

OnPLS-Based Multi-Block Data Integration: A Multivariate Approach to Interrogating Biological Interactions in Asthma

Stacey N. Reinke,^{*,†,‡} Beatriz Galindo-Prieto,^{§,||,⊥} Tomas Skotare,[§] David I. Broadhurst,[‡] Akul Singhania,^{#,∇} Daniel Horowitz,[○] Ratko Djukanović,^{#,◆} Timothy S.C. Hinks,^{#,◆,††} Paul Geladi,^{‡‡} Johan Trygg,^{§,Ⓛ} and Craig E. Wheelock^{*,†,§§,Ⓛ}

[†]Division of Physiological Chemistry 2, Department of Medical Biochemistry and Biophysics, Karolinska Institute, SE-171 77 Stockholm, Sweden

[‡]Centre for Integrative Metabolomics and Computational Biology, School of Science, Edith Cowan University, Perth 6027, Australia

[§]Computational Life Science Cluster, Department of Chemistry (KBC) and ^{||}Industrial Doctoral School (IDS), Umeå University, SE-901 87 Umeå, Sweden

[⊥]Department of Engineering Cybernetics (ITK), Norwegian University of Science and Technology (NTNU), 7491 Trondheim, Norway

[#]Clinical and Experimental Sciences, University of Southampton Faculty of Medicine and [◆]NIHR Southampton Respiratory Biomedical Research Unit, Southampton University Hospital, Southampton SO16 6YD, U.K.

[∇]Laboratory of Immunoregulation and Infection, The Francis Crick Institute, London NW1 1AT, U.K.

[○]Janssen Research and Development, High Wycombe HP12 4DP, Buckinghamshire, U.K.

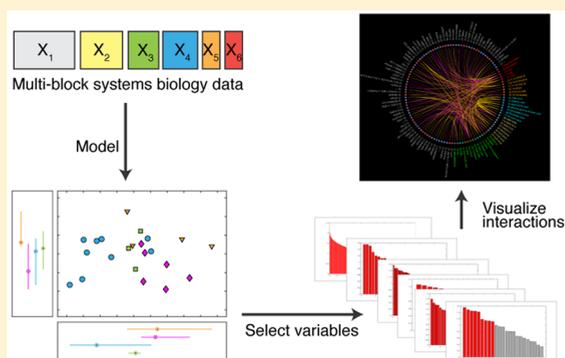
^{††}NIHR Oxford Biomedical Research Centre/Respiratory Medicine Unit, NDM Experimental Medicine, University of Oxford, Level 7, John Radcliffe Hospital, Oxford OX3 9DU, U.K.

^{‡‡}Forest Biomass and Technology, Swedish University of Agricultural Sciences, SE 90183 Umeå, Sweden

^{§§}Gunma University Initiative for Advanced Research (GIAR), Gunma University, Maebashi 371-8510, Japan

Supporting Information

ABSTRACT: Integration of multiomics data remains a key challenge in fulfilling the potential of comprehensive systems biology. Multiple-block orthogonal projections to latent structures (OnPLS) is a projection method that simultaneously models multiple data matrices, reducing feature space without relying on a priori biological knowledge. In order to improve the interpretability of OnPLS models, the associated multi-block variable influence on orthogonal projections (MB-VIOP) method is used to identify variables with the highest contribution to the model. This study combined OnPLS and MB-VIOP with interactive visualization methods to interrogate an exemplar multiomics study, using a subset of 22 individuals from an asthma cohort. Joint data structure in six data blocks was assessed: transcriptomics; metabolomics; targeted assays for sphingolipids, oxylipins, and fatty acids; and a clinical block including lung function, immune cell differentials, and cytokines. The model identified seven components, two of which had contributions from all blocks (globally joint structure) and five that had contributions from two to five blocks (locally joint structure). Components 1 and 2 were the most informative, identifying differences between healthy controls and asthmatics and a disease–sex interaction, respectively. The interactions between features selected by MB-VIOP were visualized using chord plots, yielding putative novel insights into asthma disease pathogenesis, the effects of asthma treatment, and biological roles of uncharacterized genes. For example, the gene *ATP6 V1G1*, which has been implicated in osteoporosis, correlated with metabolites that are dysregulated by inhaled corticoid steroids (ICS), providing insight into the mechanisms underlying bone density loss in asthma patients taking ICS. These results show the potential for OnPLS, combined with MB-VIOP variable selection and interaction visualization techniques, to generate hypotheses from multiomics studies and inform biology.



In the postgenomic era, data-driven science has become increasingly necessary because of the vast array of instrumentation that is capable of generating thousands of data points for a single analytical observation.^{1,2} In addition to

Received: July 17, 2018

Accepted: October 18, 2018

Published: October 18, 2018

49 using classical univariate statistical methods, machine-learning
50 techniques have become routinely used to interrogate and
51 understand vast amounts of data.^{3,4} Two common character-
52 istics of -omics data are that the number of measured variables
53 is vastly greater than the number of observations⁵ and that
54 there is a degree of multicollinearity between variables.⁶ As
55 such, computational methods that project high dimensional
56 data into a smaller number of component variables have
57 become commonplace.⁷ Multivariate projection methods such
58 as principal components analysis (PCA),⁸ partial least squares
59 discriminant analysis (PLS-DA),⁹ and canonical variate
60 analysis (CVA),⁸ together with hierarchical cluster analysis
61 (HCA),¹⁰ random forests,¹¹ and support vector machines
62 (SVM),¹² are all used to analyze -omics data.^{3,4} PLS-DA and
63 its extension, orthogonal projection to latent structures
64 discriminant analysis (OPLS-DA),^{13,14} have become popular
65 projection methods in the metabolomics community.¹⁵ As
66 modeling methods become increasingly complicated, they have
67 also become concomitantly difficult to interpret. Assignment of
68 the variable importance often becomes an a posteriori
69 statistical process based on either permutation testing or
70 random resampling (e.g., confidence intervals derived from
71 bootstrap/jackknife statistics).¹⁶ For methods based on a PLS
72 algorithm, the direct statistical method of variable influence on
73 projection (VIP)^{17,18} is often used to estimate variable
74 contribution to the resulting models.

75 In recent years, as the -omics sciences have matured, it has
76 become common to acquire data from multiple -omics
77 platforms in a single biological experiment. As such, each
78 biological sample is interrogated by multiple analytical
79 platforms, which in turn can be linked to multiple sources of
80 experimental metadata. Data from each platform (or measure-
81 ment context) can be considered a discrete *block*, with multiple
82 blocks making up the complete data set of the experiment.
83 Multivariate projection methods such as OPLS-DA have
84 proven successful in modeling the underlying latent biological
85 structure within a single high dimensional data block; however,
86 they are theoretically unsuitable for modeling multiple data
87 blocks simultaneously. There are two reasons for this issue.
88 First, if multiple data blocks are concatenated into a single
89 matrix, with no accounting for measurement context, then the
90 subsequent model can be considered as a single projection
91 model, where the weighting of each variable is governed by the
92 total sum of squares.¹⁹ This, in principle, demands that each
93 block is normalized to the same size, to avoid a projection
94 model that is biased toward the impact of the data set with the
95 most variables. In practice, this can be problematic, particularly
96 when there are a mix of blocks of vastly different sizes. For
97 example, in a model concatenating 20 000 transcripts, 200
98 metabolites, and 20 clinical variables, the transcripts would
99 over-represent the global data structure and thus have a larger
100 contribution to the resulting model. In multi-block modeling,
101 this is not an issue, as each block is treated independently. This
102 approach leaves flexibility to scale individual variables
103 according to importance and also to keep variables in their
104 original unit. Second, each individual data set is associated with
105 its own underlying structure,^{19,20} describing the true biological
106 variance and also platform-specific measurement error.
107 Covariance of biological latent structure across multiple data
108 blocks is implicit; however, it is a fair assumption that the
109 measurement error across multiple blocks will be independent
110 and thus easily ignored at this block-interaction level.
111 Conversely, if multiple data blocks are concatenated into a

single data set before projection, the model will struggle to
effectively separate true biological structure from block-specific
noise and result in erroneous interpretation of the conglom-
erate projection model.

