

Genome Biology

DIABLO: identifying key molecular drivers from multi-omic assays, an integrative approach --Manuscript Draft--

Manuscript Number:		
Full Title:	DIABLO: identifying key molecular drivers from multi-omic assays, an integrative approach	
Article Type:	Method	
Funding Information:	National Institute Of Allergy And Infectious Diseases of the National Institutes of Health (U19AI118608)	Dr Casey P Shannon
	the National Institute Of Allergy And Infectious Diseases of the National Institutes of Health (U19AI118608)	A/Prof Scott J Tebbutt
	National Health and Medical Research Council (GNT1087415)	Dr Kim-Anh Lê Cao
	Canadian Institutes of Health Research Doctoral Award (NA)	Dr Amrit Singh
Abstract:	Systems biology approaches, leveraging multi-omics measurements, are needed to capture the complexity of biological networks while identifying the key molecular drivers of disease mechanisms. We present DIABLO, a novel integrative method to identify multi-omics biomarker panels that can discriminate between multiple phenotypic groups. In the multi-omics analyses of simulated and real-world datasets, DIABLO resulted in superior biological enrichment compared to other integrative methods, and achieved comparable predictive performance with existing multi-step classification schemes. DIABLO is a versatile approach that will benefit a diverse range of research areas, where multiple high dimensional datasets are available for the same set of specimens. DIABLO is implemented along with tools for model selection, and validation, as well as graphical outputs to assist in the interpretation of these integrative analyses (http://mixomics.org/).	
Corresponding Author:	Kim-Anh Lê Cao, Ph.D University of Melbourne Brisbane, Queensland AUSTRALIA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	University of Melbourne	
Corresponding Author's Secondary Institution:		
First Author:	Amrit Singh	
First Author Secondary Information:		
Order of Authors:	Amrit Singh Casey P Shannon Benoît Gautier Florian Rohart Michaël Vacher Scott J Tebbutt Kim-Anh Lê Cao, Ph.D	

Order of Authors Secondary Information:	
Suggested Reviewers:	<p>Marylyn D Ritchie Pennsylvania State University marylyn@pennmedicine.upenn.edu Prof Ritchie is a computational biologist, and an expert in multi omics data integration and publishes methods in this specific area of research</p> <p>Susan Holmes Stanford University susan@stat.stanford.edu Prof Holmes is a statistician internationally renowned for the type of multivariate methods we have developed in this manuscript, and their application to biological data.</p> <p>Aedin Culhane Dana-Farber Cancer Institute aedin@jimmy.harvard.edu Prof Culhane is a statistician with expertise in the type of methods we developed in this manuscript for multi omics data integration.</p> <p>Levi Waldron City University of New York - Hunter College levi.waldron@sph.cuny.edu A/Prof Waldron is a statistician expert in omics data analysis and multi omics data integration in particular.</p>
Additional Information:	
Question	Response
Has this manuscript been submitted to this journal before?	Yes
Please provide the manuscript identification number from your previous submission. If you no longer have the identification number, please specify this in the text box below. as follow-up to "Has this manuscript been submitted to this journal before?"	GBIO-D-16-01112
Experimental design and statistics	Yes
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends. Have you included all the information requested in your manuscript?	
Resources	Yes
A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly	

encouraged to cite [Research Resource Identifiers](#) (RRIDs) for antibodies, model organisms and tools, where possible.

Have you included the information requested as detailed in our [Minimum Standards Reporting Checklist](#)?

[Click here to view linked References](#)

1
2
3
4 1 **DIABLO: identifying key molecular drivers from multi-omic assays, an integrative**
5 2 **approach**
6
7

8 4 Amrit Singh^{1,2,3}, Casey P. Shannon³, Benoît Gautier⁴, Florian Rohart⁵, Michaël Vacher^{6,9}, Scott
9 5 J. Tebbutt^{1,3,7}, Kim-Anh Lê Cao⁸
10
11

12 7 ¹Centre for Heart Lung Innovation, St. Paul's Hospital, University of British Columbia,
13 8 Vancouver, BC, Canada;
14

15 10 ²Department of Pathology and Laboratory Medicine, University of British Columbia, Vancouver,
16 BC, Canada;

17 11 ³Prevention of Organ Failure (PROOF) Centre of Excellence, Vancouver, BC, Canada.

18 12 ⁴The University of Queensland Diamantina Institute, Translational Research Institute,
19 13 Woolloongabba, QLD 4102, Australia

20 14 ⁵Institute for Molecular Bioscience, The University of Queensland, St Lucia, QLD 4072,
21 15 Australia

22 16 ⁶Australian Research Council Centre of Excellence in Plant Energy Biology, The University of
23 17 Western Australia, Crawley, Western Australia, Australia

24 18 ⁷Department of Medicine (Respiratory Division), University of British Columbia, Vancouver,
25 19 BC, Canada.

26 20 ⁸Melbourne Integrative Genomics, School of Mathematics and Statistics, The University of
27 21 Melbourne, Melbourne, Australia

28 22 ⁹current address: Australian eHealth Research Centre, Commonwealth Scientific and Industrial
29 23 Research Organisation, Brisbane, Queensland, Australia

30 24 Corresponding author:
31 25

32 26 Dr Kim-Anh Lê Cao

33 27 Melbourne Integrative Genomics and School of Mathematics and Statistics, The University of
34 28 Melbourne, Melbourne, Australia

35 29 T: +61 (0)3834 43971

36 30 kimanh.lecao@unimelb.edu.au

37 31

38 32

39 33

40 34

41 35

42 36

43 37

44 38

45 39

46 40

47 41

48 42

49 43

50 44

51 45

52 46

53 47

54 48

55 49

56 50

57 51

58 52

59 53

60 54

61 55

62 56

63 57

64 58

65 59

1
2
3
4 40
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 **41 Abstract**
5
6
7 **42 Systems biology approaches, leveraging multi-omics measurements, are needed to capture the**
8
9 **43 complexity of biological networks while identifying the key molecular drivers of disease**
10
11 **44 mechanisms. We present DIABLO, a novel integrative method to identify multi-omics**
12
13 **45 biomarker panels that can discriminate between multiple phenotypic groups. In the multi-omics**
14
15 **46 analyses of simulated and real-world datasets, DIABLO resulted in superior biological**
16
17 **47 enrichment compared to other integrative methods, and achieved comparable predictive**
18
19 **48 performance with existing multi-step classification schemes. DIABLO is a versatile approach**
20
21 **49 that will benefit a diverse range of research areas, where multiple high dimensional datasets are**
22
23 **50 available for the same set of specimens. DIABLO is implemented along with tools for model**
24
25 **51 selection, and validation, as well as graphical outputs to assist in the interpretation of these**
26
27 **52 integrative analyses (<http://mixomics.org/>).**

32 **53**
33
34
35

36 **54 Keywords: Systems biology, biomarkers, data integration, data visualization, asthma,**
37
38 **55 classification, breast cancer, multi-omics, network analysis**

40 **56**
41
42

43 **57**
44

45 **58**
46

47
48

49
50

51
52

53
54

55
56

57
58

59
60

61
62

63
64

1
2
3
4 59 **Background**
5
6

7 60 Technological improvements have allowed for the collection of data from different molecular
8 61 compartments (*e.g.*, gene expression, methylation status, protein abundance) resulting in multiple
9 62 omics (multi-omics) data from the same set of biospecimens (*e.g.*, transcriptomics, proteomics,
10 63 metabolomics). The large number of omic variables compared to the limited number of available
11 64 biological samples presents a computational challenge when identifying the key drivers of
12 65 disease. Further, technological limitations differ with respect to different omic platforms (*e.g.*,
13 66 sequencing *vs.* mass spectrometry), and biological effect sizes differ with respect to different
14 67 omic variable-types (*e.g.*, methylation status *vs.* protein expression). Effective integrative
15 68 strategies are needed, to extract common biological information spanning multiple molecular
16 69 compartments that explains phenotypic variation. Already, systems biology approaches which
17 70 incorporated data from multiple biological compartments, have shown improved biological
18 71 insights compared to traditional single omics analyses [1–3]. This may be because single omics
19 72 analyses cannot account for the interactions between omic layers and, consequently, are unable
20 73 to reconstruct accurate molecular networks. These molecular networks are dynamic, changing
21 74 under perturbed conditions such as disease, response to therapy, and environmental exposures.
22 75 Therefore, adopting a holistic approach by integrating multi-omics data may bridge this
23 76 information gap, and uncover networks that are representative of the underlying molecular
24 77 mechanisms [4,5].

25 78 Preliminary approaches to data integration included multi-step approaches that leveraged
26 79 existing single-omics methods: multi-omics data were concatenated, or ensembles of single
27 80 omics models created [6]. These approaches can be biased towards certain omics data types,
28 81 however, and do not account for interactions between omic layers [7,8]. Recently, more
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 82 sophisticated integrative approaches have been proposed (**Supplementary Fig. 1**) [4,9–12].
5
6 83 They can be broadly divided into unsupervised analyses, which identify coherent relationships
7
8 84 across multi-omics datasets when samples are unlabeled, and supervised analyses, which identify
9
10 85 multi-omics patterns that discriminate between known phenotypic sample groups. However these
11
12 86 supervised strategies are unable to capture the shared information across multiple biological
13
14 87 domains when identifying the key molecular drivers associated with a phenotype. Such methods
15
16 88 are needed to capture the dynamic nature of molecular networks under various disease conditions
17
18 89 and ultimately provide robust biomarkers that are both biologically and clinically relevant.
19
20
21
22
23
24 90 To address these knowledge gaps, we introduce DIABLO, a method that incorporates
25
26 91 information across high dimensional multi-omics data while discriminating phenotypic groups.
27
28 92 DIABLO uncovers robust biomarkers of dysregulated disease processes that span multiple
29
30 93 functional layers. We demonstrate the capabilities and versatility of DIABLO both in simulated
31
32 94 and real-world data, integrating multi-omics datasets to identify relevant biomarkers of various
33
34 95 diseases. DIABLO is available through the mixOmics data integration toolkit
35
36 96 (www.mixomics.org [12]) which contains a wide range of multivariate methods for the
37
38 97 exploration and integration of high dimensional biological datasets.
39
40
41
42
43 98
44
45
46 99 **Results**
47
48 100 DIABLO (**D**ata **I**ntegration **A**nalysis for **B**iomarker discovery using **L**atent **c**OMPONENTs)
49
50 101 maximizes the common or correlated information between multiple omics (multi-omics) datasets
51
52 102 while identifying the key omics variables (mRNA, miRNA, CpGs, proteins, metabolites, *etc.*)
53
54 103 and characterizing the disease sub-groups or phenotypes of interest. DIABLO uses Projection to
55
56 104 Latent Structure models (PLS) [13], and extends both sparse PLS-Discriminant Analysis [14] to
57
58
59
60
61
62
63
64
65

1
2
3
4 105 multi-omics analyses and sparse Generalized Canonical Correlation Analysis [15] to a
5
6 106 supervised analysis framework. In contrast to existing penalized matrix decomposition methods
7
8 107 [16], DIABLO is a component-based method (or a dimension reduction technique) that
9
10 108 transforms each omic dataset into latent components and maximizes the sum of pairwise
11
12 109 correlations between latent components (user-defined) and a phenotype of interest [17].
13
14 110 DIABLO is, therefore, an integrative classification method that builds predictive multi-omics
15
16 111 models that can be applied to multi-omics data from new samples to determine their phenotype.
17
18 112 Users can specify the number of variables to select from each dataset and visualize the omics
19
20 113 data and the multi-omics panel into a reduced data. The method is highly flexible in the type of
21
22 114 experimental design it can handle, ranging from classical single time point to cross-over and
23
24 115 repeated measures studies. Modular-based analysis can also be incorporated using pathway-
25
26 116 based module matrices [18] instead of the original omics matrices, as illustrated in one of our
27
28 117 case studies.

35
36 118
37
38 119 **DIABLO selects correlated and discriminatory variables**

40
41 120 Briefly, three omic datasets consisting of 200 samples (split equally over two groups) and 260
42
43 121 variables were generated by modifying the degree of correlation and discrimination, resulting in
44
45 122 four types of variables: 30 correlated-discriminatory (corDis) variables, 30 uncorrelated-
46
47 123 discriminatory (unCorDis) variables, 100 correlated-nondiscriminatory (corNonDis) variables,
48
49 124 and 100 uncorrelated-nondiscriminatory (unCorNonDis) variables (**Supplementary Note**,
50
51 125 **Supplementary Fig. 2**). Three integrative classification methods were applied to generate multi-
52
53 126 omic biomarkers panels of 90 variables each (30 variables from each omic dataset): a DIABLO
54
55 127 model with either a full design (where the correlation between all pairwise combinations of
56
57
58
59
60
61
62
63
64
65

1
2
3
4 128 datasets, as well as between each dataset and the phenotypic outcome, were maximised) or the
5
6 129 null design (where only the correlation between each dataset and the phenotypic outcome was
7
8 130 maximised, see **Methods**), a concatenation-based sPLSDA classifier which consists of naively
9
10 131 combining all datasets into one, and an ensemble of sPLSDA classifiers where a separate
11
12 132 sPLSDA classifier was fitted for each omics dataset and the consensus predictions were
13
14 133 combined using a majority vote scheme (see **Supplementary Fig. 3**). The purpose of the
15
16 134 simulation study was to compare DIABLO models with existing multi-step integrative classifiers
17
18 20
21 135 with respect to the error rate and types of variables selected as part of the multi-omic biomarker
22
23 136 panels. A secondary aim was to determine the effect of design matrix on the resulting multi-omic
24
25 137 biomarker panels identified using DIABLO.
26
27
28 138 The concatenation, ensemble and DIABLO_null classifiers performed similarly across
29
30 139 the various noise and fold-change thresholds. At lower noise levels (simulated using a
31
32 140 multivariate normal distribution with mean of zero and standard deviation of 0.2 or 0.5) the
33
34 141 DIABLO_full classifier had a slightly higher error rate compared to the other approaches (**Fig.**
35
36 142 **1a**), but consistently selected mostly correlated and discriminatory (corDis) variables, unlike the
37
38 143 other integrative classifiers (**Fig. 1b**). All methods behaved similarly with respect to the error
39
40 144 rate and types of variables selected at higher noise thresholds (simulated using a multivariate
41
42 145 normal distribution with mean of zero and standard deviation of 1 or 2). This simulation
43
44 146 highlights how the design (connection between datasets) affects the flexibility of the DIABLO
45
46 147 model, resulting in a trade-off between discrimination or correlation. DIABLO_null focused on
47
48 148 selecting discriminatory variables and disregarded most of the correlation between datasets (null
49
50 149 design), whereas DIABLO_full selected highly correlated variables across all datasets. Since the
51
52 150 variables selected by DIABLO_full reflect the correlation structure between biological

53
54
55
56
57
58
59
60
61
62
63
64
65

1 compartments, we hypothesized that they might provide a balance between prediction accuracy
2 and biological insight.

3

4 151

5 152

6 153

7 154 **DIABLO identifies molecular networks with superior biological enrichment**

8 155 To assess this, we turn to real biological datasets. We applied various integrative approaches to
9 156 cancer multi-omics datasets (mRNA, miRNA, and CpG) – colon, kidney, glioblastoma (gbm)
10 157 and lung – and identified multi-omics biomarker panels that were predictive of high and low
11 158 survival times (**Table 1**). We then compared the network properties and biological enrichment of
12 159 the selected features across approaches.

