

1 A novel approach to identifying marker genes and estimating the 2 cellular composition of whole blood from gene expression profiles

3 **AUTHORS:** Casey P. Shannon^{1,7}, Robert Balshaw^{1,2}, Virginia Chen^{1,7}, Zsuzsanna Hollander^{1,7}, Mustafa
4 Toma³, Bruce M. McManus^{1,4,7,8}, J. Mark FitzGerald^{6,8}, Don D. Sin^{6,7,8}, Raymond T. Ng^{1,5,7,8}, Scott J.
5 Tebbutt^{1,6,7,8}

6
7 **AFFILIATIONS:** ¹PROOF Centre of Excellence, Vancouver, BC, Canada; ²BC Centre for Disease Control,
8 Vancouver, BC, Canada; ³Division of Cardiology, University of British Columbia, Vancouver, BC, Canada;
9 ⁴Department of Pathology and Laboratory Medicine, University of British Columbia, Vancouver, BC,
10 Canada; ⁵Department of Computer Science, University of British Columbia, Vancouver, BC, Canada;
11 ⁶Department of Medicine, Division of Respiratory Medicine, University of British Columbia, Vancouver,
12 BC, Canada; ⁷Centre for Heart Lung Innovation, University of British Columbia, Vancouver, BC, Canada;
13 ⁸Institute for Heart and Lung Health, Vancouver, BC, Canada.

14
15 *Casey P. Shannon: casey.shannon@hli.ubc.ca
16 Robert Balshaw: robert.balshaw@bccdc.ca
17 Virginia Chen: virginia.chen@hli.ubc.ca
18 Zsuzsanna Hollander: zsuzsanna.hollander@hi.ubc.ca
19 Mustafa Toma: MToma@providencehealth.bc.ca
20 Bruce M. McManus: bruce.mcmanus@hli.ubc.ca
21 J. Mark Fitzgerald: mark.fitzgerald@vch.ca
22 Don D. Sin: don.sin@hli.ubc.ca
23 Raymond T. Ng: rng@cs.ubc.ca
24 Scott J. Tebbutt: scott.tebbutt@hli.ubc.ca

25
26 *Author to whom correspondence should be addressed.

27 Key Points

- 28
- We introduce a model that accurately predicts the composition of blood from Affymetrix Gene
29 ST gene expression profiles.
 - This model outperforms existing methods when applied to Affymetrix Gene ST expression
30 profiles from blood.
31

32 Abstract

33 Measuring genome-wide changes in transcript abundance in circulating peripheral whole blood cells is a
34 useful way to study disease pathobiology and may help elucidate biomarkers and molecular mechanisms

35 of disease. The sensitivity and interpretability of analyses carried out in this complex tissue, however,
36 are significantly affected by its dynamic heterogeneity. It is therefore desirable to quantify this
37 heterogeneity, either to account for it or to better model interactions that may be present between the
38 abundance of certain transcripts, some cell types and the indication under study. Accurate enumeration
39 of the many component cell types that make up peripheral whole blood can be costly, however, and
40 may further complicate the sample collection process. Many approaches have been developed to infer
41 the composition of a sample from high-dimensional transcriptomic and, more recently, epigenetic data.
42 These approaches rely on the availability of isolated expression profiles for the cell types to be
43 enumerated. These profiles are platform-specific, suitable datasets are rare, and generating them is
44 expensive. No such dataset exists on the Affymetrix Gene ST platform. We present a freely-available,
45 and open source, multi-response Gaussian model capable of accurately predicting the composition of
46 peripheral whole blood samples from Affymetrix Gene ST expression profiles. This model outperforms
47 other current methods when applied to Gene ST data and could potentially be used to enrich the
48 >10,000 Affymetrix Gene ST blood gene expression profiles currently available on GEO.

49 **Introduction**

50 Measuring genome-wide changes in transcript abundance in circulating peripheral whole blood cells is a
51 useful way to study disease pathobiology [1]. By providing a relatively comprehensive survey of the
52 status of the immune system, peripheral whole blood transcript abundances may help elucidate
53 molecular mechanisms [2]. The sensitivity and interpretability of analyses carried out in this tissue,
54 however, are significantly affected by its dynamic heterogeneity [3]. It is therefore desirable to quantify
55 this heterogeneity, either to account for it or to model interactions that may be present between the
56 abundance of certain transcripts, some cell types, and some phenotypic indication.