To address the need for multivariate methods to
simultaneously model multiple data matrices, a number of
multi-block data integration methods have been pro-
posed.^{21–23} In 2011, Löfstedt and Trygg²⁴ proposed a novel
multi-block multivariate method called OnPLS, which utilizes
the framework of OPLS to decompose data from more than
two input matrices. Multi-block models, such as OnPLS, are
fully symmetric, meaning each data block is weighted to allow
an equal contribution to the model, regardless of the number
of variables or underlying data structure within each block.²⁵
Multi-block approaches offer further advantages over single
block or block concatenation in biomarker discovery. First, the
validity of any true biological biomarker is significantly
increased if there is a clear covariance between data blocks,
thus reducing the possibility of false discovery.²⁶ Second,
contrary to block-concatenation modeling, which is strongly
biased toward the globally joint variation, multi-block analysis
decomposes the different levels of variation (global, local,
unique)²⁷ such that relatively small but informative trends are
also identified. Recently, Galindo-Prieto et al. adapted the VIP
concept for multi-block data analysis (multi-block variable
influence on orthogonal projections method,²⁸ MB-VIOP) to
identify the variables that contribute to these different levels of
joint structure.

The aim of this study was combine OnPLS and MB-VIOP
with data visualization methods to create a workflow capable of
simultaneously modeling and investigating interactions be-
tween multiple -omics data blocks. The study chosen for this
purpose was a subset from a previously reported asthma
cohort, for which multiple -omics data sets were acquired in
isolation.^{29,30} These analyses included untargeted metabolo-
mics, targeted metabolite assays, differential immune cell
population analyses, and cytokine arrays. Additionally, for the
present study, transcriptomics of peripheral blood T cells was
performed. OnPLS modeling and MB-VIOP were then used to
integrate the disparate data blocks into a single model, which
was then interrogated to identify novel interactions between
the data blocks and disease status as well as other clinical end
points.

■ EXPERIMENTAL SECTION

Clinical Cohort. Briefly, 12 healthy controls and 10 severe
asthmatics were included from the original study.²⁹ Tran-
scriptomics was subsequently performed on peripheral blood T
cells, and metabolomics/metabolic profiling assays were
performed on serum. All participants were enrolled from the
NIHR Southampton Respiratory Biomedical Research Unit
and University Hospital Southampton outpatient clinics; all
provided written informed consent. The National Research
Ethics Service Committee South Central—Southampton B
ethics committee (UK; ref 10/H0504/2) approved this study.
Clinical classification and enrollment criteria were previously
described.^{29,31} Participant data were included in the present
study if they were classified as either healthy control or severe
asthmatic individuals in the existing cohort, and data from all
data blocks (described in next section) were collected.

Sample Collection and Analyses. Details of sample
collection and transcriptomics analyses are available in the
Supporting Information. Details of analytics, quality control, 173

174 and data cleaning for metabolomics, targeted metabolic assays,
 175 and clinical assays were performed as previously described.^{29,30}
 176 **Data Blocks and Processing.** Six data blocks were used
 177 for modeling: Transcriptomics, Sphingolipids, Metabolomics,
 178 Fatty Acids, Oxylipins, and Clinical Data (Figure 1). A

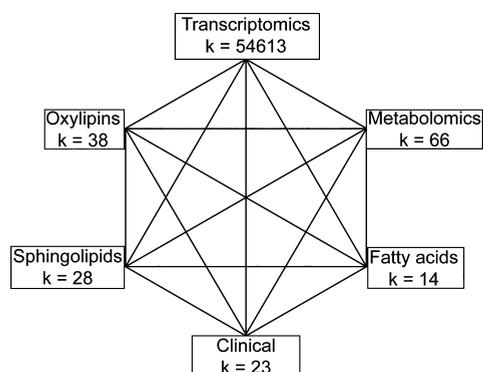


Figure 1. Schematic of potential shared structure between data blocks. The six data blocks used in this study are shown with their respective numbers of variables. The diagram shows all possible shared structure connections between the data blocks.

179 complete list of all variables included for each data block is
 180 provided in Tables S1–S6. The data blocks were defined by a
 181 priori knowledge about both the system under observation and
 182 the measurement technology.¹⁹ The primary consideration was
 183 that the underlying structure of the data could possibly
 184 confound the biological interaction between blocks. To avoid
 185 bias in combining information from different probes for one
 186 gene, all non-QC probes were included for OnPLS modeling;
 187 the Transcriptomics block included 54 613 variables. This
 188 approach is commonly used for analyzing transcriptomics
 189 data.³² Four data blocks represented serum metabolites:
 190 Sphingolipids (28 variables, targeted assay), Metabolomics
 191 (66 variables, untargeted assay screened against an in-house
 192 chemical library), Fatty Acids (14 variables, targeted assay),
 193 and Oxylipins (38 variables, targeted assay). A total of 23
 194 clinical variables were combined into the Clinical data block;
 195 these variables were derived from typical clinical assays and
 196 measurements and included lung function tests, bronchoalveo-
 197 lar lavage fluid and peripheral blood T cell populations, serum
 198 cytokines, and serum vitamin D3. For clinical data, values that
 199 were missing due to being below the limit of detection (LOD)
 200 of the respective assay were imputed with 1/10 of the lowest
 201 measured value, because the LOD was not known for each
 202 assay, and OnPLS cannot process missing values. Data that
 203 were missing for an entire subset array of the clinical data (e.g.,
 204 for individuals missing the cytokine assay) were imputed using
 205 the median value of the corresponding clinical group (control
 206 or asthma). Remaining missing values were replaced using
 207 PCA imputation. Prior to OnPLS model calculation, all data
 208 (except for transcriptomics) were log-transformed. All data
 209 were then scaled to unit variance.