13

14 160 Multi-omics biomarker panels were developed using component-based integrative

15 161 approaches that also performed variable selection: supervised methods included concatenation
16 162 and ensemble schemes using the sPLSDA classifier [14], and DIABLO with either the null or
17 163 full design (DIABLO_null, and DIABLO_full); unsupervised approaches included sparse
18 164 generalized canonical correlation analysis [15] (sGCCA), Multi-Omics Factor Analysis
19 165 (MOFA), and Joint and Individual Variation Explained (JIVE) [23] (see **Supplementary Note**
20 166 for parameter settings). Both supervised and unsupervised approaches were considered in order
21 167 to compare and contrast the types of omics-variables selected, network properties and biological
22 168 enrichment results. A distinction was made between DIABLO models in which the correlation
23 169 between omics datasets was not maximized (DIABLO_null) and those when the correlation
24 170 between omics datasets was maximized (DIABLO_full).

25

26 171 Each multi-omics biomarker panel included 180 features (60 features of each omics type
27 172 across 2 components). Approaches generally identified distinct sets of features. **Fig. 2a** depicts
28 173 the distinct and shared features between the seven multi-omics panels obtained from the

1
2
3
4 174 unsupervised (purple, sGCCA, MOFA and JIVE) and supervised (green, Concatenation,
5 Ensemble, DIABLO_null and DIABLO_full) methods. Supervised methods selected many of the
6
7 175 same features (blue), but DIABLO_full had greater feature overlap with unsupervised methods
8
9 176 (orange). The level of connectivity of each of the seven multi-omics panels was assessed by
10
11 177 generating networks from the feature adjacency matrix at various Pearson correlation coefficient
12
13 178 cut-offs (**Fig. 2b**). At all cut-offs, unsupervised approaches produced networks with greater
14
15 179 connectivity (number of edges) compared to supervised approaches. In addition, biomarker
16
17 180 panels identified by DIABLO_full, were more similar to those identified by unsupervised
18
19 181 approaches, including high graph density, low number of communities and large number of
20
21 182 triads, indicating that DIABLO_full identified discriminative sets of features that were tightly
22
23 183 correlated across biological compartments (**Supplementary Fig. 4**). For example, **Fig. 2c** (upper
24
25 184 panel) depicts the networks of all multi-omics biomarker panels for the colon cancer dataset,
26
27 185 which show higher modularity (a limited number of large clusters of variables; circled) for the
28
29 186 DIABLO_full and the unsupervised approaches as compared to the supervised ones. The
30
31 187 corresponding component plots show a clear separation between the high and low survival
32
33 188 groups for the panels derived using supervised approaches, whereas the unsupervised approaches
34
35 189 could not segregate the survival groups [**Fig. 2c** (lower panel), see **Supplementary Fig. 5 and 6**
36
37 190 for other cancer datasets].
38
39
40
41
42
43
44
45
46
47
48 192 Finally, we carried out gene set enrichment analysis on each multi-omics biomarker panel
49
50 193 (using gene symbols of mRNAs and CpGs) against 10 gene set collections (see **Methods**) and
51
52 194 tabulated the number of significant (FDR=5%) gene sets (**Table 2**). The DIABLO_full model
53
54 195 identified the greatest number of significant gene sets across the 10 gene set collections and
55
56 196 generally ranked higher than the other methods in the colon (7 collections), gbm (5 collections)
57
58
59
60
61
62
63
64
65

1
2
3
4 197 and lung (5 collections) cancer datasets, whereas JIVE outperformed all other methods in the
5
6 198 kidney cancer datasets (6 collections). Unlike all other approaches considered, DIABLO_full,
7
8 199 which aimed to explain both the correlation structure between multiple omics layers and a
9
10 200 phenotype of interest, implicated the greatest number of known biological gene sets
11
12 201 (pathways/functions/processes *etc.*).
13
14 202
15
16 203 **Case study 1: DIABLO identified known and novel multi-omics biomarkers of breast**
17
18 204 **cancer subtypes**
19
20 205 We next demonstrate that DIABLO can identify novel biomarkers in addition to biomarkers with
21
22 206 known biological associations using a case study of human breast cancer. We applied our
23
24 207 biomarker analysis workflow to breast cancer datasets to characterize and predict PAM50 breast
25
26 208 cancer subtypes (**Supplementary Fig. 7**). After preprocessing and normalization of each omics
27
28 209 data-type, the samples were divided into training and test sets (**Methods, Table 1**). The training
29
30 210 data consisted of four omics-datasets (mRNA, miRNA, CpGs and proteins) whereas the test data
31
32 211 included all remaining samples for which the protein expression data were missing. The optimal
33
34 212 multi-omics biomarker panel size was identified using a grid approach where, for any given
35
36 213 combination of variables, we assessed the classification performance using a 5-fold cross-
37
38 214 validation repeated 5 times (**Supplementary Fig. 8**). The number of variables that resulted in the
39
40 215 minimum balanced error rate were retained as previously described in [12]. The optimal multi-
41
42 216 omics panel consisted of 45 mRNA, 45 miRNAs, 25 CpGs and 55 proteins selected across three
43
44 217 components with a balanced error rate of $17.9 \pm 1.9\%$. This panel identified many variables with
45
46 218 previously known associations with breast cancer, as assessed by looking at the overlap between
47
48 219 the panel features and gene sets related to breast cancer based on the Molecular Signature
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 220 database (MolSigDB) [23], miRCancer [24], Online Mendelian Inheritance in Man (OMIM)
5
6 221 [25], and DriverDBv2 [26]. **Figure 3a** depicts the variable contributions of each omics-type
7
8 222 indicated by their loading weight (variable importance). Variables not found in any database may
9
10 223 represent novel biomarkers of breast cancer. **Figure 3b** shows the consensus and individual
11
12 224 omics component plots based on this biomarker panel, along with 95% confidence ellipses
13
14 225 obtained from the training data and superimposed with the samples from the test data. The
15
16 226 majority of the samples were within the ellipses, suggesting a reproducible multi-omics
17
18 227 biomarker panel from the training to the test set, that was predictive of breast cancer subtypes
19
20 228 (balanced error rate = 22.9%). The consensus plot corresponded strongly with the mRNA
21
22 229 component plot, depicting a strong separation of the Basal (error rate = 4.9%) and Her2 (error
23
24 230 rate = 20%) subtypes. We observed a weak separation of Luminal A (LumA, error rate = 13.3%)
25
26 231 and Luminal B (LumB, error rate = 53.3%) subtypes. Similarly, the heatmap showing the scaled
27
28 232 expression of all features of the multi-omics biomarker panel, depicted a strong clustering of the
29
30 233 Basal and Her2 samples whereas the Luminal A and B were mixed (**Fig. 3c**). Overall, the
31
32 234 features of the multi-omics biomarker panel formed a densely connected network comprising of
33
34 235 four communities where variables in each community (cluster) were densely connected with
35
36 236 themselves and sparsely connected with other clusters (**Fig. 3d**). The largest cluster in **Fig. 3d**
37
38 237 consisted of 72 variables; 20 mRNAs, 21 miRNAs, 15 CpGs and 16 proteins (red bubble) and
39
40 238 was further investigated using gene set enrichment analysis. We identified many cancer-
41
42 239 associated pathways (*e.g.* FOXM1 pathway, p53 signaling pathway), DNA damage and repair
43
44 236 pathways (*e.g.* E2F mediated regulation of DNA replication, G2M DNA damage checkpoint)
45
46 237 and various cell-cycle pathways (*e.g.* G1S transition, mitotic G1/G1S phases), demonstrating the
47
48 238 ability of DIABLO to identify a biologically plausible multi-omics biomarker panel. This panel
49
50 239
51
52 240
53
54 241
55
56 242
57
58 243
59
60
61
62
63
64
65

1
2
3
4 243 generalized to new breast cancer samples and implicated previously unknown molecular features
5
6 244 in breast cancer, which could be further validated in experimental studies.
7
8 245
9
10
11 246 **Case study 2: DIABLO for repeated measures designs and module-based analyses**
12
13 247 Next, we demonstrate the flexibility of DIABLO by extending its use to a repeated measures
14
15 248 cross-over study [27], as well as incorporating module-based analyses that incorporate prior
16
17 249 biological knowledge [28–30]. We use a small multi-omics asthma dataset, including pre and
18
19 250 post intervention timepoints, to compare a DIABLO model that can account for repeated
20
21 251 measures (multilevel DIABLO) with the standard DIABLO model as described above [20,21].
22
23
24 252 An allergen inhalation challenge was performed as we previously described in [20,21] in 14
25
26 253 subjects and blood samples were collected before (pre) and two hours after (post) challenge; cell-
27
28 254 type frequencies, leukocyte gene transcript expression and plasma metabolite abundances were
29
30 255 determined for all samples (**Table 1**). We observed a net decline in lung function after allergen
31
32 256 inhalation challenge (**Supplementary Fig. 9**), and the goal of this study was to identify
33
34 257 perturbed molecular mechanisms in the blood in response to allergen inhalation challenge. A
35
36 258 module based approach (also known as eigengene summarization [18], **see Methods**) was used
37
38 259 to transform both the gene expression and metabolite datasets into pathway datasets.
39
40 260 Consequently, each variable in those two datasets now represented the scaled pathway activity
41
42 261 expression level for each sample instead of direct gene/metabolite expression. The mRNA
43
44 262 dataset was transformed into a dataset of metabolic pathways (based on the Kyoto Encyclopedia
45
46 263 of Genes and Genomes, KEGG) whereas the metabolite dataset was transformed into a
47
48 264 metabolite pathway dataset based on annotations provided by Metabolon Inc. (Durham, North
49
50 265 Carolina, USA) (**Fig. 4a**). To account for the repeated measures experimental design, a
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 266 multilevel approach [27] was first used to isolate the within-sample variation from each dataset
5
6 267 (see **Methods**), and then DIABLO was applied to identify a multi-omics biomarker panel
7
8 268 consisting of cells, gene and metabolite modules that discriminated pre- from post-challenge
9
10 269 samples. We contrast the resulting ‘multilevel DIABLO’ (mDIABLO) with a standard DIABLO
11
12 270 model that disregards the paired nature of this study by comparing their cross-validation
13
14 271 classification performances (**Fig. 4b**). mDIABLO outperformed DIABLO (AUC=98.5% vs.
15
16 272 AUC=62.2%, leave-one-out cross-validation, see **Methods**), and we observed a greater degree of
17
18 273 separation between the pre- and post-challenge samples for mDIABLO compared to DIABLO
19
20
21 274 (**Fig. 4c**). Common features (pathways) were identified across omics-types in the mDIABLO
22
23
24 275 model, but not in the standard DIABLO model (**Fig. 4d**). Tryptophan metabolism and Valine,
25
26
27 276 leucine and isoleucine metabolism pathways were identified in both the gene and metabolite
28
29
30 277 module datasets using mDIABLO. The heatmap of pairwise associations of all features identified
31
32 278 with mDIABLO demonstrated the ability of DIABLO to select groups of correlated features
33
34
35 279 which were predictive of pre- and post-challenge samples. The Asthma pathway was also
36
37
38 280 identified [even though individual gene members were not significantly altered post-challenge
39
40
41 281 (**Supplementary Fig. 10**)] and was negatively associated with Butanoate metabolism and
42
43
44 282 positively associated with basophils, a hallmark cell-type in asthma (**Fig. 4e**). These findings
45
46
47 283 depict DIABLO’s flexibility and sensitivity to detect subtle differences between repeated
48
49
50 284 designs, and its ability to identify common molecular processes spanning different biological
51
52 285 layers. The biological pathways identified suggest a mechanistic link with response to allergen
53
54
55 286 challenge.
56
57
58 288
59
60
61
62
63
64
65

1
2
3
4 289 **Discussion**
5
6
7 290 DIABLO aims to identify coherent patterns between datasets that change with respect different
8 phenotypes. This purely data-driven, holistic, and hypothesis-free tool can be used to derive
9 291 robust biomarkers and, ultimately, improve our understanding of the molecular mechanisms that
10 drive disease.
11
12
13
14
15

16 294 We found that unsupervised methods identified features that formed strong
17 interconnected multi-omics networks, but had poor discriminative ability. In contrast, features
18 295 identified by supervised methods were discriminative, but formed sparsely connected networks.
19
20 296 This trade-off between correlation and discrimination is a fundamental challenge when trying to
21 identify biologically relevant biomarkers that are also clinically relevant [31]. DIABLO controls
22 this trade-off by incorporating *a priori* relationships between different omic domains to
23 adequately model dysregulated biological mechanisms between phenotypic conditions. This may
24 explain the superior biological enrichment of the DIABLO_full models in our benchmarking
25 experiments where the mRNA and miRNA expression as well as methylation activity were
26 assumed to be correlated (**Table 2**). Since these omic domains are known to form real regulatory
27 relationships in order to control complex biological processes, these multi-omic biomarker
28 panels may be capturing this biological complexity. In contrast, these biomarkers were not
29 uncovered when no association was assumed between omic datasets, as in the case of the
30 DIABLO_null models and existing multi-step integrative strategies. Therefore, by controlling the
31 trade-off between correlation and discrimination, DIABLO uncovered novel multi-omics
32 biomarkers that have not previously been identified using existing integrative strategies. These
33 novel biomarkers were part of densely connected clusters of omic variables which have prior
34 known biological associations, further suggesting their potential biological plausibility.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 312 There are areas of improvement that DIABLO will benefit from in the near future. The
5
6 assumption of linear relationship between the selected omics features to explain the phenotypic
7
8 response may not apply in some biological research areas, for example when integrating
9
10 distance-based metagenomics studies, where kernel approaches could be further explored [32].
11
12 315 Selecting the optimal number of variables requires repeated cross-validation to ensure unbiased
13
14 316 classification error rate evaluation. A grid approach was deemed reasonable and provided very
15
16 317 good performance results, but several iterations to refine the grid may be required depending on
17
18 318 the complexity of the classification problem. The grid search algorithm was recently improved
19
20 319 [12], but we advise using a broad filtering strategy to alleviate computational time when dealing
21
22 320 with extremely large datasets (e.g. > 50,000 features each). DIABLO was primarily developed
23
24 321 for omics-measurements on a continuous scale after normalization, and further developments are
25
26 322 needed for categorical data types, such as genotype data, as mentioned in [12]. Finally,
27
28 323 DIABLO, like other methods we benchmarked, will be affected by technical artifacts of the data,
29
30 324 such as batch effects and presence of confounding variables that may affect downstream
31
32 325 integrative analyses. Therefore, we recommend exploratory analyses be carried out in each single
33
34 326 omics dataset to assess the effect, if any, of technical factors and use of batch removal methods
35
36 327 prior to the integration analysis [33–35].
37
38
39
40
41
42
43
44

45 329 To summarize, DIABLO is a versatile, component-based method that can integrate
46
47 330 multiple high dimensional datasets and identify key variables that discriminate between
48
49 331 phenotypic groups. DIABLO identified more biologically relevant and tightly correlated features
50
51 332 across datasets when compared to existing multi-step classification schemes and integrative
52
53 333 methods. The framework is highly flexible, suitable for single point or repeated measures study
54
55 334 designs, and can accommodate various data transformations, such as feature summarization at
56
57
58
59
60
61
62
63
64
65