57 Accurate enumeration of the many component cell types that make up peripheral whole blood can be
58 costly, however, and may further complicate the sample collection process, beyond a simple complete
59 blood count and leukocyte differentials (CBC/Diffs). Further, the majority of publicly available peripheral
60 whole blood-derived gene expression profiles on the Gene Expression Omnibus [4] do not include any
61 composition information. Accurate quantification of the cellular composition of blood samples from
62 gene expression data without performing additional experiments is useful, allowing for re-analysis of
63 existing public data, for example.

64 Many approaches have been developed to infer the cellular composition of a sample from high-
65 dimensional transcriptomic [3, 5–10] and, more recently, DNA methylation data [11, 12]. Briefly, if \mathbf{X} , \mathbf{W} ,
66 and \mathbf{H} are matrices with entries X_{ij} (observed expression for sample i , gene j), w_{ik} (composition for
67 sample i , cell type k), and h_{kj} (cell type-specific contribution to the observed expression for cell type k ,
68 gene j), then the problem can be stated: having observed \mathbf{X} , we wish to estimate \mathbf{W} , based on the
69 assumed relationship between expression and composition:

70
$$X_{ij} = \sum_{k=1}^K w_{ik} h_{kj} + e_{ij}$$

71 where e_{ij} represents the expression information for sample i , gene j that is not predictable by the cell
72 composition.

73 We further assume that, for each component cell type k , there exists a subset of features X_{ij}^k in \mathbf{X} whose
74 observed expression in sample i is proportional to the relative abundance of cell type k in sample i .

75 More formally:

76
$$X_{ij}^k \propto w_i^k$$

77 These composition-discriminating features are termed marker genes. For such genes, the elements of
78 the H can be derived from omics profiles of isolated cell types obtained from reference datasets, and W
79 estimated by regression [5, 7–13]. Importantly, mapping such marker genes across technology platforms
80 is not always tractable. Not all genes can be readily mapped across gene expression platforms and the
81 values derived from reference datasets may be specific to the platform on which the gene expression
82 was measured. This limits application of these techniques to platforms on which suitable reference
83 datasets exist. Unfortunately, generating such datasets is costly, and they are correspondingly rare.
84 More recently, so-called reference-free approaches have been proposed to address this issue [6, 14].
85 When applied to transcriptomic data, these approaches still require the identification of suitable marker
86 genes for the cell types to be quantified. This selection is of paramount importance to achieve optimal
87 performance. The general strategy for marker selection is to identify genes whose expression in one cell
88 type exceeds that of all other cell types being considered [6], a process that itself relies on reference
89 datasets. In fact, all approaches discussed thus far rely on one of a handful of publicly available
90 reference datasets to derive a basis matrix or identify suitable marker genes [12, 15, 16]. No suitable
91 reference dataset exists on the newer Affymetrix Gene ST platform.

92 Here we propose a new approach that leverages a multi-task learning algorithm to construct a statistical
93 model able to predict the composition of peripheral whole blood from Affymetrix Gene ST expression
94 profiles. We further show that the coefficients of this model can be used to identify suitable marker
95 genes directly, without the need for a reference dataset. Our strategy is readily applicable to other
96 tissues and/or platforms, which would allow for the development of tools to accurately segment and
97 quantify a variety of admixed tissues from their gene expression profiles, to account for cellular
98 heterogeneity across indications or model interactions between gene expression, some cell types and
99 the indication under study. The described model is freely-available and open source, outperforms other

100 current methods when applied to Gene ST data, and could significantly improve our ability to study
101 disease pathobiology in blood by allowing a more complete study of the various components of the
102 immune compartment of blood from whole blood gene expression.

103 **Patients, material, and methods**

104 **Availability of data and materials**

105 The datasets supporting the conclusions of this article are available on the Gene Expression Omnibus
106 (GEO): repositories GSE77344 (RTP cohort samples) and GSE77343 (CHFP samples). The model is made
107 available as a package for the R statistical programming language, distributed under the GNU General
108 Public License v3.0, and is hosted on GitHub: <https://www.github.com/cashoes/enumerateblood>.