210 **OnPLS Model Calculation and Visualization.** The
 211 OnPLS model simultaneously analyzed the data matrices,
 212 returning output matrices of shared information (compo-
 213 nents), as described.²⁷ These output matrices reveal shared
 214 data structure on three levels for each data matrix, which can
 215 be summarized as

216 Globally joint components reveal structure that is shared by
 217 all input data matrices. Locally joint components reveal

$$X_i = \underbrace{X_G}_{\text{globally joint}} + \underbrace{X_L}_{\text{locally joint}} + \underbrace{X_U}_{\text{unique}} + \underbrace{E}_{\text{residual noise}} \quad (1)$$

218 structure shared by two or more, but not all, of the input
 219 matrices. Finally, unique components identify latent structure
 220 that is present in only one input matrix. The OnPLS model
 221 returned separate score vectors for each data block in each
 222 component. To identify the sources of biological variance
 223 explained by the OnPLS components, the component scores
 224 for each block were correlated with metadata variables not
 225 included in the clinical data block: clinical class (control vs
 226 asthma), sex, age, BMI, dose of inhaled and oral cortico-
 227 steroids, and smoking (current/former smoker vs never
 228 smoked). The resulting Pearson correlation coefficients were
 229 presented as a metadata correlation plot.³³ To visualize the
 230 overall OnPLS model, hierarchical principal component
 231 analysis (PCA)³⁴ was used to summarize the 30 OnPLS
 232 score vectors, resulting in 2 PCA components describing the
 233 relationships in the OnPLS model. Prior to calculating the
 234 PCA model, the score vectors were scaled to unit variance. The
 235 PCA score plot showed individual participants, and the
 236 loadings plot displayed the score vectors from the OnPLS
 237 model, labeled by block type and OnPLS model component
 238 number.

239 **MB-VIOP Concept, Motivation, and Calculation.**
 240 Multi-block variable influence on orthogonal projections
 241 (MB-VIOP) is a feature selection method that (i) sorts the
 242 input variables by importance for data interpretation in OnPLS
 243 models, either for the total model (all variation types together)
 244 or per component (global, local, or unique variations
 245 separately), and (ii) explores the connections between the
 246 variables (either in the same or a different data matrix) that
 247 contribute to explain the same component (latent variable) in
 248 the multi-block system. Multi-block-VIOP is a model-based
 249 variable selection method, because it uses the *n* preprocessed
 250 data matrices, the score vectors, and the normalized loading
 251 vectors from an OnPLS model. OnPLS regression can relate
 252 the data matrices according to the model component; however,
 253 it must be emphasized that not all input variables of these
 254 related matrices will connect among themselves to explain the
 255 variation contained in a specific model component. The MB-
 256 VIOP algorithm is necessary to sort the input variables
 257 according to their connections for interpreting the variation
 258 contained in one or more specific components. Furthermore,
 259 MB-VIOP finds the degree of importance of each variable in
 260 the correct proportion for a multi-block system, which cannot
 261 be achieved by the OnPLS normalized loadings plot.³⁵

262 The calculation of the MB-VIOP values can be summarized
 263 as the Hadamard products of the normalized loadings
 264 multiplied by the ratio of the variation explained by a model
 265 component and the cumulated variation. After a block- and
 266 component-wise iterative algorithm with all input variables
 267 from the six data matrices involved, the resulting MB-VIOP
 268 vectors were normalized by Euclidean norm and by the
 269 number of original (input) variables raised to the 1/2 power.
 270 The variables of interest that were identified by MB-VIOP
 271 were selected as a subset for further multivariate analysis as
 272 shown below. For additional details about the MB-VIOP
 273 fundamentals and algorithm, readers are referred to the original
 274 references.²⁸

275 **Data Visualization.** The between-block covariance of the
 276 subset of variables contributing to Components 1 and 2 of the
 277 OnPLS model were visualized using chord plots.³⁶ Using the

278 variables reaching a defined MB-VIOP threshold, a chord plot
 279 was constructed by first calculating the Spearman rank
 280 correlation coefficient (r) for each pairwise combination of
 281 variables with MB-VIOP values above a threshold. Those
 282 variables where a significant ($p < 0.001$) between-block
 283 correlation existed were presented as nodes in a circle
 284 (grouped by block), and the correlation represented as a
 285 colored arc (yellow being a positive correlation and purple a
 286 negative correlation). The number of arcs associated with a
 287 given node is recorded in parentheses next to the name of the
 288 variable. Each chord plot was constrained such that within-
 289 block correlations were ignored.

290 Data modeling (OnPLS), variable selection (MB-VIOP), a
 291 posteriori analyses, and creation of plots were performed using
 292 MATLAB 2018a (Mathworks, Natick, MA, USA). Correlation
 293 coefficients for the metadata correlation plots were calculated
 294 using functions from SciPy (<http://www.scipy.org/>), and the
 295 plot was created using the Matplotlib library.³⁷ SIMCA v15
 296 (Umetrics, Umeå, Sweden) was used to perform OPLS-DA
 297 analysis.

298 ■ RESULTS AND DISCUSSION

299 **Study Population.** A total of 22 participants from a
 300 previously described cohort^{29,30} were included in this study
 301 (12 healthy control individuals and 10 individuals with severe
 302 asthma). Clinical information is presented in Table 1. Age and

Table 1. Clinical Data

	healthy control ($N = 12$)	severe asthma ($N = 10$)
age (years)	26.5 (24.8, 30.8)	63 (43.5, 63)
sex (M/F)	9/3	4/6
BMI (kg/m^2)	24.1 (22.6, 43.2)	34.0 (27.4, 43.2)
smoking status		
never smoker (#)	11	6
current/former smoker (#)	1	4
treatment		
inhaled corticosteroids (#, median dose ^b)	0	10 (1280)
oral corticosteroids (#)	0	3

^aValues are medians (interquartile range) or numbers. ^bBeclomethasone dipropionate equivalent μg .

303 BMI were significantly higher in the severe asthmatic group
 304 and thus represented confounders in the study. Furthermore,
 305 all individuals in the severe asthmatic group were treated with
 306 inhaled and/or oral corticosteroids (ICS/OCS). Although the
 307 sex ratio and proportion of smokers were also different, they
 308 were not significantly altered between the two groups.

OnPLS Model. The OnPLS model calculated seven 309
 components that shared joint structure between at least two 310
 of the data blocks (Table 2). Two components (1 and 4) had 311
 globally joint structure, with contributions from all six blocks. 312
 The remaining components had locally joint structure, with 313
 between 2 and 5 data blocks contributing to the joint structure. 314
 The model did not identify any unique components. 315

The amount of variance explained in each component, for 316
 each data block, as well as cumulative variance explained by the 317
 model is reported in Table 2. Only 37% of the total variance in 318
 the Transcriptomics data block was explained, indicating that 319
 the majority of the information contained in this block is not 320
 descriptive for describing asthma. This could be due to the 321
 global and unbiased nature of the platform and/or the fact that 322
 the transcriptomics data were derived from the entire 323
 peripheral blood CD3+ T cell population. It would be of 324
 more clinical relevance to target specific cell subpopulations in 325
 a single-cell transcriptomics approach.³⁸ The clinical data 326
 described only 55% of the variance in the Clinical block; 327
 however, 16 of the 22 variables were either differential immune 328
 cell counts/subpopulation frequencies or cytokines produced 329
 by immune cells. Given the pathophysiological heterogeneity 330
 of asthma, traditional cell population and cytokine measures 331
 alone are insufficient to describe the disease.³⁹ 332