1
2
3
4 335 the pathway level to enhance biological interpretability. DIABLO's implementation includes
5
6 336 intuitive graphical outputs to facilitate the interpretation of integrative analyses.
7
8 337
9
10
11 338
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 339 **Online Methods**
5
6 340 **Code availability and software tool requirements.** The DIABLO framework is implemented in
7
8
9 341 the mixOmics R package [12]. mixOmics currently includes 19 multivariate methodologies, for
10
11 342 single-omics and integrative analyses. All scripts and tutorials are provided in our companion
12
13
14 343 web-page <http://www.mixomics.org/mixDIABLO>. All analyses were performed using the R
15
16 344 statistical computing program (version 3.4.1) and the mixOmics package (version 6.3.0).
17
18
19 345
20
21 346 **Statistical methods and analysis**
22
23
24 347 ***General multivariate framework to integrate multiple datasets measured on the same samples.***
25
26 348 DIABLO extends sparse generalized canonical correlation analysis (sGCCA) [15] to a
27
28 349 classification (supervised) framework. sGCCA is a multivariate dimension reduction technique
30
31 350 that uses singular value decomposition and selects co-expressed (correlated) variables from
32
33 351 several omics datasets in a computationally and statistically efficient manner. sGCCA maximizes
35
36 352 the covariance between linear combinations of variables (latent component scores) and projects
37
38 353 the data into the smaller dimensional subspace spanned by the components. The selection of the
39
40 354 correlated molecules across omics levels is performed internally in sGCCA with l_1 -penalization
42
43 355 on the variable coefficient vector defining the linear combinations. *Note that since all latent*
44
45 356 *components are scaled in the algorithm, sGCCA maximizes the correlation between components.*
47
48 357 *However, we will retain the term ‘covariance’ instead of ‘correlation’ throughout this section to*
49
50 358 *present the general sGCCA framework.*
52
53 359 Denote K normalized, centered and scaled datasets $X_1 (n \times p_1), \dots, X_K (n \times p_K)$, measuring the
54
55 360 expression levels of p_1, p_2, \dots, p_K omics variables on the same n samples, $k = 1, \dots, K$, sGCCA
57
58 361 solves the optimization function:
59
60
61
62
63
64
65

$$\max_{\mathbf{a}^1, \dots, \mathbf{a}^K} \sum_{k,j=1, k \neq j}^K c_{jk} \text{cov}(X_k \mathbf{a}^k, X_j \mathbf{a}^j), \quad \text{s.t. } \left\| \mathbf{a}^k \right\|_2 = 1 \text{ and } \left\| \mathbf{a}^k \right\|_1 < \lambda_k \quad (1)$$

where c_{jk} indicates whether to maximize the covariance between the datasets X_k and X_j according to the design matrix, with c_{jk} values ranging from 0 (no relationship modelled between the datasets) to 1 otherwise, \mathbf{a}^k is the variable coefficient vector for each dataset X_k , λ_k is a non-negative parameter that controls the amount of shrinkage and thus the number of non-zero coefficients in \mathbf{a}^k . Similar to Lasso [36] and other l_1 – penalized multivariate models developed for single omics analysis [14], the l_1 penalization improves the interpretability of the component scores $X_k \mathbf{a}^k$ that is now only defined on a subset of variables with non-zero coefficients in X_k . The result is the identification of variables that are highly correlated between and within omics datasets.

Equation (1) describes the sGCCA model for the first dimension. Once the first set of coefficient vectors \mathbf{a}_1^k and associated component scores $\mathbf{t}_1^k = X_k \mathbf{a}_1^k$ are obtained, residual matrices are calculated during the ‘deflation’ step for the second dimension, such that $X_k^2 = X_k^1 - \mathbf{t}_1^k \mathbf{a}_1^k$, where X_k^1 is the original centered and scaled data matrix. The subsequent set of components scores and coefficient vectors are then obtained by substituting X_k by X_k^2 in (1). This process is repeated until a sufficient number of dimensions (or set of components) is obtained.

The underlying assumption of the sGCCA model is that the major source of common biological variation can be extracted via the first sets of component scores $X_k \mathbf{a}^k$, while any unwanted variation due to heterogeneity across the datasets X_K does not impact the statistical model. The optimization problem (1) is solved using a monotonically convergent algorithm [15].

1
 2
 3
 4 384
 5
 6
 7 385 ***DIABLO for supervised analysis and prediction.***
 8
 9 386 To extend sGCCA for a classification framework, we substitute one omics dataset X_k in (1) with
 10
 11 387 a dummy indicator matrix Y of size ($n \times G$), where G is the number of phenotype groups that
 12
 13 388 indicate the class membership of each sample. In addition, and for easier use of the method, we
 14
 15 389 replaced the l_1 penalty parameter λ_k by the number of variables to select in each dataset and each
 16
 17 390 component, as there is a direct correspondence between both parameters.
 18
 19
 20
 21 391 Denote a new sample i which is measured across the different types of omics datasets x_k^i ,
 22
 23
 24 392 its class membership is predicted by the fitted sGCCA model with the estimated variable
 25
 26 393 coefficients vectors $\widehat{\alpha}^k$ to obtain the predicted scores $t^{k,i} = x_k^i \widehat{\alpha}^k$, $k = 1, \dots, K$. Therefore, to
 27
 28
 29 394 each dataset k corresponds a predicted continuous score $t^{k,i}$. The predicted class of sample i for
 30
 31 395 each dataset is obtained from the predicted score using one of the distances Maximum, Centroids
 32
 33
 34 396 or Mahalanobis [37] as described in [12]. The consensus class membership is determined using
 35
 36 397 either a majority vote, or by averaging all $t^{k,i}$ across all K datasets before using the prediction
 37
 38
 39 398 distance of choice ('average prediction' scheme). In case of ties in the majority vote scheme,
 40
 41 399 'NA' is allocated as a prediction but is counted as a misclassification error during the
 42
 43
 44 400 performance evaluation. As the class prediction relies on individual vote from each omics set,
 45
 46 401 DIABLO allows for some missing datasets X_k during the prediction step, as illustrated in the
 47
 48
 49 402 Breast Cancer case study. We used the centroid distance for the weighted majority vote scheme
 50
 51 403 (breast cancer study) and the maximum distance for the average vote scheme (asthma study) as
 52
 53 404 those led to best performance (see [12] for details about distance measures and voting schemes
 54
 55
 56 405 that can be used).
 57
 58 406
 59
 60
 61
 62
 63
 64
 65

1
2
3
4 407 ***Design matrix in DIABLO.*** The design matrix C is a ($K \times K$) matrix with values ranging from 0
5 to 1 which specifies whether the covariance between two datasets should be maximized
6
7 408 to 1 which specifies whether the covariance between two datasets should be maximized
8
9 409 DIABLO (see equation (1)). In our simulation study, we evaluated two scenarios: a null design
10
11 410 (DIABLO_null) when no omics datasets are connected, and a full design when all datasets are
12
13
14 411 connected (DIABLO_full):
15
16
17
18
19
20 412
$$C_{\text{null}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad C_{\text{full}} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

21
22
23 413 However, every dataset is connected to the outcome Y internally in the method. For the two case
24
25 414 studies (breast cancer and asthma) the design matrix was chosen based on our proposed method
26
27
28 415 (see below ***Parameters tuning***). Note that the design matrix is not restricted to 0 and 1 values
29
30
31 416 only and a compromise between correlation and discrimination can also be modelled as
32
33 417 described in [12].
34
35 418
36
37
38 419 ***Input data in DIABLO.*** While DIABLO does not assume particular data distributions, all
39
40 420 datasets should be normalized appropriately according to each omics platform and preprocessed
41
42
43 421 if necessary (see normalization steps described below for each case study). Samples should be
44
45
46 422 represented in rows in the data matrices and match the same sample across omics datasets. The
47
48
49 423 phenotype outcome Y is a factor indicating the class membership of each sample. The R function
50
51
52 424 in mixOmics will internally center and scale each variable as is conventionally performed in
53
54 425 PLS-based models and will create the dummy matrix outcome from Y . A multilevel variance
55
56
57 426 decomposition option is available for repeated measures study designs (see below).
58
59
60 428 ***Parameters tuning.***
61
62
63
64
65

1
2
3
4 429 The first parameter to tune is the design matrix C, which can be determined using either prior
5 biological knowledge, or a data-driven approach. The latter approach uses PLS method
6
7 430 implemented in mixOmics that models pair-wise associations between omics datasets. If the
8 correlation between the first component of each omics dataset is above a given threshold (e.g.
9 431 0.8) then a connection between those datasets is included in the DIABLO design as a 1 value.
10
11 432
12
13 433
14
15

16 434 The second parameter to tune is the total number of components. In several analyses we
17 found that $G - 1$ components were sufficient to extract sufficient information to discriminate all
18 phenotype groups [14], but this can be assessed by evaluating the model performance across all
19 specified components (described below) as well as using graphical outputs such as sample plots
20
21 436
22
23 437
24
25 438 to visualize the discriminatory ability of each component.
26
27

28 439 Finally, the third set of parameters to tune is the number of variables to select per dataset
29 and per component. Such tuning can rapidly become cumbersome, as there might be numerous
30 combinations of selection sizes to evaluate across all K datasets. For the breast cancer study, we
31 440 used 5-fold cross-validation repeated 50 times to evaluate the performance of the model over a
32 grid of different possible values of variables to select (**Supplementary Fig. 8**). The performance
33 441 of the model for a given set of parameters (including number of component and number of
34
35 442 variables to select) was based on the balanced classification error rate using majority vote or
36
37 443 average prediction schemes with centroids distance. The balanced classification error rate is
38
39 444 useful in the case of imbalanced class sizes, where the majority classes can have strong influence
40
41 445 on the overall error rate. The balanced error rate measure calculates the weighted average of the
42
43 446 individual class error rates with respect to their class sample size. In our experience, the number
44
45 447 of variables to select in each dataset provided less of an improvement on the error rate compared
46
47 448 to tuning the number of components. Therefore, even a grid composed of a small number of
48
49 449
50
51 448
52
53 449
54
55 450
56
57 448
58 451 to tuning the number of components. Therefore, even a grid composed of a small number of
59
60
61
62
63
64
65

1
2
3
4 452 variables (<50 with steps of 5 or 10) may suffice as it does not substantially change the
5
6 classification performance. This is because of the use of regularization constraints which reduces
7
8 the variability in the variable coefficients and thus maintains the predictive ability of the model.
9
10
11 455 Further, the variable selection size can also be guided according to the downstream biological
12
13 interpretation to be performed. For example, a gene-set enrichment analysis may require a larger
14
15 set of features than a literature-search interpretation.
16
17
18
19 458
20
21 459 ***Visualization outputs with DIABLO.*** To facilitate the interpretation of the integrative analysis,
22
23 460 several types of graphical outputs were implemented in mixOmics.
24
25
26 461 *Sample plots.* The consensus plot which depicts the samples is computed by calculating the
27
28 average of the components from each dataset. Omics specific samples plots can also be obtained
29
30 by plotting components associated to each data set. The sample plot are useful to visualize the
31
32 ability of the DIABLO model to extract common information at the sample level for each
33
34 dataset, and the discriminatory power of each data type to separate the phenotypic groups. The
35
36 scatterplot matrix represents the correlation between components for the same dimension but
37
38 across all omics datasets. This plot assesses the model's ability to maximize the correlation as
39
40 indicated in the design matrix. Separation of subjects according to their phenotypic groups can
41
42
43 468 be visualized.
44
45
46
47
48 470 *Variable plots.* To visualize selected variables, we proposed circos plot to represent correlations
49
50 between and within variables from each dataset at the variable level. The association between
51
52 variables is computed using a similarity score that is analogous to a Pearson correlation
53
54 coefficient, as previously described in [38]. For each omics dataset, DIABLO produces a
55
56 variable coefficient matrix of size ($p_k \times H$), where H is the total number of components in the
57
58
59
60
61
62
63
64
65

1
2
3
4 475 model. The product of any two matrices approximates the association score between variables of
5
6 the two omics datasets. The association between variables is displayed as a color-coded link
7
8 inside the plot to represent a positive or negative correlation above a user-specified threshold.
9
10 477 The selected variables are represented on the side of the circos plot, with side colors indicating
11
12 each omics type, optional line plots represent the expression levels in each phenotypic group.
13
14 479
15
16 480 *Clustered Image Map (CIM)*. A clustered image map [38] based on the Euclidean distance and
17
18 the complete linkage displays an unsupervised clustering between the selected variables
19
20 (centered and scaled) and the samples. Color bars represent the sample phenotypic groups
21
22 (columns) and the type of omics (rows) variables.
23
24 483
25
26 484
27
28 485 **Gene-set enrichment analyses**
29
30
31 486 Significance of enrichment was determined using a hypergeometric test of the overlap between
32
33 the selected features (mapped to official HUGO gene symbols or official miRNA symbols) and
34
35 the various gene sets contained in the collections. In order to carry out the comparison, each
36
37 feature set was mapped back to official HUGO gene symbols. This was done as follows across
38
39 the respective data types: mRNA, CpGs and proteins. The following collections were used as
40
41 gene-sets for the enrichment analysis [39]: C1 - positional gene sets for each human chromosome
42
43 and cytogenetic band. C2 – curated gene sets (Pathway Interaction DB [PID], Biocarta
44
45 [BIOCARTA], Kyoto Encyclopedia of Genes and Genomes [KEGG], Reactome [REACTOME],
46
47 and others), C3 - motif gene sets based on conserved cis-regulatory motifs from a comparative
48
49 analysis of the human, mouse, rat, and dog genomes. C4 – computational gene sets (from the
50
51 Cancer Gene Neighbourhoods [CGN] and Cancer Modules [CM] – citation available via the
52
53 MolSigDB [23]. C5 - GO gene sets consist of genes annotated by the same GO terms. C6 –
54
55
56
57
58 497
59
60
61
62
63
64
65

1
2
3
4 498 ontologic gene sets (Gene sets represent signatures of cellular pathways which are often dis-
5
6 regulated in cancer). C7 - immunologic gene sets defined directly from microarray gene
7
8 expression data from immunologic studies. H - hallmark gene sets are coherently expressed
9 500 signatures derived by aggregating many MSigDB gene sets to represent well-defined biological
10
11 501 states or processes. & A. BTM - Blood Transcriptional Modules [40]. B. TISSUES - cell-specific
12
13 502 states or processes. & A. BTM - Blood Transcriptional Modules [40]. B. TISSUES - cell-specific
14
15 503 expression from Benita *et al.* [41].
16
17 504
18
19 505

20
21 505 **Modular analysis:** Eigengene summarization is a common approach to decompose a $n \times p$
22
23 dataset (where n is the number of samples and p is the number of variables in a module), to a
24
25 component (linear combination of all p variables) that represents the summarized expression of
26
27 genes in the module [18]. For the asthma study, 15,683 genes were reduced to 229 KEGG
28
29 pathways and 292 metabolites were reduced to 60 metabolic pathways using eigengene
30
31 509 summarization.
32
33 510
34
35
36 511
37
38 512 **Multilevel transformation:** For multivariate analyses, A multilevel approach separates the
39
40 within subject variation matrix (X_w) and the between subject variation (X_b) for a given dataset (X)
41
42 513 [42], ie. $X = X_w + X_b$. In the case of a two-repeated measured problem (e.g. pre vs post
43
44 challenge), the within subject variation matrix is similar to calculating the net difference for each
45
46 515 challenge, the within subject variation matrix is similar to calculating the net difference for each
47
48 516 individual between the data obtained for pre and post challenge. For each omics dataset, the
49
50 517 within-subject variation matrix was extracted prior to applying DIABLO. In the asthma study,
51
52
53 518 the multilevel approach (called variance decomposition step) was applied to the cell-type, gene
54
55 519 and metabolite module datasets.
56
57
58 520
59
60
61
62
63
64
65

1
2
3
4 521 **Declarations**
5
6
7 522• **Acknowledgements**
8
9 523• The authors would like to thank Mr Kevin Chang (University of Auckland) for some preliminary
10
11 524 exploratory analyses of the breast cancer dataset. We would also like to thank Dr Chao Liu
12
13
14 525 (University of Queensland) for obtaining the PAM50 phenotypic information for the TCGA
15
16 526 datasets.
17
18
19 527
20
21 528• **Competing interests**
22
23
24 529 The authors declare no competing interests.
25
26 530
27
28
29 531• **Funding**
30
31 532 AS is the recipient of the Canadian Institutes of Health Research Doctoral Award – Frederick
32
33 533 Banting and Charles Best Canada Graduate Scholarship and the Michael Smith Foreign Study
34
35
36 534 Supplement award. Research reported in this publication was supported in part by the National
37
38 535 Institute Of Allergy And Infectious Diseases of the National Institutes of Health under Award
39
40
41 536 Number U19AI118608 (CPS and SJT). The content is solely the responsibility of the authors and
42
43 does not necessarily represent the official views of the National Institutes of Health. KALC is
44
45
46 538 supported in part by the National Health and Medical Research Council (NHMRC) Career
47
48 539 Development fellowship (GNT1087415).
49
50
51 540•
52
53 541• **Authors' contributions**
54
55 542• AS performed the data pre-processing, the statistical analyses and developed the DIABLO
56
57
58 543 method. BG implemented the R scripts for DIABLO and graphical outputs, CPS performed the
59
60
61
62
63
64
65

1
2
3
4 544 gene enrichment analyses, MV implemented the circos plots, FR and BG implemented the R
5
6 scripts in mixOmics along with the S3 functions, SJT supervised AS and participated in the
7
8 design of the study. KALC supervised AS, BG, MV and FR, participated in the development of
9
10
11 547 the DIABLO method and provided statistical advice. AS and KALC edited the manuscript, with
12
13
14 548 editorial input from SJT and CPS.
15
16 549
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 550 **Tables**
5
6

7 551 **Table 1. Overview of multi-omics datasets analyzed for method benchmarking and in two**
8 552 **case studies.** The breast cancer case study includes training and test datasets for all omics types
9 553 except proteins.