109 **Cohorts**

110 We used two large clinical cohorts to train and validate the new statistical model. The Rapid Transition
111 Program (RTP) included prospectively enrolled patients with chronic obstructive pulmonary disease
112 (COPD), presenting either to St. Paul's Hospital or Vancouver General Hospital (Vancouver, Canada).
113 Subjects presenting to the emergency department or those visiting the COPD clinic were approached for
114 consent to participate in the study. Matched genome-wide transcript abundance and DNA methylation
115 profiles were available for 172 samples from this cohort. This data was used for training the model and
116 cross-validation. Complete blood counts, including leukocyte differentials (CBC/Diffs) were available for
117 all blood samples and used as an independent measure of blood composition (excluding lymphocyte
118 subtypes).

119 The chronic heart failure (HF) program (CHFP) included prospectively enrolled HF patients presenting to
120 St. Paul's Hospital or Vancouver General Hospital (Vancouver, Canada). Subjects were approached
121 during their visit to the heart function, pre-transplant, or maintenance clinics, and those who consented

122 were enrolled in the study. A blood sample was collected at the time of enrollment. Genome-wide
123 transcript abundance profiles and complete blood count, including leukocyte differential (CBC/Diffs)
124 were available for 197 HF patients. This data was used to independently validate the performance of the
125 statistical model.

126 Both studies were approved by the University of British Columbia Clinical Research Ethics Board and
127 Providence Health Care Research Ethics Board and confirm to the principles outlined in the Declaration
128 of Helsinki.

129 **Sample processing**

130 For all subjects, blood was collected in PAXgene (PreAnalytix, Switzerland) and EDTA tubes. The EDTA
131 blood was spun down (200 x *g* for 10 minutes at room temperature) and the buffy coat aliquoted out.
132 Both PAXgene blood and buffy coat samples were stored at -80°C.

133 **Transcript abundance**

134 Total RNA was extracted from PAXgene blood on the QIAcube (Qiagen, Germany), using the PAXgene
135 Blood miRNA kit from PreAnalytix, according to manufacturer's instructions. Human Gene 1.1 ST array
136 plates (Affymetrix, United States) were used to measure mRNA abundance. This work was carried out at
137 The Scripps Research Institute DNA Array Core Facility (TSRI; La Jolla, CA). The resulting CEL files were
138 processed using the 'oligo' R package [17].

139 **Dna methylation**

140 For the RTP cohort samples only, DNA was extracted from buffy coat using Qiagen's QIAamp DNA Blood
141 Mini kits. DNA was bisulphate-converted using the Zymo Research EZ DNA methylation conversion kit,
142 and Infinium HumanMethylation450 BeadChips (Illumina, United States) were used to measure
143 methylation status at >485,000 sites across the genome. This work was carried out at The Centre for

144 Applied Genomics (TCAG; Toronto, Canada). The resulting IDAT files were processed using the ‘minfi’ R
145 package [18].

146 **Statistical analysis**

147 Following preprocessing with their respective packages (‘oligo’ or ‘minfi’), the normalized data were
148 batch corrected using the ‘ComBat’ algorithm [19], as implemented in the ‘sva’ R package [20].

149 **1. Model training**

150 Next, we inferred the cellular composition of the RTP cohort blood samples from their DNA methylation
151 profiles using the ‘estimateCellCounts’ function provided by ‘minfi’. This function uses publicly available
152 DNA methylation profiles obtained from isolated leukocyte sub-types to infer the relative abundance of
153 granulocytes, monocytes, B, CD4+ T, CD8+ T and NK cells (details in **Table 1**) with very high accuracy [12,
154 18]. We compared these composition estimates to those obtained from a hematology analyzer
155 (CBC/Diffs) to assess accuracy.

156 We then fit a multi-response Gaussian model using elastic net regression via the ‘glmnet’ R package [21]
157 on the genome-wide transcript abundance data, using the DNA methylation-derived cell proportions as
158 response variables. The multi-response Gaussian model family is useful when there are a number of
159 possibly correlated responses – a so called “multi-task learning” problem – as is the case for these cell
160 proportions. Probesets with minimum \log_2 expression < 5.5 across all samples (22,251) were excluded
161 using the ‘exclude’ parameter. We set the elastic net mixing parameter ‘alpha’ at 0.1 to encourage the
162 selection of a smaller subset of genes and chose the regularization parameter ‘lambda’ using the
163 ‘cv.glmnet’ function set to minimize mean squared error (MSE).

