronic obstructive pulmonary disease and response to cigarette smoke. *Journal of Inflammation (United Kingdom)* **18**, 1–13 (2021).
- 1192
1193
1194
55. Kotlyarov, S. & Kotlyarova, A. Anti-Inflammatory Function of Fatty Acids and Involvement of Their Metabolites in the Resolution of Inflammation in Chronic Obstructive Pulmonary Disease. *Int J Mol Sci* **22**, (2021).
- 1195
1196
56. Su, Y. *et al.* Multi-Omics Resolves a Sharp Disease-State Shift between Mild and Moderate COVID-19. *Cell* **183**, 1479-1495.e20 (2020).
- 1197
1198
57. Donlan, A. N. *et al.* IL-13 is a driver of COVID-19 severity. *JCI Insight* **6**, (2021).
- 1199
1200
1201
58. Bader, S. M., Cooney, J. P., Pellegrini, M. & Doerflinger, M. Programmed cell death: the pathways to severe COVID-19? *Biochemical Journal* **479**, 609 (2022).
- 1202
1203
1204
59. Geiger, J. D., Khan, N., Murugan, M. & Boison, D. Possible Role of Adenosine in COVID-19 Pathogenesis and Therapeutic Opportunities. *Front Pharmacol* **11**, 594487 (2020).
- 1205
1206
1207
60. Meng, C. *et al.* MOGSA: Integrative single sample gene-set analysis of multiple omics data. *Molecular and Cellular Proteomics* **18**, S153–S168 (2019).
- 1208
1209
1210
61. Chowdhury, S. & Sarkar, R. R. Comparison of human cell signaling pathway databases—evolution, drawbacks and challenges. *Database* **2015**, 126 (2015).
- 1211
1212
62. Wittig, U. & De Beuckelaer, A. Analysis and comparison of metabolic pathway databases. *Brief Bioinform* **2**, 126–142 (2001).
- 1213
1214
1215
63. Mubeen, S. *et al.* The Impact of Pathway Database Choice on Statistical Enrichment Analysis and Predictive Modeling. *Front Genet* **10**, 1203 (2019).

- 1216
1217
64. Fahy, E. *et al.* Update of the LIPID MAPS comprehensive classification system for lipids. *J Lipid Res* **50**, S9–S14 (2009).
- 1218
1219
1220
1221
65. Wörheide, M. A., Krumsiek, J., Kastenmüller, G. & Arnold, M. Multi-omics integration in biomedical research – A metabolomics-centric review. *Analytica Chimica Acta* vol. 1141 144–162 Preprint at <https://doi.org/10.1016/j.aca.2020.10.038> (2021).
- 1222
1223
1224
1225
66. Karp, P. D., Midford, P. E., Caspi, R. & Khodursky, A. Pathway size matters: the influence of pathway granularity on over-representation (enrichment analysis) statistics. *BMC Genomics* *2021 22:1* **22**, 1–11 (2021).
- 1226
1227
1228
67. Mubeen, S., Tom Kodamullil, A., Hofmann-Apitius, M. & Domingo-Fernández, D. On the influence of several factors on pathway enrichment analysis. *Brief Bioinform* (2022) doi:10.1093/BIB/BBAC143.
- 1229
1230
68. Martens, H. & Martens, M. *Multivariate analysis of quality: an introduction*. (2001).
- 1231
1232
69. Flores, J. E. *et al.* Missing data in multi-omics integration: Recent advances through artificial intelligence. *Front Artif Intell* **6**, (2023).
- 1233
1234
1235
70. Pang, Z. *et al.* MetaboAnalyst 5.0: narrowing the gap between raw spectra and functional insights. *Nucleic Acids Res* (2021) doi:10.1093/nar/gkab382.
- 1236
1237
1238
71. Wold, S., Sjöström, M. & Eriksson, L. PLS-regression: A basic tool of chemometrics. in *Chemometrics and Intelligent Laboratory Systems* vol. 58 109–130 (Elsevier, 2001).
- 1239
1240
1241
1242
72. Farrés, M., Platikanov, S., Tsakovski, S. & Tauler, R. Comparison of the variable importance in projection (VIP) and of the selectivity ratio (SR) methods for variable selection and interpretation. *J Chemom* **29**, 528–536 (2015).
- 1243
1244
1245
1246
73. Mendez, K. M., Broadhurst, D. I. & Reinke, S. N. Migrating from partial least squares discriminant analysis to artificial neural networks: a comparison of functionally equivalent visualisation and feature contribution tools using jupyter notebooks. *Metabolomics* **16**, 17 (2020).
- 1247
- 1248
- 1249
- 1250

1251 **Acknowledgements**

1252 **COPDgene grant support and disclaimer:**

1253 This work was supported by NHLBI grants U01 HL089897 and U01 HL089856 and by NIH
1254 contract 75N92023D00011. The COPDGene study (NCT00608764) is also supported by the
1255 COPD Foundation through contributions made to an Industry Advisory Committee that has
1256 included AstraZeneca, Bayer Pharmaceuticals, Boehringer-Ingelheim, Genentech,
1257 GlaxoSmithKline, Novartis, Pfizer, and Sunovion. The content is solely the responsibility of
1258 the authors and does not necessarily represent the official views of the National Heart, Lung,
1259 and Blood Institute or the National Institutes of Health.

1260 We gratefully acknowledge Peter Castaldi and Craig Hersh for generating the COPDgene
1261 RNAseq data.

1262 **Funding acknowledgments:**

1263 This research was funded in whole, or in part, by the Wellcome Trust [222837/Z/21/Z]. For
1264 the purpose of open access, the author has applied a CC BY public copyright licence to any
1265 Author Accepted Manuscript version arising from this submission.

1266 TE acknowledges partial support from BBSRC grants BB/T007974/1 and BB/W002345/1.

1267 RPJL was supported by a UK MRC fellowship (MR/R008922/1) which is part of the EDCTP2
1268 programme supported by the European Union and a NIH-NIAID grant (R01 AI145436).

1269 JC is supported by a state-funded PhD contract (MESRI (Minister of Higher Education,
1270 Research and Innovation)).

1271 FJ - This research was funded by the Agence Nationale de la Recherche (ANR, French National
1272 Research Agency)—MetaboHUB, the national metabolomics and fluxomics infrastructure
1273 (Grant ANR-INBS-0010).

1274 KK, RB - Research reported in this publication was supported by the National Heart Lung
1275 Blood Institute of the National Institutes of Health to KK and RB under award number
1276 R01HL152735.

1277

1278 **Author Contributions**

1279 CW, RPJL, and TE conceptualised the project and designed the study. TE and RPJL supervised
1280 the project. CW performed the data analysis and developed the PathIntegrate Python
1281 package. KK and RB facilitated access to COPDgene data. CW produced the initial manuscript
1282 draft. All authors provided critical feedback and helped shape the research, analysis, and
1283 manuscript. All authors contributed to the final version of the manuscript and approved it.

1284 **Competing Interests statement**

1285 The authors declare no competing interests.

1286