Analysis	Dataset	Number of samples	Sample size in each subtype			Omics	Number of variables
Benchmark cancer datasets (Wang et al. [3])	Colon	92	High (33) Low (59)			mRNA	17,814
						miRNA	312
						CpGs	23,088
	Kidney	122	High (61) Low (61)			mRNA	17,665
						miRNA	329
						CpGs	24,960
	Glioblastoma (gbm)	213	High (105) Low (108)			mRNA	12,042
						miRNA	534
						CpGs	1,305
	Lung	106	High (53) Low (53)			mRNA	12,042
						miRNA	353
						CpGs	23,074
Case study 1 (The Cancer Genome Atlas) [19]	Breast cancer	989		Train	Test	mRNA	16,851
			Basal	76	102	miRNA	349
			Her2	38	40	CpGs	9,482
			LumA	188	346	Proteins	Train: 115 Test: 0
			LumB	77	122		
Case study 2 (Singh et al. [20,21])	Asthma	28	Pre (14) Post (14)			Cell-types	9
						mRNA-modules	229
						metabolite-modules	60

45 554
46
47 555 **Table 2. Number of significant gene sets for each integrative method and benchmarking**
48 556 **cancer dataset.** Best performing method is indicated in the shaded cell. Each row represents a
49 557 gene set collection (see Methods for details, FDR = 5%).

		Unsupervised, integrative			Supervised, non-integrative			Supervised, integrative
disease	collection	JIVE	MOFA	sGCCA	Concatenation	Ensemble	DIABLO_null	DIABLO_full
Colon	BTM	0	4	0	0	0	0	23
	C1	0	0	0	0	0	0	0
	C2	15	14	5	12	3	21	113
	C3	8	5	14	11	2	6	0

	C4	0	1	0	1	2	1	46
	C5	19	36	147	7	0	0	216
	C6	0	0	0	0	0	0	0
	C7	1	87	11	61	10	62	218
	H	0	0	0	0	0	2	7
	TISSUES	2	12	0	0	0	0	16
	TOTAL	45	159	177	92	17	92	639
	Gbm	BTM	0	0	19	10	9	30
		C1	0	0	0	0	0	0
		C2	275	337	193	258	358	312
		C3	94	64	37	14	15	15
		C4	49	43	68	47	50	62
		C5	825	708	706	526	669	776
		C6	22	25	18	30	24	21
		C7	460	82	526	432	173	147
		H	12	8	8	19	23	20
		TISSUES	18	29	21	10	12	14
		TOTAL	1755	1296	1596	1346	1333	1380
	Kidney	BTM	1	0	0	0	0	0
		C1	0	0	1	0	0	1
		C2	42	33	7	10	5	15
		C3	8	80	1	4	35	23
		C4	17	6	0	7	1	3
		C5	157	110	1	55	27	46
		C6	0	0	0	0	0	0
		C7	0	74	15	93	13	10
		H	6	3	0	1	0	1
		TISSUE S	2	0	0	0	0	0
		TOTAL	233	306	25	170	81	98
	Lung	BTM	0	0	0	0	2	0
		C1	0	0	0	1	0	1
		C2	4	17	2	0	0	1
		C3	48	20	57	50	26	21
		C4	17	0	47	0	0	18
		C5	35	127	42	0	25	22
		C6	1	0	1	3	2	5
		C7	18	13	78	0	7	72
		H	0	2	0	0	1	0

	TISSUE S	0	0	0	0	0	9	20
	TOTAL	123	179	227	54	61	150	386

8 558

9

10 559

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

1
2
3
4 560 **Figure captions**
5
6
7 561

8 562 **Figure 1. Simulation study: performance assessment and benchmarking.** Simulated datasets
9 563 included different types of variables: correlated & discriminatory (corDis); uncorrelated &
10 564 discriminatory (unCorDis); correlated & nondiscriminatory (corNonDis) and uncorrelated &
11 565 nondiscriminatory (unCorNonDis) for different fold-changes between sample groups and
12 566 different noise levels (see **Supplementary Note**). Integrative classifiers included DIABLO with
13 567 either the full or null design, concatenation and ensemble-based sPLSDA classifiers and were all
14 568 trained to select 90 variables across three multi-omics datasets. **a)** Classification error rates (10-
15 569 fold cross-validation averaged over 50 simulations). Dashed line indicates a random performance
16 570 (error rate = 50%). All methods perform similarly with the exception of DIABLO_full which has
17 571 a higher error rate. **b)** Number of variables selected according to their type. DIABLO_full
18 572 selected mainly variables that were correlated & discriminatory (corDis, red), whereas the other
19 573 methods selected an equal number of correlated or uncorrelated discriminatory variables (corDis
20 574 and unCorDis, red and blue).

21 575
22 576 **Figure 2. Benchmarking integrative methods using multi-omics biomarker panels for**
23 577 **different cancers.** **a)** Overlap of selected features using both supervised (green) and
24 578 unsupervised approaches (purple): a strong overlap was observed between the supervised
25 579 approaches with the exception of DIABLO_full (blue bars) which showed more similarity to
26 580 unsupervised methods (dark orange bars). **b)** Number of edges within each panel network at
27 581 various Pearson correlation cut-offs: unsupervised approaches panels were more connected than
28 582 those from supervised approaches, with the exception of DIABLO_full which led to a highly-
29 583 connected panel. An edge is present if the association between two omic variables is greater than
30 584 a given correlation cut-off. **c)** Upper panel: network modularity of each multi-omic biomarker
31 585 panel for colon cancer showed that unsupervised approaches and DIABLO_full resulted in a few
32 586 groups of highly connected features, whereas supervised approaches identified networks with
33 587 many groups of sparsely connected features. Lower panel: component plots depicting the clear
34 588 separation of subjects in the high and low survival groups for supervised methods as opposed to
35 589 the unsupervised methods.

36 590
37 591 **Figure 3. Identification of a multi-omics biomarker panel predictive of breast cancer**
38 592 **subtypes.** **a)** Variable contributions of each omics-type biomarker that are important to
39 593 discriminate breast cancer subtypes. **b)** DIABLO component plots and the derived biomarker
40 594 panel: 95% confidence ellipses were calculated from the training data set and points depict
41 595 samples from the test set. **c)** Heatmap of the scaled expression of variable from the biomarker
42 596 panel. **d)** Network visualization of the biomarker panel highlights correlated variables (Pearson
43 597 correlation $> |0.4|$) and four communities based on edge betweenness scores. **e)** A gene set
44 598 enrichment analysis was conducted on the largest community from d (red cluster) where many
45 599 cancer related pathways were identified.

46 600
47 601 **Figure 4. Asthma study: cross-over design and module-based analysis with DIABLO.**
48 602 **a)** DIABLO design includes a module-based decomposition approach to discriminate pre-and
49 603 post-inhalation challenge samples. **b)** Receiver operating characteristic curves comparing the
50 604 performance of the standard DIABLO and ‘multilevel DIABLO’ for repeated measures

1
2
3
4 605 (mDIABLO) using leave-one-out cross-validation. **c**) Component plots depicting the separation
5 606 of the pre- and post-challenge samples based on DIABLO and mDIABLO. **d**) Overlapping
6 607 features selected from either DIABLO or mDIABLO. **e**) Heatmap of the Pearson correlation
7 608 values between the features selected with mDIABLO. **f**) Circos plot depicting the strongest
8 609 correlations between different omics features from the mDIABLO panel.
9 610
10 611
11 612
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

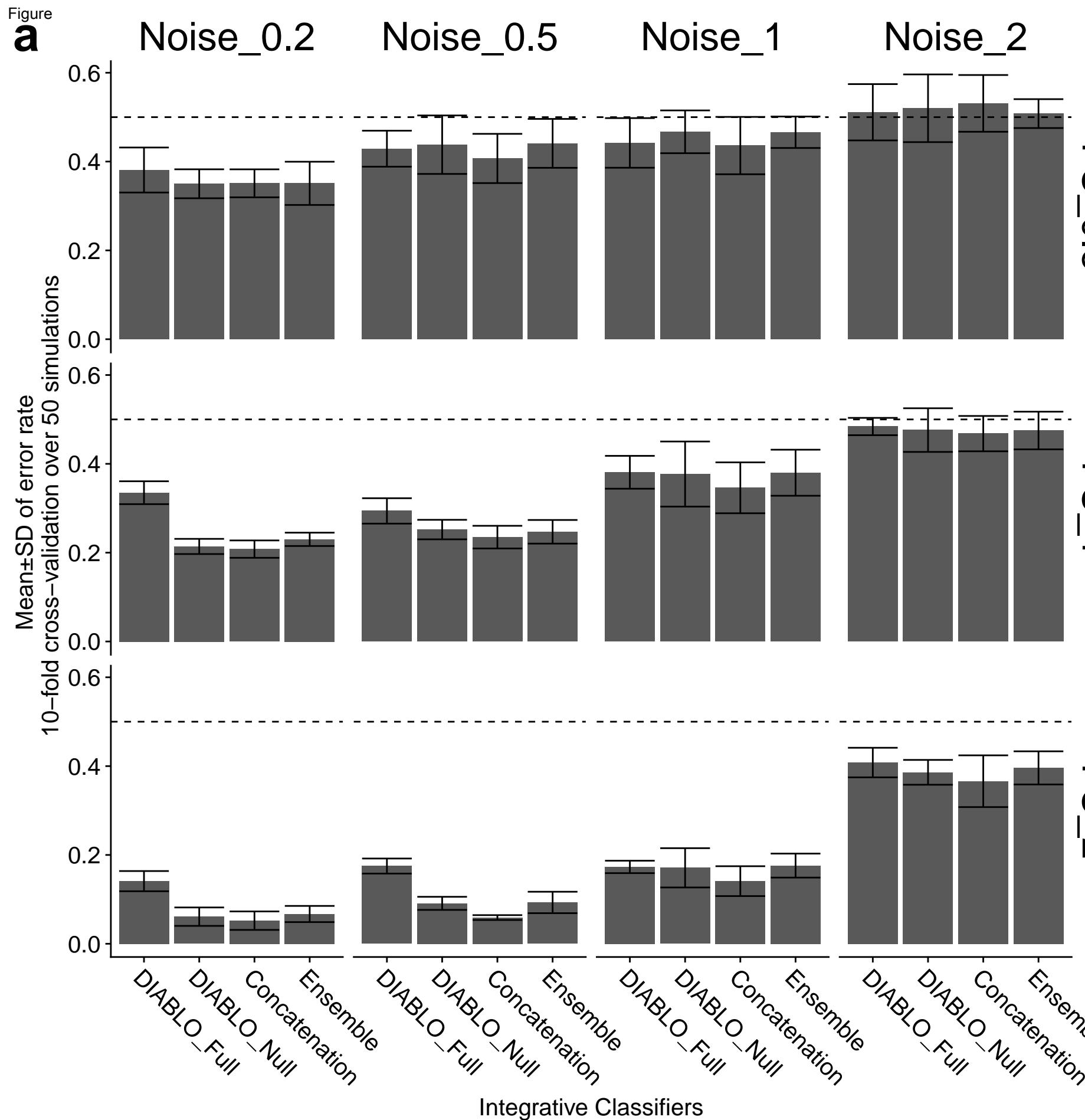
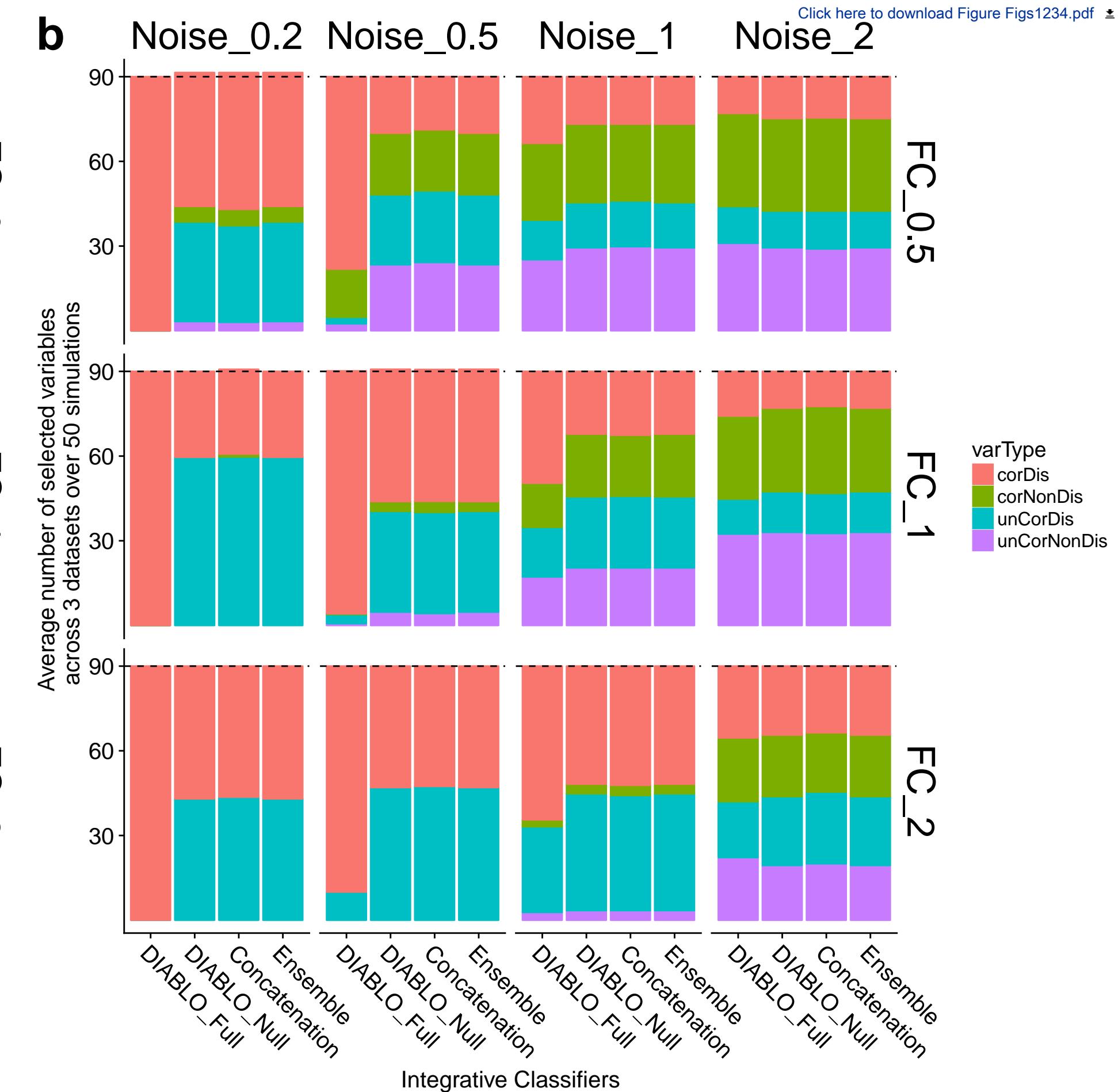
1
2
3
4
5 613 **References**
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