164 **2. Estimating out-of-sample performance**

165 Out-of-sample performance of our model was evaluated using 20 x 10-fold cross-validation (not to be
166 confused with the cross-validation performed by ‘cv.glmnet’ in order to choose an effective

167 regularization parameter ‘lambda’). We then validated the accuracy and calibration of our model by
168 comparing its predicted cell proportions to the available CBC/Diffs data in the CHFP cohort.
169 Unfortunately, a more complete enumeration of the lymphocyte compartment (e.g., by flow cytometry)
170 was not available in any of our cohorts, so we could not independently validate performance in the
171 various lymphocyte sub-types. Instead, the sum of the predicted B, CD4+ T, CD8+ T and NK cell
172 proportions was compared to total lymphocyte proportions from the CBC/Diffs.

173 **3. Performance compared to other current approaches**

174 Finally, we compared the performance of our model to two alternative approaches for determining the
175 composition of blood samples from their gene expression profiles, described by Abbas *et al.* [5] and
176 Chikina *et al.* [6], in this independent heart failure cohort. First, the basis matrix from Abbas *et al.*,
177 derived from the IRIS (Immune Response In Silico) reference dataset, was used to predict the cell
178 proportions of neutrophils, monocytes, B, CD4+ T, CD8+ T and NK cells [15]. Again, the Abbas predicted
179 proportions for B, CD4+ T, CD8+ T and NK cells were summed to obtain a predicted lymphocyte
180 proportion. The Abbas predicted neutrophil, monocyte and lymphocyte proportions were compared to
181 CBC/Diffs.

182 **4. Model features as marker genes for use with reference-free approaches**

183 Next, we evaluated whether our approach could be used to identify more suitable marker gene sets
184 compared to a reference dataset approach. The reference-free approach described by Chikina *et al.*
185 does not require a basis matrix, relying instead on a set of putative marker genes. These are used to
186 guide the decomposition of the dataset’s covariance structure into separate variance components, using
187 singular value decomposition (SVD). Marker genes for each cell type are summarized in this manner, a
188 technique known as eigengene summarization [22]. Given a good set of marker genes, these
189 summarized values, termed surrogate proportion variables, should track with mixture proportions. We

190 used the reference-free approach described by Chikina *et al.* (as implemented in the ‘CellCODE’ R
191 package) and marker genes derived either from the IRIS reference dataset, as recommended by Chikina
192 *et al.*, or from the coefficients of our model. We then compared the surrogate proportion variables
193 produced by ‘CellCODE’, using either marker gene sets, to those obtained from CBC/Diffs in order to see
194 whether we could identify better marker genes. Spearman’s rank correlation (ρ) was used to summarize
195 association between predictions and root mean squared error of prediction (rmse) was used to
196 summarize accuracy and precision.

197 Results

198 DNA methylation derived predictions of the cellular composition of the RTP cohort blood samples were
199 accurate when compared to those obtained from CBC/Diffs (root mean squared error [rmse] = 0.01 –
200 0.08, Spearman’s ρ = 0.85 – 0.94; **Supplementary Figure S1**). The observed error rates were consistent
201 with those previously reported [11, 12]. These predictions were used as the response variables in a
202 multi-response Gaussian model fit to the RTP cohort gene expression data using an elastic net
203 regression. The model selected by ‘cv.glmnet’ retained 491 features. Its fit to the data is visualized in
204 **Figure 2**, against both the DNA methylation derived composition estimates (**Figure 2A**), and CBC/Diffs
205 (**Figure 2B**). Model fit was good (DNA methylation composition: rmse = 0.01 to 0.04; ρ = 0.86 to 0.97;
206 CBC/Diffs: rmse = 0.01 to 0.06; ρ = 0.91 to 0.97) across all cell types, with the exception, perhaps, of
207 CD8+ T cells. When considering the model fit to the CBC/Diffs data, we noted slight bias, with
208 granulocyte proportions tending to be under-predicted and lymphocyte proportions over-predicted.
209 To characterize the potential performance of this model on new data, we carried out a 20 x 10-fold
210 cross-validation. Estimated out-of-sample performance varied across cell types (**Figure 3**). We report the
211 mean rmse (scaled to the expected cell abundance) and Spearman’s ρ across all 200 generated models.
212 Scaled rmse was lowest for granulocytes (0.08) and monocytes (0.24), higher in B, NK and CD4+ T cells