The OnPLS model explained >70% of the variance in each 333
 metabolic profiling data block (Sphingolipids, Fatty Acids, and 334
 Oxylipins) and 56% of the Metabolomics block. This higher 335
 degree of explained variance can be attributed to the selective 336
 association between these variables and asthma. These targeted 337
 assays were performed to confirm findings from the initial 338
 metabolomics screen.³⁰ While not all targeted metabolites were 339
 originally detected using metabolomics, they represent bio- 340
 logical processes known to be involved in inflammation. This 341
 point is of particular relevance in that it is not the number of 342
 variables in a given data block that is the primary driver but 343
 rather the inherent biological content.^{4,9} This facet makes it 344
 meaningful to combine disparate -omics blocks of varying 345
 structure into a single OnPLS model and demonstrates the 346
 utility of this approach for data modeling. However, there is 347
 the expected caveat that data blocks that contain higher levels 348
 of biological structure will have a concomitant increase in 349
 contribution to the overall OnPLS model. 350

To determine the biological factors associated with each 351
 OnPLS component and data block, model score vectors were 352
 correlated with a number of known biological factors (Figure 353
 2). Component 1 scores from all blocks positively and 354
 significantly ($p < 0.05$) correlated with disease status (healthy 355
 vs asthma), age, and BMI. All blocks, except Fatty Acids, 356
 positively and significantly ($p < 0.05$) correlated with ICS and 357
 OCS dose. Transcriptomics, Fatty Acids, and Clinical scores 358

Table 2. OnPLS Model Summary

component	connection	Transcriptomics	Sphingolipids	Metabolomics	Fatty Acids	Oxylipins	Clinical
1	global	8%	33%	16%	25%	8%	17%
2	local	7%	11%	7%	-	40%	15%
3	local	-	14%	-	-	7%	-
4	global	8%	9%	10%	39%	7%	14%
5	local	6%	6%	9%	-	-	-
6	local	4%	13%	7%	-	-	-
7	local	5%	-	7%	10%	9%	9%
sum		37%	86%	56%	74%	71%	55%

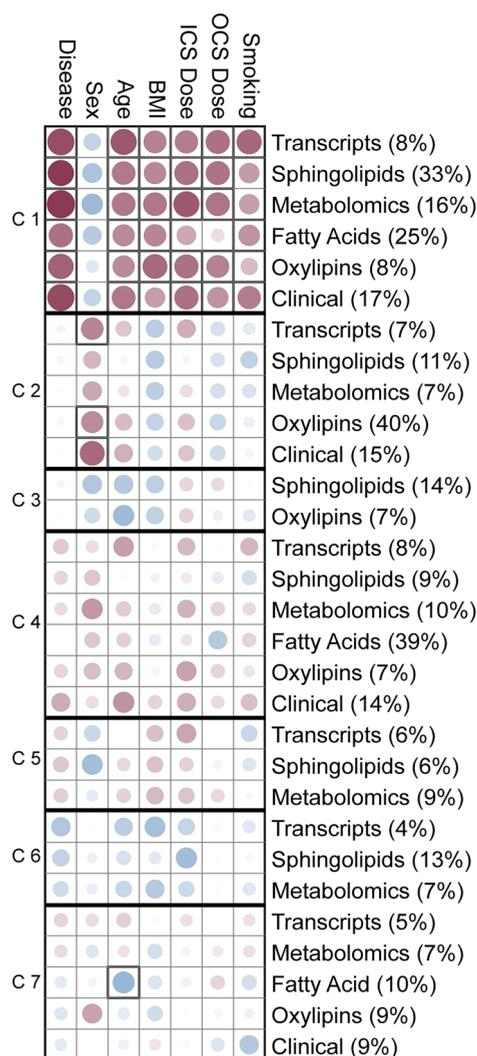


Figure 2. Correlation between model scores and metadata. Circle size and color intensity are proportional to strength of correlation (larger and darker indicates strong correlation). Red, positive correlation; blue, negative correlation. Thick outline around box, significant correlation ($p < 0.05$). The amount of variance that is explained by each data block, in each component, is shown in parentheses. Components are listed as C1–C7 on the left side of the figure.

359 correlated with smoking status (nonsmoker vs has ever
360 smoked). As expected, age, BMI, and corticosteroid treatment
361 were all confounded with disease status (Table 1); thus,
362 explained variation in the model because of these factors was
363 not distinguished from that of disease. The Component 2
364 scores for the Transcriptomics, Oxylipins, and Clinical blocks
365 significantly ($p < 0.05$) and positively correlated with sex.
366 While sex was not a significant confounder in this study, the
367 distribution between the two classes was different. This
368 highlights the utility for OnPLS to identify biological sources
369 for variation in -omics data. The scores for Components 3–6
370 did not correlate significantly with any of the listed biological
371 factors and likely describe either a combination of recorded
372 biological factors or biological factors that were either not
373 observed or not recorded in this study. As such, this highlights
374 the importance of strict experimental design measures and
375 extensive record keeping in data-driven sciences. Despite being
376 a confounder in the study, age negatively correlated with
377 Component 7 Fatty Acid scores and highlights the potential for

OnPLS to identify underlying biology associated with data 378
379 blocks.

PCA of OnPLS Score Vectors. To visualize the entire 380
OnPLS model, principal components analysis (PCA) was 381
performed on the scaled OnPLS score vectors (hierarchical 382
PCA, Figure 3). The first principal component (PC1) showed 383 383

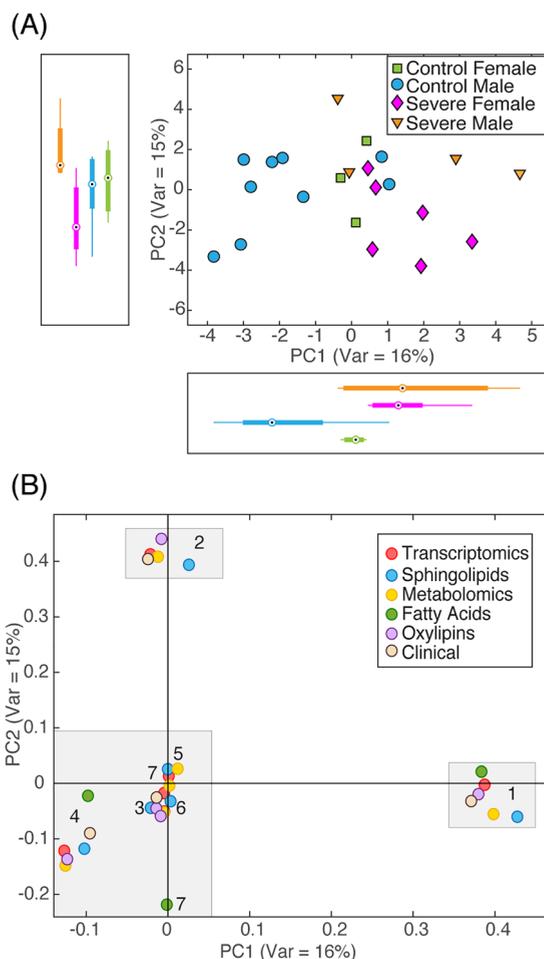


Figure 3. PCA visualization of OnPLS model score vectors. Score vectors from the OnPLS model were scaled to unit variance before performing H-PCA. (A) Score plot. Green squares, control females; blue circles, control males; purple diamonds, severe females; orange inverted triangles, severe males. Bar graphs on axes show distribution of each group along the respective axis. (B) Loadings plot. Red, Transcriptomics; blue, Sphingolipids; yellow, Metabolomics; green, Fatty Acids; purple, Oxylipins; tan, Clinical. Numbers represent OnPLS components, from which score vectors originate. Shaded boxes are for visualization purposes only.