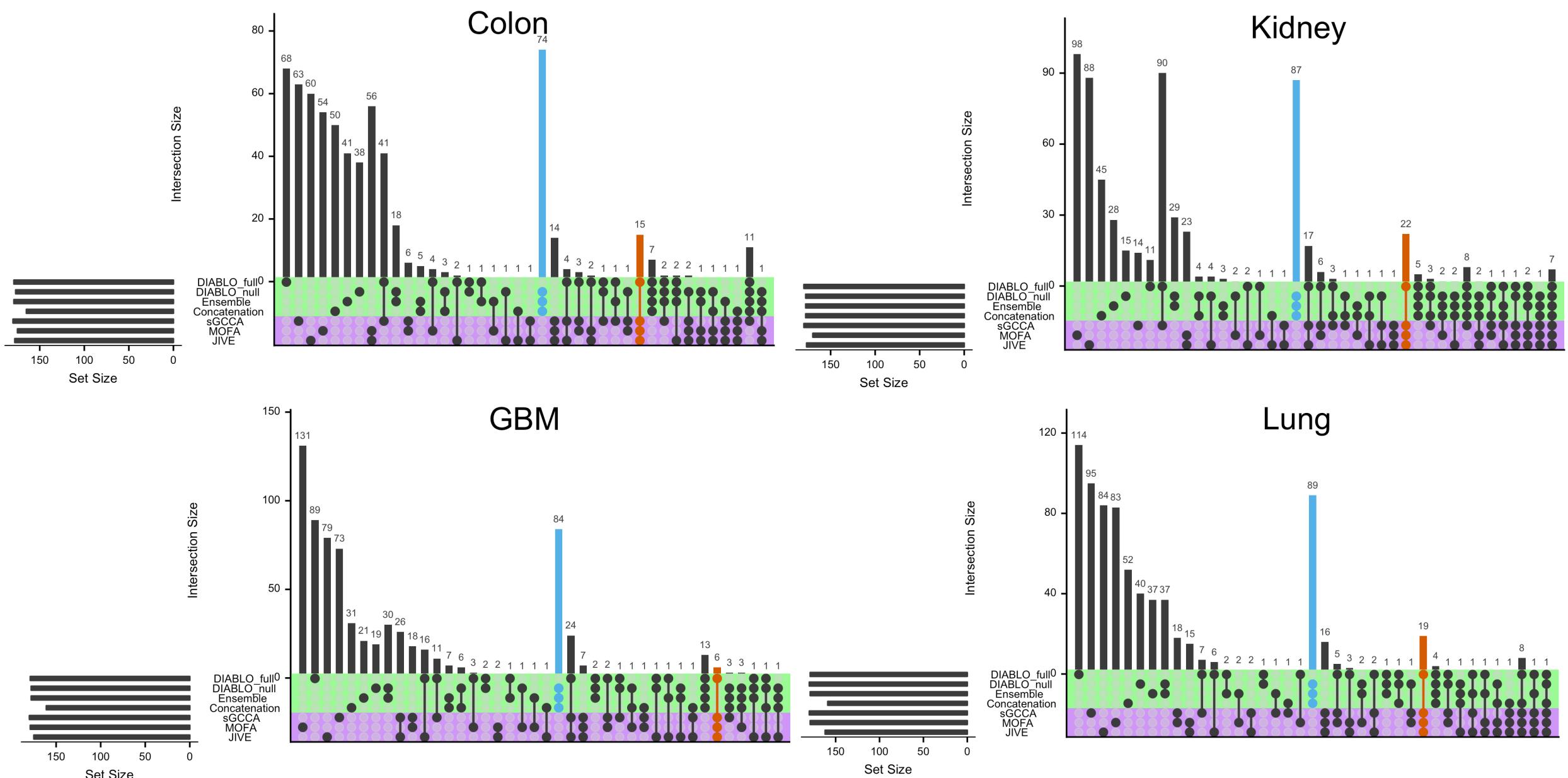
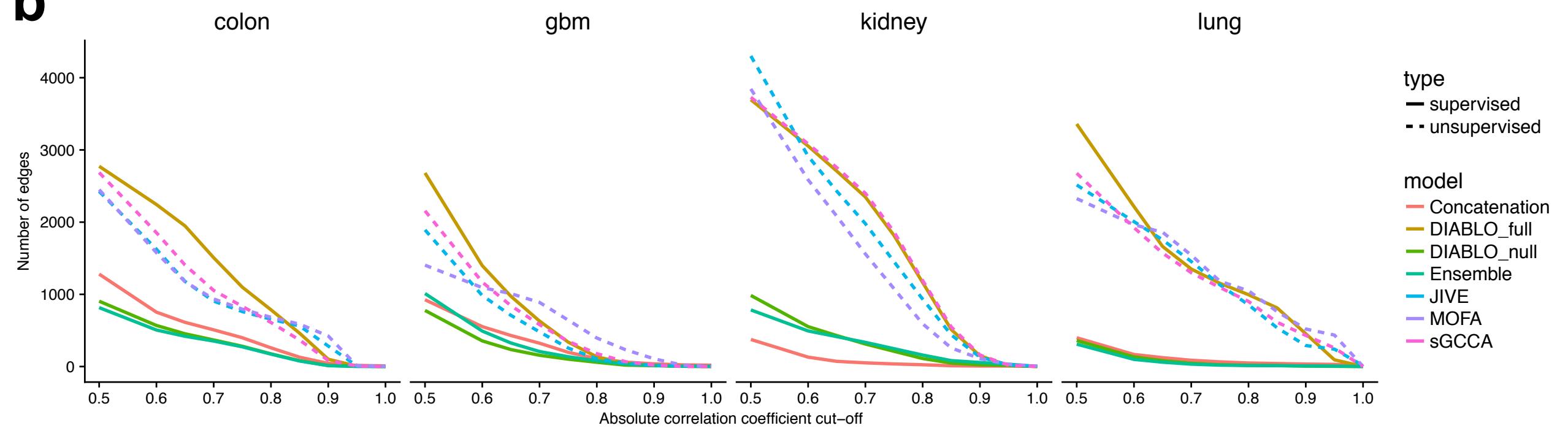
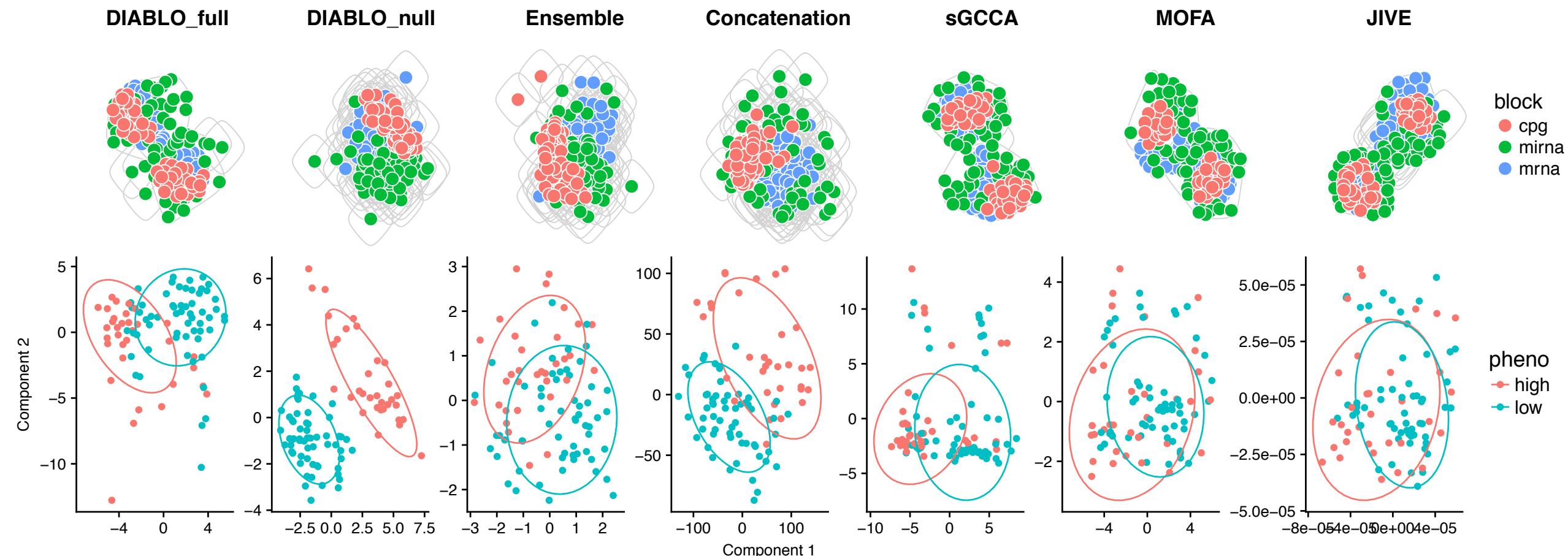
1. Zhu J, Sova P, Xu Q, Dombek KM, Xu EY, Vu H, et al. Stitching together multiple data dimensions reveals interacting metabolomic and transcriptomic networks that modulate cell regulation. Levchenko A, editor. *PLoS Biol* [Internet]. 2012 [cited 2016 Jan 19];10:e1001301. Available from: <http://dx.plos.org/10.1371/journal.pbio.1001301>
2. Kim D, Li R, Dudek SM, Ritchie MD. ATHENA: Identifying interactions between different levels of genomic data associated with cancer clinical outcomes using grammatical evolution neural network. *BioData Min*. 2013;6:23.
3. Wang B, Mezlini AM, Demir F, Fiume M, Tu Z, Brudno M, et al. Similarity network fusion for aggregating data types on a genomic scale. *Nat Methods* [Internet]. 2014 [cited 2016 Jan 19];11:333–7. Available from: <http://www.nature.com/doifinder/10.1038/nmeth.2810>
4. Ritchie MD, Holzinger ER, Li R, Pendergrass SA, Kim D. Methods of integrating data to uncover genotype–phenotype interactions. *Nat Rev Genet* [Internet]. 2015 [cited 2015 Jul 10];16:85–97. Available from: <http://www.nature.com/doifinder/10.1038/nrg3868>
5. Yugi K, Kubota H, Hatano A, Kuroda S. Trans-Omics: How To Reconstruct Biochemical Networks Across Multiple ‘Omic’ Layers. *Trends Biotechnol* [Internet]. 2016 [cited 2018 Feb 21];34:276–90. Available from: <http://linkinghub.elsevier.com/retrieve/pii/S0167779915002735>
6. Günther O, Chen V, Freue GC, Balshaw R, Tebbutt S, Hollander Z, et al. A computational pipeline for the development of multi-marker bio-signature panels and ensemble classifiers. 2012 [cited 2016 Jan 19];13:326. Available from: <http://summit.sfu.ca/item/13303>
7. Aben N, Vis DJ, Michaut M, Wessels LFA. TANDEM: a two-stage approach to maximize interpretability of drug response models based on multiple molecular data types. *Bioinformatics* [Internet]. 2016 [cited 2017 Aug 2];32:i413–20. Available from: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btw449>
8. Ma S, Ren J, Fenyö D. Breast cancer prognostics using multi-omics data. *AMIA Summits Transl Sci Proc* [Internet]. 2016 [cited 2017 May 30];2016:52. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5001766/>
9. Bersanelli M, Mosca E, Remondini D, Giampieri E, Sala C, Castellani G, et al. Methods for the integration of multi-omics data: mathematical aspects. *BMC Bioinformatics* [Internet]. 2016 [cited 2016 May 8];17. Available from: <http://www.biomedcentral.com/1471-2105/17/S2/15>
10. Meng C, Zelezniak OA, Thallinger GG, Kuster B, Gholami AM, Culhane AC. Dimension reduction techniques for the integrative analysis of multi-omics data. *Brief Bioinform* [Internet]. 2016 [cited 2018 Feb 21];17:628–41. Available from: <https://academic.oup.com/bib/article-lookup/doi/10.1093/bib/bbv108>