213 (0.51, 0.52, and 0.58, respectively), and highest in CD8+ T cells (1.21). Absolute rmse (0.02 – 0.06)
214 compared favorably to other methods for inferring cellular composition of samples from gene
215 expression data [5, 7, 8, 13]. Results for Spearman’s ρ were consistent: highest in granulocytes (0.926),
216 followed by monocytes (0.824), NK cells (0.812), CD4+ T cells (0.785), B cells (0.731), and CD8+ T cells
217 (0.671).

218 Next, we applied the model to gene expression profiles from the CHFP cohort blood samples in order to
219 independently validate the model’s performance. Performance remained good (rmse = 0.02 to 0.09; ρ =
220 0.69 to 0.91; **Figure 4A**), though the bias we previously noted was more pronounced. Prediction of
221 monocyte proportions was significantly worse than that seen in-sample (ρ = 0.69 vs. 0.91) and expected
222 out-of-sample (from cross-validation; ρ = 0.80 vs. 0.91). Comparing performance of this model against
223 another available approach for inferring the composition of whole blood samples from microarray gene
224 expression data [5], we find that our model performs better, with both correlation to CBC/Diffs data and
225 prediction error markedly improved, especially for monocytes (lymphocyte $\text{rmse}_{\text{Abbas}} = 0.28$ vs. $\text{rmse}_{\text{glmnet}}$
226 = 0.09; monocyte $\text{rmse}_{\text{Abbas}} = 0.07$ vs. $\text{rmse}_{\text{glmnet}} = 0.02$; $\rho_{\text{Abbas}} = 0.31$ vs. $\rho_{\text{glmnet}} = 0.69$; **Figure 4B**).

227 Marker genes derived from the coefficients of our model outperformed those derived from the IRIS
228 reference dataset when used to predict cellular composition using the approach proposed by Chikina *et*
229 *al.* (granulocytes $\rho = 0.87$ vs. 0.67, lymphocytes $\rho = 0.84$ vs. 0.78, and monocytes $\rho = 0.73$ vs. 0.32; **Figure**
230 **5**). The marker gene sets showed minimal overlap (granulocytes = 3/51, monocytes = 4/58, B cells =
231 0/55, CD4+ T cells = 0/11, CD8+ T cells = 1/15, NK cells = 6/22).

232 Finally, we applied the model to predict the composition of the RTP cohort blood samples from their
233 gene expression. This is a contrived example, as this information was already available to us, but it
234 serves to illustrate a possible application of the approach. As expected, large differences exist in the
235 proportions of the various cell types between patients given prednisone or not. Patients on prednisone

236 had proportionally lower amounts of monocytes ($p = 2.9 \times 10^{-4}$), B ($p = 6.6 \times 10^{-5}$), CD4+ T ($p = 6.6 \times 10^{-7}$),
237 CD8+ T ($p = 5.0 \times 10^{-10}$) and NK cells ($p = 9.3 \times 10^{-10}$), and proportionally higher amounts of granulocytes
238 ($p = 2.3 \times 10^{-8}$).

239 Discussion

240 We introduce a statistical model for predicting the composition of blood samples from Affymetrix Gene
241 ST gene expression profiles. We demonstrate that this model has suitable performance across all
242 included cell types in cross-validation, and validate its performance in an independent cohort. The
243 training and validation cohorts represent 2 major clinical indications, COPD and CHF, and include
244 patients with various comorbidities, on various drugs, some with strong effects on blood gene
245 expression (e.g., prednisone), suggesting that our model may generalize well and be broadly applicable.
246 All training and validation samples were from older individuals, however, and it may be that this model
247 will not generalize well to pediatric populations. A loss of performance in pediatric population has been
248 noted when using a similar approach with DNA methylation data [23].