a separation between healthy controls and asthmatic 384
individuals in the score plot (Figure 3A). Aligning with the 385
results of the correlation analysis, this separation was driven by 386
the OnPLS Component 1 score vectors (Figure 3B). It was 387
then expected that PC2 would solely describe a sex difference, 388
as OnPLS Component 2 score vectors drove the separation. 389
Interestingly, PC2 actually described an interaction between 390
disease and sex (Figure 3A). While there was a sex difference 391
among asthmatics, this was not observed in the controls. 392
Investigating the interaction between sex and disease was not 393
an aim of the original cohort study; however, this interaction 394
was identified by simultaneously modeling all the data in 395

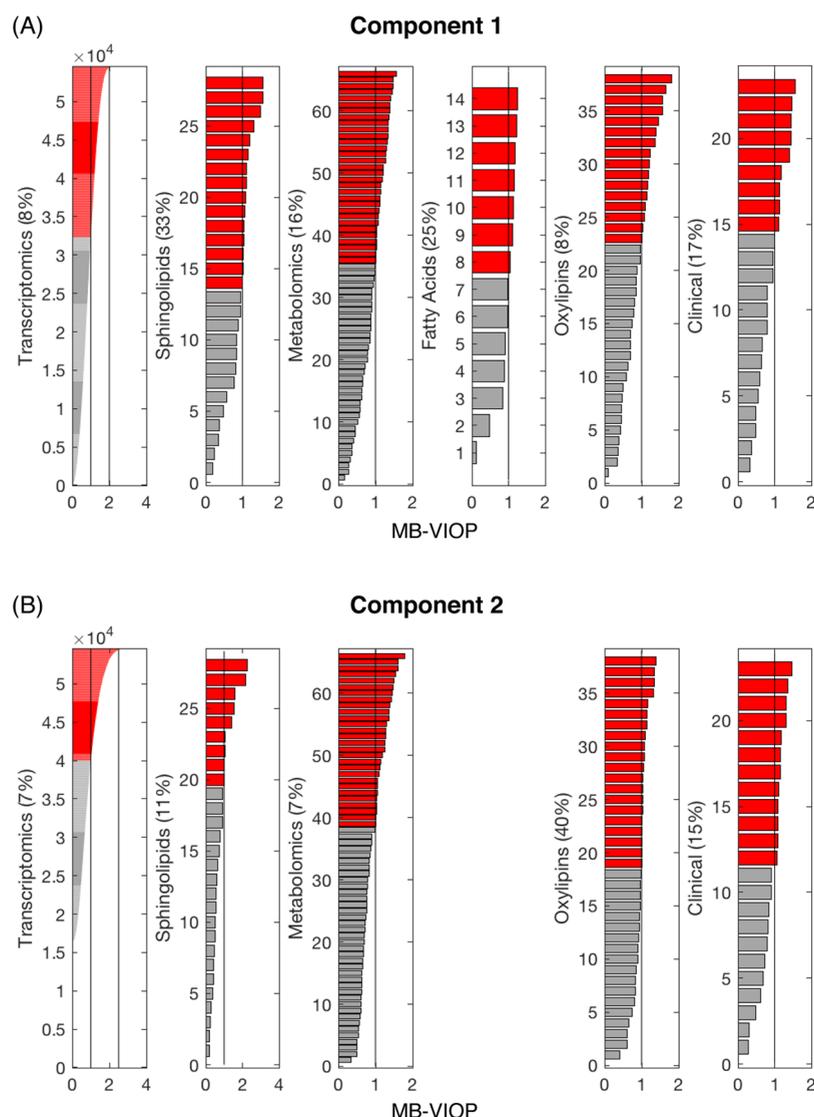


Figure 4. MB-VIOP variable selection for OnPLS Components 1 and 2. The MB-VIOP values are shown for each block in Components 1 and 2. Gray bars, variables with MB-VIOP ≤ 1.0 ; red bars, variables with MB-VIOP > 1.0 . Vertical lines are drawn to show MB-VIOP > 1.0 threshold for all blocks in addition to the increased MB-VIOP thresholds of >2.0 and >2.5 for Transcriptomics in Components 1 and 2, respectively. Percentages reflect the amount of variance described by each component, for each data block. (A) Component 1. (B) Component 2.

396 combination with integrative visualization. In addition, the
 397 hierarchical PCA model corroborates the correlation analysis,
 398 showing that OnPLS Components 1 and 2 contain the most
 399 structural information. Therefore, these components were
 400 selected for further exploration with MB-VIOP analysis.

401 **Multi-block Variable Influence on Orthogonal Pro-**
 402 **jections (MB-VIOP).** To further investigate the variables and
 403 their interactions underlying the shared structure of OnPLS
 404 components 1 and 2, MB-VIOP variable selection and
 405 subsequent correlation analysis were applied. A MB-VIOP
 406 threshold of >1.0 was used to select the variables of interest
 407 from each component. For Component 1, 22 297 transcripts,
 408 31 metabolites, 15 sphingolipids, 7 fatty acids, 16 oxylipins,
 409 and 9 clinical variables contributed to explaining the shared
 410 structure describing disease separation (Figure 4A). For
 411 visualization purposes, the MB-VIOP threshold was increased
 412 to 2.0 for the Transcriptomics data block, leaving 151 variables.
 413 For Component 2, 14 618 transcripts, 28 metabolites, 9
 414 sphingolipids, 20 oxylipins, and 12 clinical variables con-
 415 tributed to explaining the shared structure describing the

interaction between sex and disease (Figure 4B). The
 Transcriptomics block appeared to have a strong influence
 on the disease–sex interaction, with 1487 transcripts passing
 the higher MB-VIOP threshold of 2.0; thus, the threshold was
 further increased to >2.5 to identify only the strongest
 contributions, leaving 203 transcripts. The complete list of
 all MB-VIOP values calculated for Components 1 and 2 is
 presented in Tables S1–S6.