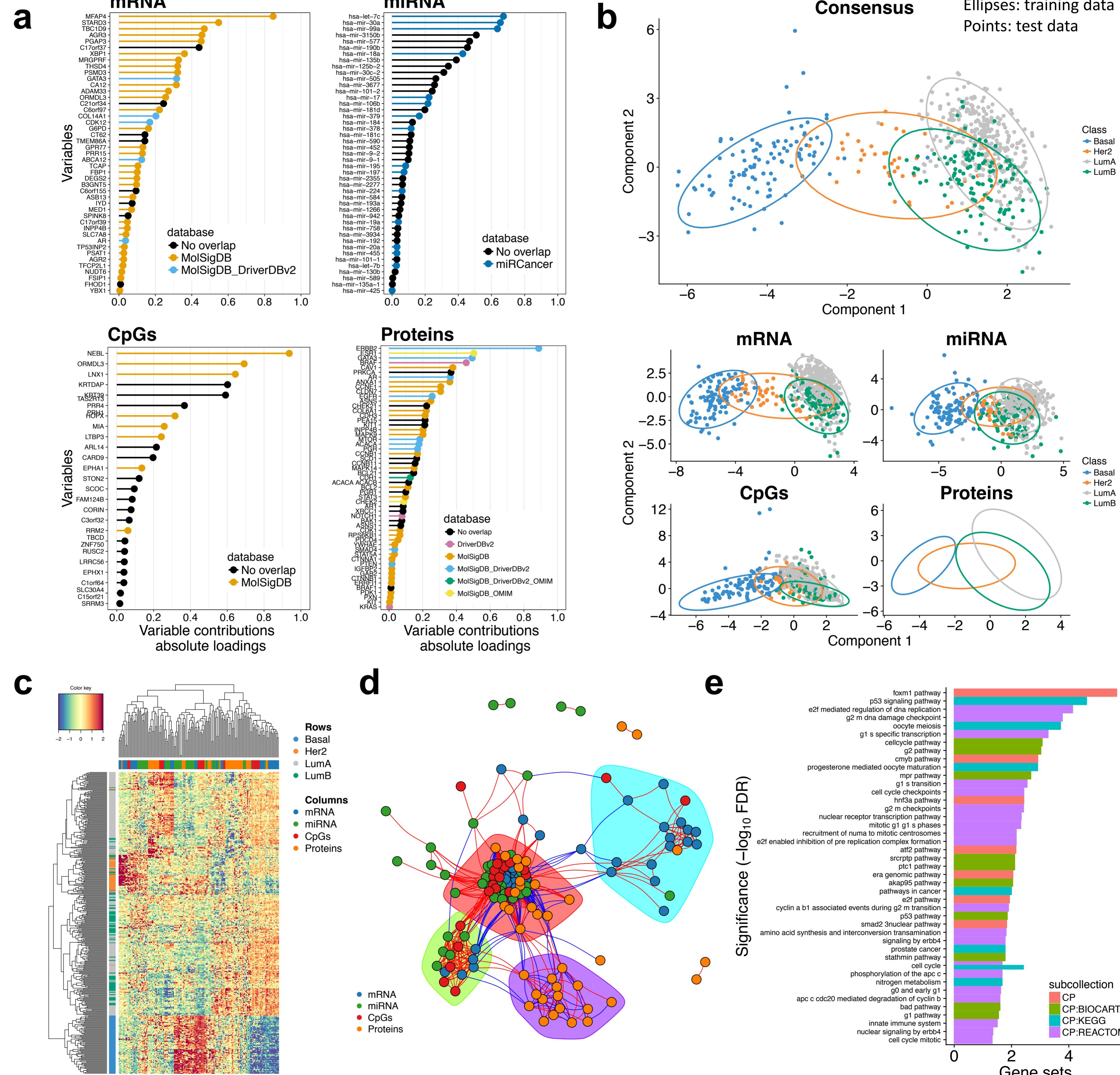
- 1
2
3
- 4 647 11. Huang S, Chaudhary K, Garmire LX. More Is Better: Recent Progress in Multi-Omics Data
5 648 Integration Methods. *Front Genet* [Internet]. 2017 [cited 2018 Feb 21];8. Available from:
6 649 <http://journal.frontiersin.org/article/10.3389/fgene.2017.00084/full>
- 7
8
9 650 12. Rohart F, Gautier B, Singh A, Cao K-AL. mixOmics: An R package for ‘omics feature
10 651 selection and multiple data integration. *PLOS Comput Biol* [Internet]. 2017 [cited 2018 Jan
11 652 29];13:e1005752. Available from:
12 653 <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1005752>
- 13
14
15 654 13. Wold H. Estimation of Principal Components and Related Models by Iterative Least squares.
16 655 *Multivar Anal.* 1966;391–420.
- 17
18 656 14. Lê Cao K-A, Boitard S, Besse P. Sparse PLS discriminant analysis: biologically relevant
19 657 feature selection and graphical displays for multiclass problems. *BMC Bioinformatics* [Internet].
20 658 2011 [cited 2015 Jul 15];12:253. Available from: <http://www.biomedcentral.com/1471-2105/12/253/>
- 21
22 659
23
24 660 15. Tenenhaus A, Philippe C, Guillemot V, Le Cao K-A, Grill J, Frouin V. Variable selection for
25 661 generalized canonical correlation analysis. *Biostatistics* [Internet]. 2014 [cited 2015 Jul
26 662 15];15:569–83. Available from:
27 663 <http://biostatistics.oxfordjournals.org/cgi/doi/10.1093/biostatistics/kxu001>
- 28
29
30 664 16. Witten DM, Tibshirani R, Hastie T. A penalized matrix decomposition, with applications to
31 665 sparse principal components and canonical correlation analysis. *Biostatistics* [Internet]. 2009
32 666 [cited 2016 Jul 27];10:515–34. Available from:
33 667 <http://biostatistics.oxfordjournals.org/cgi/doi/10.1093/biostatistics/kxp008>
- 34
35
36 668 17. Lee HK, Hsu AK, Sajdak J, Qin J, Pavlidis P. Coexpression analysis of human genes across
37 669 many microarray data sets. *Genome Res* [Internet]. 2004 [cited 2016 Mar 30];14:1085–1094.
38 670 Available from: <http://genome.cshlp.org/content/14/6/1085.short>
- 39
40
41 671 18. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis.
42 672 *BMC Bioinformatics* [Internet]. 2008 [cited 2016 Apr 4];9:559. Available from:
43 673 <http://www.biomedcentral.com/1471-2105/9/559>
- 44
45
46 674 19. The TCGA Research Network. The Cancer Genome Atlas [Internet]. Available from:
47 675 <http://cancergenome.nih.gov/>
- 48
49 676 20. Singh A, Yamamoto M, Kam SHY, Ruan J, Gauvreau GM, O’Byrne PM, et al. Gene-
50 677 metabolite expression in blood can discriminate allergen-induced isolated early from dual
51 678 asthmatic responses. Hsu Y-H, editor. *PLoS ONE* [Internet]. 2013 [cited 2015 Jul 18];8:e67907.
52 679 Available from: <http://dx.plos.org/10.1371/journal.pone.0067907>
- 53
54
55 680 21. Singh A, Yamamoto M, Ruan J, Choi JY, Gauvreau GM, Olek S, et al. Th17/Treg ratio
56 681 derived using DNA methylation analysis is associated with the late phase asthmatic response.
57 682 *Allergy Asthma Clin Immunol* [Internet]. 2014 [cited 2016 Mar 2];10:32. Available from:
58 683 <http://www.biomedcentral.com/content/pdf/1710-1492-10-32.pdf>
- 59
60
61
62
63
64
65

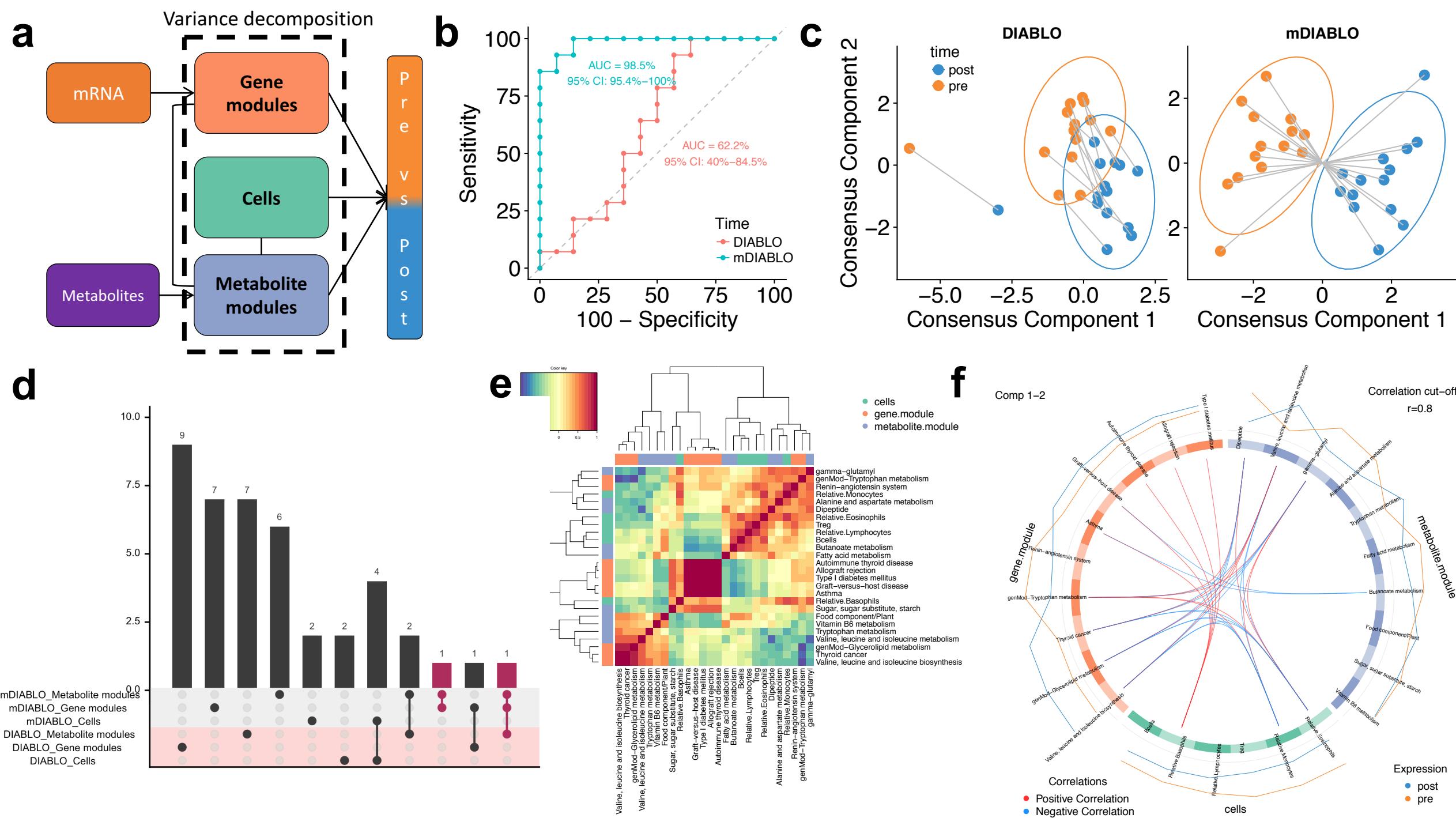
- 1
2
3
- 4 684 22. Lock EF, Hoadley KA, Marron JS, Nobel AB. Joint and individual variation explained
5 685 (JIVE) for integrated analysis of multiple data types. *Ann Appl Stat* [Internet]. 2013 [cited 2018
6 686 Jan 24];7:523–42. Available from: <http://projecteuclid.org/euclid.aoas/1365527209>
- 7
8
9 687 23. Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The Molecular
10 688 Signatures Database Hallmark Gene Set Collection. *Cell Syst* [Internet]. 2015 [cited 2018 Jan
11 689 30];1:417–25. Available from: <http://linkinghub.elsevier.com/retrieve/pii/S2405471215002185>
- 12
13
14 690 24. Xie B, Ding Q, Han H, Wu D. miRCancer: a microRNA-cancer association database
15 691 constructed by text mining on literature. *Bioinformatics* [Internet]. 2013 [cited 2018 Jan
16 692 30];29:638–44. Available from: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btt014>
- 17 693
18
19
20 694 25. Hamosh A. Online Mendelian Inheritance in Man (OMIM), a knowledgebase of human
21 695 genes and genetic disorders. *Nucleic Acids Res* [Internet]. 2004 [cited 2018 Jan 30];33:D514–7.
22 696 Available from: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gki033>
- 23
24 697 26. Chung I-F, Chen C-Y, Su S-C, Li C-Y, Wu K-J, Wang H-W, et al. DriverDBv2: a database
25 698 for human cancer driver gene research. *Nucleic Acids Res* [Internet]. 2016 [cited 2018 Jan
26 699 30];44:D975–9. Available from: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkv1314>
- 27
28 700
29
30 701 27. Liquet B, Lê Cao K-A, Hocini H, Thiébaut R. A novel approach for biomarker selection and
31 702 the integration of repeated measures experiments from two assays. *BMC Bioinformatics*
32 703 [Internet]. 2012 [cited 2015 Jul 18];13:325. Available from:
33 704 <http://www.biomedcentral.com/1471-2105/13/325/>
- 34
35
36 705 28. Allahyar A, de Ridder J. FERAL: network-based classifier with application to breast cancer
37 706 outcome prediction. *Bioinformatics* [Internet]. 2015 [cited 2018 Feb 1];31:i311–9. Available
38 707 from: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btv255>
- 39
40
41 708 29. Cun Y, Fröhlich H. Network and data integration for biomarker signature discovery via
42 709 network smoothed t-statistics. Boccaletti S, editor. *PLoS ONE* [Internet]. 2013 [cited 2017 May
43 710 30];8:e73074. Available from: <http://dx.plos.org/10.1371/journal.pone.0073074>
- 44
45
46 711 30. Sokolov A, Carlin DE, Paull EO, Baertsch R, Stuart JM. Pathway-based genomics prediction
47 712 using generalized elastic net. *PLoS Comput Biol* [Internet]. 2016 [cited 2017 May
48 713 30];12:e1004790. Available from:
49 714 <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004790>
- 50
51
52 715 31. Wang TJ. Assessing the Role of Circulating, Genetic, and Imaging Biomarkers in
53 716 Cardiovascular Risk Prediction. *Circulation* [Internet]. 2011 [cited 2018 Feb 23];123:551–65.
54 717 Available from: <http://circ.ahajournals.org/cgi/doi/10.1161/CIRCULATIONAHA.109.912568>
- 55
56
57 718 32. Mariette J, Villa-Vialaneix N. Unsupervised multiple kernel learning for heterogeneous data
58 719 integration. *Bioinformatics* [Internet]. 2017 [cited 2018 Mar 6]; Available from:
59 720 <http://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/btx682/4565592>
- 60 721
61
62
63
64
65

- 1
2
3
- 4 722 33. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using
5 723 empirical Bayes methods. *Biostatistics* [Internet]. 2007 [cited 2016 May 12];8:118–27. Available
6 724 from: <http://biostatistics.oxfordjournals.org/cgi/doi/10.1093/biostatistics/kxj037>
- 7
8
9 725 34. Gagnon-Bartsch JA, Speed TP. Using control genes to correct for unwanted variation in
10 726 microarray data. *Biostatistics* [Internet]. 2012 [cited 2018 Mar 6];13:539–52. Available from:
11 727 <https://academic.oup.com/biostatistics/article-lookup/doi/10.1093/biostatistics/kxr034>
- 12
13
14 728 35. Parker HS, Corrada Bravo H, Leek JT. Removing batch effects for prediction problems with
15 729 frozen surrogate variable analysis. *PeerJ* [Internet]. 2014 [cited 2016 May 12];2:e561. Available
16 730 from: <https://peerj.com/articles/561>
- 17
18
19 731 36. Tibshirani R. Regression shrinkage and selection via the lasso. *J R Stat Soc Ser B Methodol.*
20 732 1996;58:267–88.
- 21
22
23 733 37. Le Cao K-A, Gonzalez I, Dejean S. integrOmics: an R package to unravel relationships
24 734 between two omics datasets. *Bioinformatics* [Internet]. 2009 [cited 2016 Apr 3];25:2855–6.
25 735 Available from: <http://bioinformatics.oxfordjournals.org/cgi/doi/10.1093/bioinformatics/btp515>
- 26
27 736 38. González I, Lê Cao K-A, Davis MJ, Déjean S. Visualising associations between paired
28 737 ‘omics’ data sets. *BioData Min* [Internet]. 2012 [cited 2015 Jul 15];5:1–23. Available from:
29 738 <http://link.springer.com/article/10.1186/1756-0381-5-19>
- 30
31
32 739 39. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set
33 740 enrichment analysis: a knowledge-based approach for interpreting genome-wide expression
34 741 profiles. *Proc Natl Acad Sci* [Internet]. 2005 [cited 2016 Jul 26];102:15545–15550. Available
35 742 from: <http://www.pnas.org/content/102/43/15545.short>
- 36
37
38 743 40. Chaussabel D, Quinn C, Shen J, Patel P, Glaser C, Baldwin N, et al. A Modular Analysis
39 744 Framework for Blood Genomics Studies: Application to Systemic Lupus Erythematosus.
40 745 Immunity [Internet]. 2008 [cited 2016 Jul 22];29:150–64. Available from:
41 746 <http://linkinghub.elsevier.com/retrieve/pii/S1074761308002835>
- 42
43
44 747 41. Benita Y, Cao Z, Giallourakis C, Li C, Gardet A, Xavier RJ. Gene enrichment profiles reveal
45 748 T-cell development, differentiation, and lineage-specific transcription factors including ZBTB25
46 749 as a novel NF-AT repressor. *Blood* [Internet]. 2010 [cited 2018 Mar 5];115:5376–84. Available
47 750 from: <http://www.bloodjournal.org/cgi/doi/10.1182/blood-2010-01-263855>
- 48
49
50 751 42. Westerhuis JA, van Velzen EJJ, Hoefsloot HCJ, Smilde AK. Multivariate paired data
51 752 analysis: multilevel PLSDA versus OPLSDA. *Metabolomics* [Internet]. 2010 [cited 2016 Jul
52 753 27];6:119–28. Available from: <http://link.springer.com/10.1007/s11306-009-0185-z>
- 53
54 754
- 55
56
57
58
59
60
61
62
63
64
65

a**b**

a**b****c**







Click here to access/download
Supplementary Material
Supplementary Figures.docx





Click here to access/download
Supplementary Material
Supplementary Note.docx

THE UNIVERSITY OF
MELBOURNEFACULTY OF
SCIENCE

Dr. Kim-Anh Lê Cao
Senior Lecturer, Statistical Genomics
NHMRC Career Development Fellow
School of Mathematics and Statistics
Melbourne Integrative Genomics
The University of Melbourne

9th April 2018

Dear Editor of *Genome Biology*,

We wish to submit our manuscript "**DIABLO: identifying key molecular drivers from multi-omic assays, an integrative approach**" for consideration as a research article in your journal.