In order to identify between-block biological interactions in
 Components 1 and 2, chord plots were used to visualize
 correlations of variables passing the specified MB-VIOP
 thresholds (Figure 5). This approach revealed a number of
 interesting interactions, of which a selected few are discussed
 as examples of the application of the proposed workflow. Five
 metabolites that correlated with ICS dose³⁰ (cortisol;
 cortisone; dehydroepiandrosterone sulfate, DHEA-S; N-
 palmitoyltaurine, pipercolate) passed the MB-VIOP threshold
 criteria for Component 1. These metabolites correlated with
 the transcripts of 21 unique genes (Figure 5A), of which *ATP6*
VIG1 was particularly interesting. *ATP6 VIG1* has been

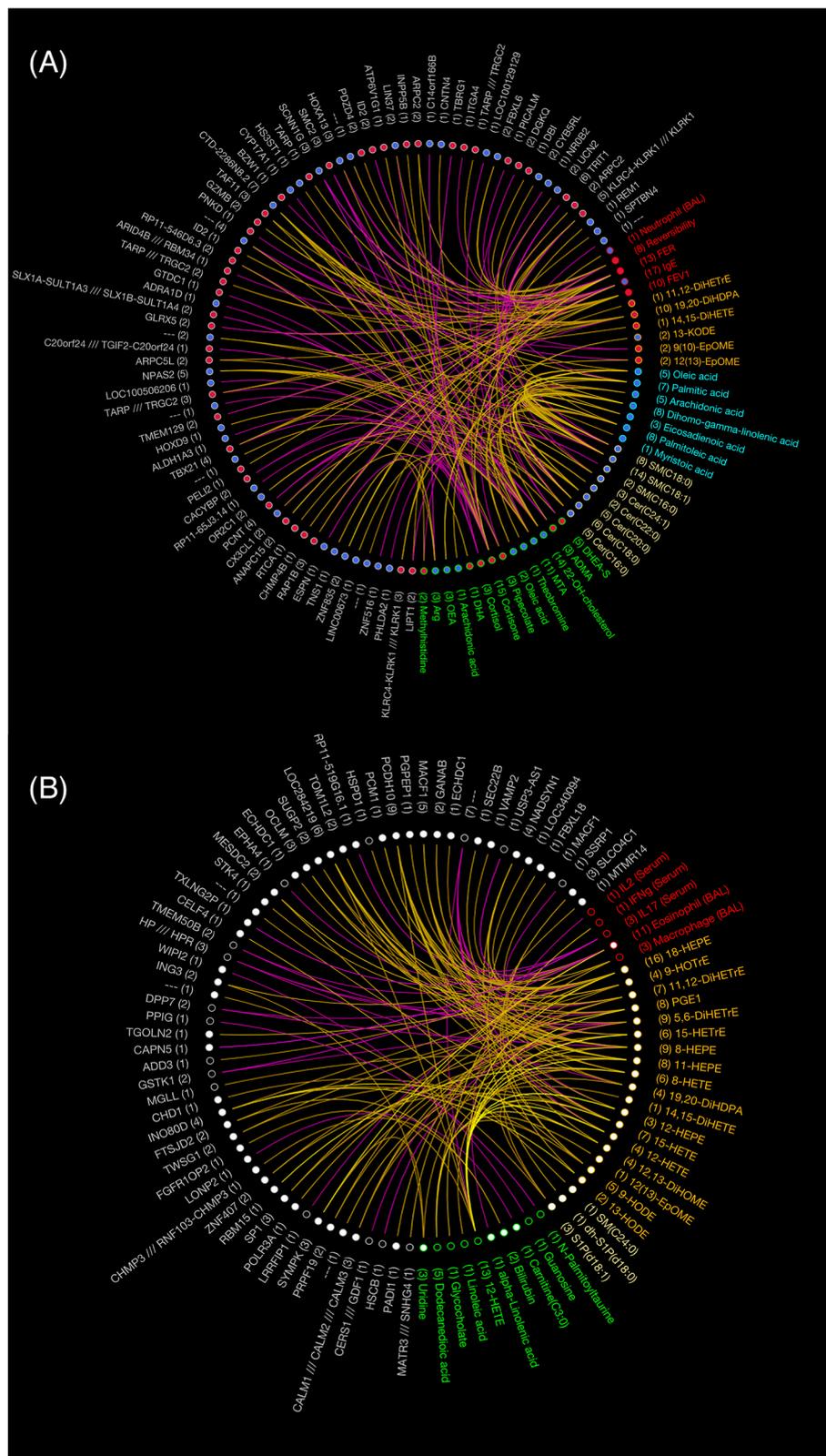


Figure 5. Chord plots showing between-block correlations. (A) Component 1. (B) Component 2. Chord plots were made by calculating the Spearman rank correlations for each pairwise comparison of variables meeting the MB-VIOP thresholds. Variables with a significant ($p < 0.01$) between-block correlation were presented in the chord plots. Nodes represent variables. Text color is associated with block: gray, Transcriptomics; green, Metabolomics; yellow, Sphingolipids; blue, Fatty Acids; orange, Oxylipins; red, Clinical. The number of correlations associated with a given node is noted in parentheses next to the name of the variable. Node color represents direction of change. Component 1: blue, increased in asthma; red, decreased in asthma. Component 2: white, increased in females; black, increased in males. Chords represent correlations: yellow, positive correlation; purple, negative correlation. Each chord plot was constrained such that within-block correlations were ignored. (---) denotes noncoding gene transcripts.

436 implicated in osteoporosis and specifically osteoclast func-
437 tion,⁴⁰ which is a known side-effect of ICS treatment.⁴¹ This
438 novel link may provide insights to the mechanisms underlying
439 bone density loss in asthma patients taking ICS. In addition,
440 *NPAS2*, a transcription factor involved in mediating circadian
441 rhythm,⁴² correlated with five metabolites, four of which were
442 ceramides (Figure 5A). Evidence suggests that ceramide levels
443 fluctuate diurnally;^{43,44} however, to our knowledge, this is the
444 first time an association has been made between *NPAS2* and
445 ceramides. More importantly, as all samples were collected at
446 the same time of day (between 09:00 and 11:00), this supports
447 emerging evidence of dysregulated circadian rhythm gene
448 expression in asthma.⁴⁵ Indeed, experiencing symptoms more
449 than once per week was a classification criterion of severe
450 asthma.²⁹ The disease–sex interaction identified by Compo-
451 nent 2 was largely driven by differential bronchoalveolar lavage
452 cell profiles (eosinophils, macrophages) and oxylipins (Figure
453 5B). It also identified a high degree of correlation between the
454 oxylipins and both *PCDH10* and the uncharacterized gene
455 locus LOC284219, suggesting that these genes may play a
456 previously unidentified role in oxylipin metabolism. Together,
457 these examples highlight the value of this method for
458 interrogating biology and generating hypotheses from
459 multiomics data.

460 By combining OnPLS multi-block modeling with MB-VIOP
461 variable selection and various visualization methods, the
462 composite of data derived from this study could be
463 interrogated. Where methods such as OPLS are useful for
464 identifying covariance in isolated data blocks, OnPLS offers the
465 advantage of identifying combined covariance, thus offering a
466 more complete understanding of the whole system. For
467 example, when OPLS was applied to the Metabolomics data
468 block in isolation, 21 variables had a $VIP_{OPLS} > 1.0$ with
469 dehydroepiandrosterone-sulfate (DHEA-S) being the strongest
470 driver of the control–asthma difference (Supplemental
471 Tables). Component 1 of OnPLS had 31 variables with a
472 MB-VIOP > 1.0 , 15 of which were unique to OnPLS modeling.
473 Whereas DHEA-S was a major driver in the covariance in the
474 single-block analysis, it was less important in the combined
475 covariance of the OnPLS model. The Transcriptomics,
476 Oxylipin, and Clinical data blocks showed similar trends,
477 with OPLS and OnPLS revealing different biological insights
478 (data not shown).

479 While the present study shows the potential for OnPLS-
480 based modeling to be useful for simultaneously modeling
481 multiple data blocks and generating hypotheses, it is limited by
482 sample size and study power. Furthermore, OnPLS is currently
483 unable to derive a block weighting such that MB-VIOP values
484 can be scaled and directly compared across all blocks.
485 Accordingly, MB-VIOP values can only be directly compared
486 within a given data block and not between blocks. In
487 interpreting the results, one must consider the overall
488 contribution, not only of the block per se but also of the
489 individual variables, to the respective component.

490 ■ CONCLUSIONS

491 The multi-block OnPLS method combined with MB-VIOP
492 variable selection and interaction visualization techniques
493 yielded putative novel insights into asthma disease patho-
494 genesis, the effects of asthma treatment, and biological roles of
495 genes. The current study was performed in a worst-case
496 scenario approach using a small sample set, with unbalanced
497 groups and multiple study confounders. While these issues

limit the ability of the different components of the OnPLS 498
model to identify unique biological sources of variation, it 499
demonstrates the potential for this method for identifying key 500
structure in -omics data integration. It is likely that in large 501
well-designed studies, the different components would be able 502
to identify and explain other sources of biological and/or 503
experimental variability (e.g., therapeutics, center bias, diet). It 504
is also possible that this approach would be useful in 505
identifying subphenotypes of disease, with different subgroups 506
and/or mechanisms described by different components. We 507
therefore propose that OnPLS modeling can be incorporated 508
into large-scale molecular phenotyping studies for stratified 509
medicine. Given that the -omics technologies detect molecules 510
that function in a highly interdependent and dynamic manner 511
within a living system, multi-block methods such as OnPLS, 512
together with MB-VIOP and interaction visualization, provide 513
a logical approach to investigating systems biology. 514

■ ASSOCIATED CONTENT

Supporting Information

The Supporting Information is available free of charge on the 517
ACS Publications website at DOI: 10.1021/acs.anal- 518
chem.8b03205. 519

Supporting methods for transcriptomics (PDF) 520

MBVIOP_{OnPLS} and VIP_{OPLS-DA} values for Block 1 521
(Transcriptomics) variables, Table S1; MBVIOP_{OnPLS} 522
and VIP_{OPLS-DA} values for Block 2 (Sphingolipids) 523
variables, Table S2; MBVIOP_{OnPLS} and VIP_{OPLS-DA} 524
values for Block 3 (Metabolomics) variables, Table S3; 525
MBVIOP_{OnPLS} and VIP_{OPLS-DA} values for Block 4 (Fatty 526
Acids) variables, Table S4; MBVIOP_{OnPLS} and VI- 527
P_{OPLS-DA} values for Block 5 (Oxylipins) variables, 528
Table S5; MBVIOP_{OnPLS} and VIP_{OPLS-DA} values for 529
Block 6 (Clinical) variables, Table S6; Rho values for 530
correlations between variables presented in Figure 5A 531
(Component 1), Table S7; *P* values for correlations 532
between variables presented in Figure 5A (Component 533
1), Table S8; Rho values for correlations between 534
variables presented in Figure 5B (Component 2), Table 535
S9; *P* values for correlations between variables presented 536
in Figure 5B (Component 2), Table S10 (XLSX) 537

■ AUTHOR INFORMATION

Corresponding Authors

*E-mail: craig.wheelock@ki.se (C.E.W.) 540

*E-mail: stacey.n.reinke@ecu.edu.au (S.N.R.) 541

ORCID

Johan Trygg: 0000-0003-3799-6094 543

Craig E. Wheelock: 0000-0002-8113-0653 544

Notes

The authors declare no competing financial interest. 545
546

■ ACKNOWLEDGMENTS

The authors wish to thank Rickard Sjögren for providing the 548
script to generate the metadata correlation plot. S.N.R. was 549
supported by a Canadian Institutes of Health Research 550
(CIHR) Fellowship (MFE-135481). T.S.C.H. was supported 551
by Wellcome Trust Research Fellowships (088365/z/09/z and 552
104553/z/14/z), by the Academy of Medical Sciences, and by 553
the National Institute for Health Research (NIHR) Oxford 554
Biomedical Research Centre (BRC). A.S. was supported by the 555

556 Faculty of Medicine, University of Southampton, UK. B.G.P.
557 was supported by MKS Instruments AB, by IDS/KBC of
558 Umeå University (Sweden) for 2016–2017, and by an ERCIM
559 “Alain Bensoussan” Fellowship Programme at the Department
560 of Engineering Cybernetics (ITK) of the Norwegian University
561 of Science and Technology (Norway) for 2017–2018. C.E.W.
562 was supported by the Swedish Heart Lung Foundation (HLF
563 20170603). We acknowledge the support of the Swedish Heart
564 Lung Foundation (HLF 20170734), the Swedish Research
565 Council (2016-02798), the Karolinska Institutet, and the
566 ChAMP (Centre for Allergy Research Highlights Asthma
567 Markers of Phenotype) consortium, which is funded by the
568 Swedish Foundation for Strategic Research, the Karolinska
569 Institutet, AstraZeneca & Science for Life Laboratory Joint
570 Research Collaboration, and the Vårdal Foundation.