In the omics era, computational solutions to integrate different types of biological data measured on the same specimens or samples are trailing behind data generation. Our manuscript aims to fill this gap by proposing an efficient, flexible and easy-to-use computational framework to integrate multiple omics data generated from emerging high-throughput technologies.

The main challenge we face in multi-omics data integration is the large heterogeneity and difference in scales between omics platforms. Statistical integrative methods for biomarker discovery are still at their infancy and provide limited insight into complex biological processes. They are built on existing methods that either concatenate or combine the independent analyses from each data set, and do not model the correlation structure between the different molecular levels. This is highly problematic as important information can be missed, leading to incorrect conclusions. DIABLO maximises the correlation between data sets whilst identifying the key molecular features that explain and reliably classify a phenotype of interest. The dimension reduction process enables intuitive visualisations of the samples and selected multi-omics signatures. We benchmarked and demonstrated the ability of our method to select relevant correlated and discriminative biomarkers in a comprehensive simulation studies and in six multi-omics studies including two case studies in human breast cancer and asthma. In each of those studies we integrated various omics datasets ranging from transcriptomics (mRNA, miRNA), epigenomics (CpGs), proteomics and cell-type frequencies.

DIABLO facilitates the integration of large and heterogeneous data sets to identify relevant biomarker candidates in a wide range of biological settings. The method will be of significant interest to the scientifically diverse readership of *Genome Biology* to capitalise on fastly generated multi-omics data and push novel biological discoveries to an unprecedented level.

We are fervent advocates of open data and open science. All analyses are available in R markdown format as supplementary material, and the method is implemented in the open source R package mixOmics, with detailed tutorials on our companion website <http://www.mixOmics.org/mixDIABLO>.

The submitted manuscript has been approved by all authors and has not been submitted to any other journal. This manuscript is a substantial revised version to our previous submission to *Genome Biology*, **GBIO-D-16-01112**. We improved the method and added

four more case studies to benchmark the method to address the reviewers' comments. We provide a point-by-point response to reviewers in the next section. We look forward to your reply.

Yours sincerely,
Dr. Kim-Anh LÊ CAO

A handwritten signature in black ink, appearing to read "Lê Caо". The signature is somewhat stylized and includes a small checkmark or flourish at the end.

Reviewer #1:

The article has several strengths:

- a) *The article is very well written and provides a good overview of various statistical methods for analyzing genomic data.*
- b) *I think it presents an honest analysis of the data. The authors resist the temptation to oversell their method.*
- They acknowledge that their method does not outperform existing methods when it comes to accuracy.*
- c) *The authors have implemented the method in an R package*
- d) *This is a multi-omic method that integrates data.*
- e) *The authors apply their method to both empirical data and to simulated data.*

We appreciate the positive comments from the reviewer and the careful review. In the previous iteration of the manuscript our main focus was on the classification performance using a single breast cancer case study. However, one clear benefit of our approach is that the molecular signatures identified bring **superior biological enrichment** compared to other methods that we benchmarked on an **additional four multi-omics cancer datasets** (lung, kidney, colon and glioblastoma), each with three types of omics data (mRNA, miRNA and CpGs).

There are a few weaknesses.

The method is quite complicated and involves several parameter choices surrounding the underlying correlation structure. Why use a complicated method when simpler methods have similar predictive accuracy?

We believe our revision had addressed these weaknesses. We agree that the depiction of the method was lacking important details that led to misleading interpretations. In the revised manuscript, we have extensively benchmarked other multi-omics integration methods and demonstrated that DIABLO does not require as many parameters settings (see **Supplement**) as compared to existing methodologies. Briefly, our method requires 3 parameters, 1) number of variables to select from each omic dataset, 2) number of components to select from each omics dataset and 3) whether the correlation between certain omics datasets should be maximised (e.g. mRNA and miRNA). We provide a tuning function to choose parameters 1 and 2. For parameter 3 we provide guidelines that either rely on biological assumptions or a data-driven approach. Our method not only focuses on extracting the correlation structure across omic datasets but also discriminates between phenotypic groups. Such integrative approach is the first of its kind to identify molecular signatures with biological relevance and led to superior biological enrichment across various collections of gene set databases.

Why measures different types of data when a single data source (e.g. mRNA) already leads to good accuracy?

We agree with the reviewer that if the focus is on biomarker discovery, why not use the simplest, and cheapest strategy to identify biomarkers? But the focus of this manuscript is rather to capitalise on multi-omics studies and extract complementary information across omics data. Therefore, our focus is not only on identifying strong biomarkers, but also markers correlated across functional levels to give more insight into disease mechanisms. Therefore, we have changed the title to "**DIABLO: identifying key molecular drivers from multi-omic assays, an integrative approach**", to reflect the focus on key molecular drivers rather than biomarkers only.

I am not convinced that the method helps to elucidate the underlying biology. I understand that the latent structure might uncover interesting biology but I would never use this method to learn biology. Rather, I would use cluster analysis or unsupervised learning methods.

In the previous version of the manuscript we did not compare the biological enrichment of the various methods that were used. However, based on the reviewers' comments, we extensively explored this area, using multiple cancer multi-omics datasets, multiple gene-set databases with both unsupervised and supervised integrative methods that can perform variable selection. We demonstrate that our method outperforms unsupervised methods with respect to biological enrichment thus elucidating more known biology. In the human breast cancer study, we show that DIABLO can also detect novel biomarkers that have not been previously associated with breast cancer.

We also researched the literature to give an overview of the current state in integrative methods either supervised or unsupervised, and with or without variable selection, to highlight where the gaps are in terms of methods development (see **Supplementary Fig. 1**). In the revised version of the manuscript we have included unsupervised methods used for multi-omics data integration as well as supervised multi-step approaches.

Overall, I am not sure how much biology can be learnt by applying this method. Bottom line: this predictive method does not seem to improve predictive accuracy.

We have refocused our manuscript on biological insights primarily, rather than prediction performance as the former was our main motivation in driving methodological developments. By extending our analyses with six multi-omics studies including two case studies in human breast cancer and asthma we believe we have demonstrated that data integration performed using appropriate computational methods generate new biological insights and novel hypotheses to be further tested in the laboratory. The important contribution of DIABLO is its resulting molecular signatures that both explain the correlation structure across multiple biological domains and discriminate multiple phenotypic groups, with increased biological enrichment compared to other methods.

Reviewer #2:

This is a well written article which addresses an important need in the field. 1) In the introduction, the longer intro to sparse CCA should be provided. In the methods the actual method is more clearly stated "DIABLO extends sparse gCCA to a classification framework".

We thank the reviewer for their appreciative comments. In the revised version of the manuscript we provide a clearer explanation of our method DIABLO (see lines 99-116).

*"DIABLO (**D**ata **I**ntegration **A**nalysis for **B**iomarker **d**iscovery using **L**atent **c**Omponents) maximizes the common or correlated information between multiple omics (multi-omics) datasets while identifying the key omics variables (mRNA, miRNA, CpGs, proteins, metabolites, etc.) and characterizing the disease sub-groups or phenotypes of interest. DIABLO uses Projection to Latent Structure models (PLS) [1], and extends both sparse PLS-Discriminant Analysis [2] to multi-omics analyses and sparse Generalized Canonical Correlation Analysis [3] to a supervised analysis framework. In contrast to existing penalized matrix decomposition methods [4], DIABLO is a component-based method (or a dimension reduction technique) that transforms each omic dataset into latent components and maximizes the sum of pairwise correlations between latent components (user-defined) and a phenotype of interest [5]. DIABLO is, therefore, an integrative classification method that builds predictive multi-omics models that can be applied to multi-omics data from new samples to determine their phenotype. Users can specify the number of variables to select from each dataset and visualize the omics data and the multi-omics panel into a reduced data. The method is highly flexible in the type of experimental design it can handle, ranging from classical single time point to cross-over and repeated measures studies. Modular-based analysis can also be incorporated using pathway-based module matrices [6] instead of the original omics matrices, as illustrated in*

one of our case studies.”

The mathematical formulas such as the sGCCA algorithm, and its extension to a discriminant framework is detailed in the Methods section.

2) Can the approach handle missing data, that is missing row or column observations or is it only missing datasets. I presume, the later, as the intersection of tumors with complete data was used in training real data. This is important and should be made clear in the intro, abstract and discussion.

Currently, our method does not account for completely missing observations or variables.

Random missing values are allowed in the dataset matrices as local regressions are fitted in the model and missing values will be omitted when calculating the latent components and loading vectors. The prediction step however, as highlighted in the Breast Cancer case study can be performed with an entire dataset missing.

I 401: ‘*As the class prediction relies on individual vote from each omics set, DIABLO allows for some missing datasets X_k during the prediction step, as illustrated in the Breast Cancer case study.*’

I209: ‘*The training data consisted of four omics-datasets (mRNA, miRNA, CpGs and proteins) whereas the test data included all remaining samples for which the protein expression data were missing.*’

As such we do not think this is an information that should appear all throughout the document.

3) Can PLS DA be applied to multi class classification. Was this tested?

Yes, the revised version of the manuscript uses sparse Partial Least Squares Discriminant Analysis (sPLS-DA), in various multi-step classification schemes such as concatenation and ensemble-based schemes. Generally speaking SPLS-DA can handle multiple classes (Lê Cao et al., 2011, BMC Bioinformatics 22:253). However, this is not highlighted in this study as we used sPLS-DA for the cancer benchmark data sets that only include 2 classes.

4) The order of pair comparisons appears important. (discussion page 17, 18 and methods).

Those unfamiliar with their data may specify a suboptimal Design Matrix. Could there be some tools that provide guidance? For example, multiple factorial analysis or one of many tensor decompositions could be used to compute an RV coefficient. Alternative, can datasets be weighted in the analyses? In multi dataset approaches, data are often weighted by quality/size, the first eigenvector etc (reviewed by Meng et al., Brief Bioinform (2016) doi: 10.1093/bib/bbv108). If data has a batch effect, and this data were used to seed the analysis (aka in the first pair of data analyzed), would that skew the results? Could this please be tested.

We thank the reviewer for their suggestions In fact, there is no order to the pairwise comparisons in the DIABLO framework: it is the **sum of pairwise correlations** that is maximized, see **Methods**. Therefore the pairwise correlations are considered simultaneously in the SGCCA algorithm[3].

In the revised manuscript we got inspiration from the multiblock literature such as multiblock partial least squares (MBPLS[7]) whose datasets (also called blocks) are weighted based on their correlation with the response variable. In the new DIABLO implementation, we have used a weighted majority vote scheme based on the correlation between the latent component of each omics dataset with the latent component from the response matrix. This has significantly improved our classification error rates, as the strongest discriminatory datasets is given a higher weight in the overall class prediction for a new sample. Further, the weighted majority vote option in our function overcome the case where an equal number of voting classifiers and no consensus can be achieved.

In the discussion, we underlined the influence of batch effects on the multivariate modelling performed by our method (see lines 322-327), as this is outside the focus of this manuscript (but developments are in progress for other types of data).

“Finally, DIABLO, like other methods we benchmarked, will be affected by technical artifacts of the data, such as batch effects and presence of confounding variables that may affect downstream integrative analyses. Therefore, we recommend exploratory analyses be carried out in each single omics dataset to assess the effect, if any, of technical factors and use of batch removal methods prior to the integration analysis [8–10].”

5) On page 8, "validation of the Diablo methods on synthetic data". Three different criteria are explored 1) CorNonDis 2) CorDis 3) NonCorDis. Please explain the rational behind nonCorDis should be explained. In a 2 class system, methods such as CCA or PLS extract eigenvectors of correlated variables. Therefore a discriminate eigenvector will represent a set of correlated variables. Gene expression and 'omics data, measure genes which work in pathways, and therefore data has considerable correlation structure. Discriminatory non-correlated vectors, may reflect system noise.

The rationale for including four types of variables was to determine the influence of the correlation structure between datasets as well as discrimination between phenotypic groups. This is why the simulation framework includes different combinations of discrimination (discriminatory, non-discriminatory) and correlation (correlated, uncorrelated) variables. Unlike CCA and PLS which can only maximize the correlation between at most two data matrices, DIABLO can simultaneously maximize the correlation between any number of data matrices (see **Methods**). In our previous simulation study, we generated four types variables by controlling the correlation between variables or discrimination between groups. In our revised version, we have instead generated the correlation structure first by controlling the different relationships between latent components of different datasets (see **Supplementary Fig. 2**). The latent components are than used to compute the four-types of variables based on different correlation structures (see **Supplement for complete details**). The relevant variables include 30 corDis (correlated and discriminatory) and 30 unCorDis (uncorrelated and discriminatory) variables, in order to determine the effect of the design matrix on the types of variables selected. We also simulated 100 corNonDis (correlated and non-discriminatory) and 100 unCorNonDis (uncorrelated and non-discriminatory) variables. Therefore it is the corNonDis and unCorNonDis variables that represent noise and irrelevant variables, although by chance some of these variables might be correlated with the response. The purpose of the simulation was to determine whether any of multi-step classification schemes and DIABLO model happen to (wrongly) select these irrelevant variables (corNonDis, unCorNonDis) and relevant variables (corDis, unCorDis).

6) please provide a discussion on filtering data. In each case, data were filtered and reduced. Is this to reduce "noise" or for computational efficiency. Please discuss and comments on the computational cost of larger datasets.

We provide a discussion about filtering data in the revised version of the manuscript (lines 319-320). "...we advise using a broad filtering strategy to alleviate computational time when dealing with extremely large datasets (e.g. > 50,000 features each)." However, for this revised manuscript, we have not performed any filtering for the benchmark datasets and retained all the variables that were downloaded from their respective websites. For the breast cancer case study, some filtering was involved to remove low abundance variables, mostly to reduce some amount of noise rather than saving on computational time. We also provide additional guidance on filtering in the mixOmics article [11], which we refer to in this manuscript. **Table 1** lists the size of the datasets we analysed.

7) On page 11 the acronym BER (balanced error rate) is used before it is defined.

Since we do not use the acronym BER many times in the revised manuscript, we have removed it altogether and explicitly stated 'balanced error rate'.

8) on Page 13, please describe "eigengene summarization" in more detail. Please describe how to interpret the results, saying it is "common approach" on page 29 is insufficient;

We provide a description of eigengene summarization in the revised manuscript (lines 504-509).

"Modular analysis: Eigengene summarization is a common approach to decompose a $n \times p$ dataset (where n is the number of samples and p is the number of variables in a module), to a component (linear combination of all p variables) that represents the summarized expression of genes in the module [6]. For the asthma study, 15,683 genes were reduced to 229 KEGG pathways and 292 metabolites were reduced to 60 metabolic pathways using eigengene summarization."

9) page 16. Why were 9 variables (36 in total) selected in analysis of the BRCA data? Is there any guidance as to how many variables should be selected. For example, mRNA and protein were more informative in gsea, therefore it might be better to select more variables from these datasets?

Yes, we provide a tuning function which is implemented along with DIABLO in the mixOmics R package [11]. The function uses a grid approach to select an optimal number of variables to select from each omics dataset. A section on parameter tuning discusses the grid approach to identifying the optimal number of variables and components to select (lines 438-456).

*"Finally, the third set of parameters to tune is the number of variables to select per dataset and per component. Such tuning can rapidly become cumbersome, as there might be numerous combinations of selection sizes to evaluate across all K datasets. For the breast cancer study, we used 5-fold cross-validation repeated 50 times to evaluate the performance of the model over a grid of different possible values of variables to select (**Supplementary Fig. 8**). The performance of the model for a given set of parameters (including number of component and number of variables to select) was based on the balanced classification error rate using majority vote or average prediction schemes with centroids distance. The balanced classification error rate is useful in the case of imbalanced class sizes, where the majority classes can have strong influence on the overall error rate. The balanced error rate measure calculates the weighted average of the individual class error rates with respect to their class sample size. In our experience, the number of variables to select in each dataset provided less of an improvement on the error rate compared to tuning the number of components. Therefore, even a grid composed of a small number of variables (<50 with steps of 5 or 10) may suffice as it does not substantially change the classification performance. This is because of the use of regularization constraints which reduces the variability in the variable coefficients and thus maintains the predictive ability of the model. Further, the variable selection size can also be guided according to the downstream biological interpretation to be performed. For example, a gene-set enrichment analysis may require a larger set of features than a literature-search interpretation."*

10) page 17, 3rd line from top. "known cell-types and pathways in the" typo. Insert space
Thank you, the typo has been corrected in the revised version of the manuscript.

11) The statement on p 18, "To our knowledge DIABLO is the only integrative classification methods that models the correlation structure between omics data" is inaccurate. See the Meng et al., review. There are many many methods that use dimensions reduction to extract correlated structure in multiple 'omics data. Also Jeffrey et al., 2007 Bioinformatics. 2007 Feb 1;23(3):298P305. Epub 2006 Nov 24. used discriminative analysis with coinertia analysis and described a supervised integrative latent variable approach, which is related to this work but did not employ sparse methods.

We have removed this erroneous statement from the revised manuscript, thank you.

12) In the methods, p21. In the abstract/intro, the method is describes as a supervised PLS/DA approach, but on p21 it appears to be cluster to partitioning around centroids, with majority voting. Please describe the approach clearly and consistently.

In the revised version of the manuscript we have added additional details regarding the development of DIABLO, which extends sparse Generalised Canonical Correlation Analysis. We believe the confusion from the reviewer may come from the fact that multiple-types of prediction distances can be used in DIABLO such as centroids, max distance, and Mahalanobis distance. Please see lines 347-383 for a general description of the sGCCA algorithm, lines 386-397 for the classification implementation and lines 399-417 for the implementation for the different types of error rate that can be computed.

13) is the analysis effective by the number of variables. For example if dataset A has several thousand variables and dataset B has less than 50, would this impact the analysis?

The difference in the number of variables in each dataset should not impact the analysis as each dataset is summarised by its own set of latent components, so that components across data sets are maximally correlated, irrespective of how many variables there are in each dataset. This makes DIABLO a much more attractive solution than a concatenation method, where datasets than include a large number of variables tend to be more ‘favoured’ in the molecular signature compared to smaller data sets.

14) p25, The de-duplication effort in GSEA is important and should be clear to users, If a more stringent assignment were used, would this impact results?

GSEA is impacted by the number of features that are input into the analysis and the types of gene sets that are used to determine biological enrichment. In the revised manuscript we include a benchmarking experiment where we constructed multi-omic biomarker panels of equivalent number of features with a total of 180 features. Further we tested 10 different gene set databases, from Molecular signature database[12], blood transcriptional modules[13] and cell-specific expression from Benita *et al.* [14].

15) p 27. Data Processing. Were 3,073 BRCA clinical variables used in this study? The PAM50 assignments for tumors (obtained from TCGA staff) should be made available together with the filtered TCGA data, such that others can reproduce this work.

The 3,073 variables listed describe the data that were obtained from TGCA. From the clinical data, only the PAM50 labels and sample-type variables were used. The complete code and data files can be found with the github repository (<https://github.com/singha53>).

16) p28. Terms in the Voom equation are not fully defined. Filtering removed "genes with counts less than 0". Does this mean the sum of the gene across all tumors was zero, or that any gene which has a zero tumor in any 1 tumor was excluded?

We have clarified the following in the revised manuscript (see **Supplementary Data file**). The count data for the mRNA dataset, X_{counts} was normalized to log2-counts per million (logCPM), X_{norm} , similar to limma voom [15]:

$$X_{norm} = \log_2 \left(\frac{(X_{counts} + 0.5)^T}{(lib.size + 1) * 10^6} \right)$$

After library size (lib.size = total number of reads per sample) normalization, genes with counts less than 0 in more than 70% of samples were removed. The PAM50 genes were also removed from the mRNA dataset prior to analyses. Similarly, the miRNA count data was normalized to logCPM and miRNA transcripts with counts less than 0 in more than 70% of the samples were also removed.

17) p28 Asthma study. Genes were reduced 229 KEGG pathways and metabolites were reduced to 60 pathways. Why were variables reduced to GeneSet. The rational and need for this is not

explained. Was it simply to aid biological interpretation of the data or was it for computational reasons?

For this specific case study we wished to incorporate modular-based analyses within the DIABLO framework to focus on pathways spanning common biological mechanisms that significantly changed in response to allergen inhalation challenge. The purpose of this analysis was only to aid in the biological interpretation and the reduction to gene sets was not performed for computational reasons. A secondary reason for including this approach was to demonstrate to potential users the benefits of combining modular-based analyses with the DIABLO framework. Other types of approaches that identify modules such as data-driven techniques like WGCNA (weighted gene co-expression networks) may also be incorporated with the DIABLO framework, since each cluster of variables can be reduced to a single variable that explains the entire cluster of features.

18) Figure 1 A) not clear if concatenation is performed on genes or tumors (rows/cols). C) The DIABLO diagram is confusing. it is not clear that DIABLO is a pairPwise approach.

The figure has been updated to clarify the integration and classification aspects of the DIABLO framework (see **Supplementary Figure 3**). Each dataset is a $n \times p_j$ matrix, where p_j is the number of variables (columns) for the j^{th} dataset. For the concatenation-based analysis, the datasets are combined row-wise since the number of samples are the same for each omics dataset, that is, the multi-omics data is obtained for the same set of samples.

Although DIABLO computes the pairwise correlation between latent components of pairs of omic datasets, similar to PLS and CCA, its objective is to maximize the sum of pairwise correlation between different omics datasets (see objective function in **Methods**), see our earlier answer to Question 4.

19) Figure 4 legend. DIABLO 1P12 are not defined, What was the difference between these models.
In the previous version of the manuscript, the concatenation and ensemble biomarker panels were tuned such that each panel consisted of a specific number of variables and DIABLO panels matched the same number of variables to keep the comparisons consistent. However, given this extra confusion, we have added 4 benchmark datasets where each method selects the same number of variables of each omic-type.

20) Figure 5. There is no scale on the ciros plot (gene level) which makes its interpretation difficult. Also please add Gene Names to the heatmap (E)

The purpose of the circosplot is to depict the inter-correlations between omic datasets. These can be observed from the red (positive) and blue (negative) lines between omics datasets (different colors). The scale for the lines surrounded the ideogram is not depicted as it is centered at zero, therefore the line height represent the average expression levels of a given variable in a given phenotypic groups compared to others. Gene names have not been added to the heatmap, due to size limitation of the figure. However, the feature plot in Figure 3a lists all the features selected by the multi-omic biomarker panel.

21). Reference 38 Gauvreau et al., is in upper case

All references have been checked and the capitalization has now been fixed.

22) Please provide more details on the computational complexity of the method as 1) the number of variables increases 2) the number of datasets increased 3) the impact of correlated datasets (eg microarray and RNAseq)

We decided to focus our comparisons mainly on the correlation and discrimination structure between datasets (simulation study), the effect of the design matrix on the types of variables selected and whether this led to superior biological enrichment compared to other integrative strategies

(benchmark study, four new real multi-omics datasets). Computational times are provided for different scenarios in our article (Rohart et al, 2017, Plos Computational Biology 13 [11] in the main and supplemental material) for various numbers of variables and data sets, see the screenshots below.

Table 2. Example of computational time for the data sets presented in the Results section with a macbook pro 2013, 2.6GHz, 16Go Ram.

Framework	Single 'omics sPLS-DA		<i>N</i> -integration DIABLO		<i>P</i> -integration MINT	
Data	srbc		breast.tcga		stemcells	
<i>N</i>	63		150		125	
<i>P</i>	2, 308		200; 184; 142		400	
function	tune	perf	tune	perf	tune	perf
#fold CV (repeated)	5(10)	5(10)	10(1)	10(10)	LOGOCV	LOGOCV
ncomp	6	3	2	2	2	2
grid length per component	39	-	13 ³	-	100	-
#cpu	1	1	2	1	1	1
run time	9min	31sec	18min	25sec	30sec	0.2sec

Figure 1. Computational time for various mixOmics methods, including DIABLO, as published in [11] in the main article.

Table 4: Example of runtime for very large data sets analysed in mixOmics. Tuning and performance assessments were performed with 5-fold CV for single 'omics and N-integration, or LOGOCV for P-integration (Rohart et al. 2017, Singh et al. 2016, cluster with 10 cpus and 50 Gb RAM).

Framework	Single 'omics sPLS-DA		<i>N</i> -integration DIABLO		<i>P</i> -integration MINT	
Data	HNSCC		Asthma (2 omics)		Stem Cell (8 studies)	
<i>N</i>	60		194		210	
<i>P</i>	82, 132		30,000; 30,000		13, 313	
function	tune	perf	tune	perf	tune	perf
#fold CV (repeated)	5(10)	5(10)	5(1)	5(10)	LOGOCV	LOGOCV
ncomp	5	3	2	2	2	2
grid length per component	40	-	22 ²	-	100	-
#cpu	10	10	10	10	1	1
runtime	15min	6min	19min	3min	17min	12sec

Figure 2 Computational time for various mixOmics methods, including DIABLO, as published in [11] in the supplement material.

References

1. Wold H. Estimation of principal components and related models by iterative least squares. *Multivar Anal.* 1966;391–420.
2. Lê Cao K-A, Boitard S, Besse P. Sparse PLS discriminant analysis: biologically relevant feature selection and graphical displays for multiclass problems. *BMC Bioinformatics* [Internet]. 2011 [cited 2015 Jul 15];12:253. Available from: <http://www.biomedcentral.com/1471-2105/12/253/>
3. Tenenhaus A, Philippe C, Guillemot V, Le Cao K-A, Grill J, Frouin V. Variable selection for generalized canonical correlation analysis. *Biostatistics* [Internet]. 2014 [cited 2015 Jul 15];15:569–83. Available from: <http://biostatistics.oxfordjournals.org/cgi/doi/10.1093/biostatistics/kxu001>
4. Witten DM, Tibshirani R, Hastie T. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics* [Internet]. 2009 [cited 2016 Jul 27];10:515–34. Available from: <http://biostatistics.oxfordjournals.org/cgi/doi/10.1093/biostatistics/kxp008>

5. Lee HK, Hsu AK, Sajdak J, Qin J, Pavlidis P. Coexpression analysis of human genes across many microarray data sets. *Genome Res* [Internet]. 2004 [cited 2016 Mar 30];14:1085–1094. Available from: <http://genome.cshlp.org/content/14/6/1085.short>
6. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics* [Internet]. 2008 [cited 2016 Apr 4];9:559. Available from: <http://www.biomedcentral.com/1471-2105/9/559>
7. BOUGEARD S, QANNARI EM, LUPO C, HANAFI M. From Multiblock Partial Least Squares to Multiblock Redundancy Analysis. A Continuum Approach. :16.
8. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* [Internet]. 2007 [cited 2016 May 12];8:118–27. Available from: <http://biostatistics.oxfordjournals.org/cgi/doi/10.1093/biostatistics/kxj037>
9. Gagnon-Bartsch JA, Speed TP. Using control genes to correct for unwanted variation in microarray data. *Biostatistics* [Internet]. 2012 [cited 2018 Mar 6];13:539–52. Available from: <https://academic.oup.com/biostatistics/article-lookup/doi/10.1093/biostatistics/kxr034>
10. Parker HS, Corrada Bravo H, Leek JT. Removing batch effects for prediction problems with frozen surrogate variable analysis. *PeerJ* [Internet]. 2014 [cited 2016 May 12];2:e561. Available from: <https://peerj.com/articles/561>
11. Rohart F, Gautier B, Singh A, Cao K-AL. mixOmics: An R package for ‘omics feature selection and multiple data integration. *PLOS Comput Biol* [Internet]. 2017 [cited 2018 Jan 29];13:e1005752. Available from: <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1005752>
12. Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The molecular signatures database hallmark gene set collection. *Cell Syst* [Internet]. 2015 [cited 2018 Jan 30];1:417–25. Available from: <http://linkinghub.elsevier.com/retrieve/pii/S2405471215002185>
13. Chaussabel D, Baldwin N. Democratizing systems immunology with modular transcriptional repertoire analyses. *Nat Rev Immunol* [Internet]. 2014 [cited 2016 Jul 22];14:271–80. Available from: <http://www.nature.com/doifinder/10.1038/nri3642>
14. Benita Y, Cao Z, Giallourakis C, Li C, Gardet A, Xavier RJ. Gene enrichment profiles reveal T-cell development, differentiation, and lineage-specific transcription factors including ZBTB25 as a novel NF-AT repressor. *Blood* [Internet]. 2010 [cited 2018 Mar 5];115:5376–84. Available from: <http://www.bloodjournal.org/cgi/doi/10.1182/blood-2010-01-263855>
15. Law CW, Chen Y, Shi W, Smyth GK. Voom: precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol* [Internet]. 2014 [cited 2016 Mar 2];15:R29. Available from: <http://www.biomedcentral.com/content/pdf/gb-2014-15-2-r29.pdf>