571 ■ REFERENCES

- 572 (1) Ideker, T.; Galitski, T.; Hood, L. *Annu. Rev. Genomics Hum.*
573 *Genet.* **2001**, *2*, 343–72.
- 574 (2) Kell, D. B.; Oliver, S. G. *BioEssays* **2004**, *26* (1), 99–105.
- 575 (3) Brown, M.; Dunn, W. B.; Ellis, D. I.; Goodacre, R.; Handl, J.;
576 Knowles, J. D.; O’Hagan, S.; Spasic, I.; Kell, D. B. *Metabolomics* **2005**,
577 *1* (1), 39–51.
- 578 (4) Gromski, P. S.; Muhamadali, H.; Ellis, D. I.; Xu, Y.; Correa, E.;
579 Turner, M. L.; Goodacre, R. *Anal. Chim. Acta* **2015**, *879*, 10–23.
- 580 (5) Wheelock, A. M.; Wheelock, C. E. *Mol. BioSyst.* **2013**, *9* (11),
581 2589–96.
- 582 (6) Nørgaard, L.; Bro, R.; Westad, F.; Engelsen, S. B. *J. Chemom.*
583 **2006**, *20* (8–10), 425–435.
- 584 (7) Broadhurst, D. I.; Kell, D. B. *Metabolomics* **2007**, *2* (4), 171–
585 196.
- 586 (8) Krzanowski, W. J. *Principles of Multivariate Analysis: A User’s*
587 *Perspective*; Clarendon Press, 1988.
- 588 (9) Wold, S.; Sjöström, M.; Eriksson, L. *Chemom. Intell. Lab. Syst.*
589 **2001**, *58* (2), 109–130.
- 590 (10) Hastie, T.; Tibshirani, T.; Friedman, J. *The Elements of*
591 *Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed.;
592 Springer-Verlag: New York, 2009; p 745.
- 593 (11) Breiman, L. *Mach Learn* **2001**, *45* (1), 5–32.
- 594 (12) Cortes, C.; Vapnik, V. *Mach Learn* **1995**, *20* (3), 273–297.
- 595 (13) Bylesjö, M.; Rantalainen, M.; Cloarec, O.; Nicholson, J. K.;
596 Holmes, E.; Trygg, J. *J. Chemom.* **2006**, *20* (8–10), 341–351.
- 597 (14) Trygg, J.; Wold, S. *J. Chemom.* **2002**, *16* (3), 119–128.
- 598 (15) Madsen, R.; Lundstedt, T.; Trygg, J. *Anal. Chim. Acta* **2010**,
599 *659* (1–2), 23–33.
- 600 (16) Xia, J.; Broadhurst, D. I.; Wilson, M.; Wishart, D. S.
601 *Metabolomics* **2013**, *9* (2), 280–299.
- 602 (17) Wold, S.; Johansson, E.; Cocchi, M. PLS Partial Least Squares
603 Projections to Latent Structures. In *3D QSAR in Drug Design: Theory,*
604 *Methods, and Applications*; Kubinyi, H., Ed.; Springer, 1993; pp 523–
605 550.
- 606 (18) Galindo-Prieto, B.; Eriksson, L.; Trygg, J. *J. Chemom.* **2014**, *28*
607 (8), 623–632.
- 608 (19) Höskuldsson, A.; Svinning, K. *J. Chemom.* **2006**, *20* (8–10),
609 376–385.
- 610 (20) Cavill, R.; Jennen, D.; Kleinjans, J.; Briede, J. J. *Briefings Bioinf.*
611 **2016**, *17* (5), 891–901.
- 612 (21) Van Loan, C. F. *SIAM J. Numer Anal* **1976**, *13* (1), 76–83.
- 613 (22) Van Deun, K.; Van Mechelen, I.; Thorrez, L.; Schouteden, M.;
614 De Moor, B.; van der Werf, M. J.; De Lathauwer, L.; Smilde, A. K.;
615 Kiers, H. A. L. *PLoS One* **2012**, *7* (5), e37840.
- 616 (23) Lock, E. F.; Hoadley, K. A.; Marron, J. S.; Nobel, A. B. *Ann.*
617 *Appl. Stat* **2013**, *7* (1), 523.
- 618 (24) Löfstedt, T.; Trygg, J. *J. Chemom.* **2011**, *25* (8), 441–455.
- 619 (25) Smilde, A. K.; Westerhuis, J. A.; de Jong, S. *J. Chemom.* **2003**,
620 *17* (6), 323–337.
- (26) Li, C. X.; Wheelock, C. E.; Skold, C. M.; Wheelock, A. M. *Eur. J.*
621 *Respir. J.* **2018**, *51* (5), 1701930. 622
- (27) Löfstedt, T.; Hoffman, D.; Trygg, J. *Anal. Chim. Acta* **2013**,
623 *791*, 13–24. 624
- (28) Galindo-Prieto, B. Novel variable influence on projection (VIP)
625 methods in OPLS, O2PLS, and OnPLS models for single-and multi-
626 block variable selection: VIP_{OPLS}, VIP_{O2PLS}, and MB-VIOP methods. 627
Doctoral Dissertation, Umeå University, Umeå, Sweden, 2017. 628
- (29) Hinks, T. S.; Zhou, X.; Staples, K. J.; Dimitrov, B. D.; Manta,
629 A.; Petrossian, T.; Lum, P. Y.; Smith, C. G.; Ward, J. A.; Howarth, P.
630 H.; Walls, A. F.; Gadola, S. D.; Djukanovic, R. *J. Allergy Clin. Immunol.*
631 **2015**, *136* (2), 323–33. 632
- (30) Reinke, S. N.; Gallart-Ayala, H.; Gomez, C.; Checa, A.;
633 Fauland, A.; Naz, S.; Kamleh, M. A.; Djukanovic, R.; Hinks, T. S.;
634 Wheelock, C. E. *Eur. Respir. J.* **2017**, *49* (3), 1601740. 635
- (31) Vijayanand, P.; Seumois, G.; Pickard, C.; Powell, R. M.; Angco,
636 G.; Sammut, D.; Gadola, S. D.; Friedmann, P. S.; Djukanovic, R. *N.*
637 *Engl. J. Med.* **2007**, *356* (14), 1410–22. 638
- (32) Diez, D.; Wheelock, A. M.; Goto, S.; Haeggstrom, J. Z.;
639 Paulsson-Berne, G.; Hansson, G. K.; Hedin, U.; Gabrielsen, A.;
640 Wheelock, C. E. *Mol. BioSyst.* **2010**, *6* (2), 289–304. 641
- (33) Skotare, T.; Sjögren, R.; Surowiec, I.; Nilsson, D.; Trygg, J. *J.*
642 *Chemom.* **2018**, e3071. 643
- (34) Wold, S.; Kettaneh, N.; Tjessem, K. *J. Chemom.* **1996**, *10* (5–
644 6), 463–482. 645
- (35) Galindo-Prieto, B.; Trygg, J.; Geladi, P. *Chemom. Intell. Lab.*
646 *Syst.* **2017**, *160*, 110–124. 647
- (36) Holten, D. *IEEE Trans Vis Comp Graph* **2006**, *12* (5), 741–
648 748. 649
- (37) Hunter, J. D. *Comput. Sci. Eng.* **2007**, *9* (3), 90–95. 650
- (38) Wang, D.; Bodovitz, S. *Trends Biotechnol.* **2010**, *28* (6), 281–
651 90. 652
- (39) Holgate, S. T.; Wenzel, S.; Postma, D. S.; Weiss, S. T.; Renz,
653 H.; Sly, P. D. *Nat. Rev. Dis Primers* **2015**, *1*, 15025. 654
- (40) Tan, L. J.; Wang, Z. E.; Wu, K. H.; Chen, X. D.; Zhu, H.; Lu, S.;
655 Tian, Q.; Liu, X. G.; Papisian, C. J.; Deng, H. W. *J. Clin. Endocrinol.*
656 *Metab.* **2015**, *100* (11), E1457–66. 657
- (41) Wong, C. A.; Walsh, L. J.; Smith, C. J.; Wisniewski, A. F.;
658 Lewis, S. A.; Hubbard, R.; Cawte, S.; Green, D. J.; Pringle, M.;
659 Tattersfield, A. E. *Lancet* **2000**, *355* (9213), 1399–1403. 660
- (42) McNamara, P.; Seo, S. B.; Rudic, R. D.; Sehgal, A.; Chakravarti,
661 D.; FitzGerald, G. A. *Cell* **2001**, *105* (7), 877–89. 662
- (43) Jang, Y. S.; Kang, Y. J.; Kim, T. J.; Bae, K. *Mol. Biol. Rep.* **2012**,
663 *39* (4), 4215–21. 664
- (44) Gooley, J. J.; Chua, E. C. *J. Genet. Genomics* **2014**, *41* (5), 231–
665 50. 666
- (45) Kenfield, M.; Yu, H.; Ehlers, A.; Xie, W.; Gunsten, S.; Agapov,
667 E.; Horani, A.; Holtzman, M. J.; Brody, S. L.; Haspel, J. *Am J Respir*
668 *Crit Care Med* **2017**, *195*, A5210. 669