249 We also show that platform-specific marker gene sets can be derived without the need for reference
250 datasets of isolated gene expression profiles for the cell types we wish to enumerate. Using marker
251 genes selected from the coefficients of our model in combination with the reference-free approach
252 proposed by Chikina *et al.* resulted in better performance compared to using marker genes derived from
253 isolated leukocyte gene expression profiles obtained on another microarray platform. Interestingly, the
254 reference-free approach performed only slightly worse than our model, although with loss of scale. This
255 suggests that the non-zero coefficient weights of the model (which we used to select marker genes for
256 the various included cell types) can be estimated entirely in the data, and that these marker genes may
257 be context-independent surrogates of cell proportions.

258 More generally, the strategy we adopted to derive our model (and identify suitable marker genes) could
259 be readily applied to other platforms, or tissues of interest. The only requirements are accurate
260 quantification of the cell types of interest across a large cohort with matched omics profiling. For many
261 popular platforms (e.g., RNA-seq), this schema may be more cost effective than sorting and profiling a
262 number of replicates for all cells of interest, particularly when we consider how costs would scale with
263 additional cell types to be quantified. Moreover, for low abundance cell types, obtaining a sufficient
264 quantity to profile may not be feasible, depending on the efficiency of available separation techniques,
265 amount of admixed tissue that can be collected in practice.

266 The lack of independent validation within the lymphocyte sub-types is a limitation, though cross-
267 validation performance was good across all cell types. We believe it is unlikely that poor performance in
268 some or all lymphocyte sub-types would result in good performance when summed and compared to
269 CBC/Diffs. Model fit exhibits some degree of shrinkage (flattening of the plot of predicted vs. observed
270 away from the 45 degree line). This is expected, however, and related to the phenomenon of regression
271 to the mean. Performance in cross-validation was notably worse for CD8+ T cells. This could be because
272 of the preponderance of zero values for this particular cell type. We also note that performance in
273 monocytes drops significantly in the validation cohort. It is unclear why this is, but one possibility is the
274 difference in the distribution of values in the validation cohort (mean monocyte proportion in training:
275 0.073 vs. 0.090 in the validation; $p = 1.39 \times 10^{-7}$). We have observed poor performance of various
276 deconvolution approaches in quantifying monocytes in the past [13, 24]. It might be that circulating
277 monocyte diversity is poorly reflected in our current framework and we may be selecting poor marker
278 genes for this cell type as a result. A similar rationale could be applied to explain the poor CD8+ T cell
279 performance results in cross-validation. Certainly, it offers the opportunity for further exploration of the
280 true complexity of these cell types in peripheral blood.

281 In summary, our freely-available, open source statistical model is capable of accurately inferring the
282 composition of peripheral whole blood samples from Affymetrix Gene ST expression profiles. The
283 strategy we adopted to derive this model is readily applicable to other tissues and/or platforms, which
284 would allow for the development of tools to accurately segment and quantify a variety of admixed
285 tissues from their gene expression profiles, to account for cellular heterogeneity across indications or
286 model interactions between gene expression, some cell types and the indication under study. The
287 described model outperforms other current methods when applied to Gene ST data and significantly
288 improves our ability to study disease pathobiology in blood. We provide the opportunity to enrich the
289 >10,000 Affymetrix Gene ST blood gene expression profiles currently available on GEO, by allowing a
290 more complete study of the various components of the immune compartment of blood from whole
291 blood gene expression.

292 **Acknowledgements**

293 The authors would like to thank the research participants without whose tissue donations none of this
294 work would be possible. Additional thanks to study nurses and clinical coordinators for their
295 contributions to patient recruitment and data collection. The authors would also like to acknowledge Dr.
296 Karen Lam for her insightful comments and discussion.

297 Funding agencies: Genome Canada, Genome British Columbia, Genome Quebec, Canadian Institutes of
298 Health Research, Providence Health Care, St. Paul's Hospital Foundation, and PROOF Centre

299 **Authorship and Conflict-of-Interest Statements**

300 Collected data: MT, BMM, JMF, and DDS. Designed the research: CPS, RB, RTN, and SJT. Analyzed and
301 interpreted data: CPS, VC, and ZH. Performed the statistical analysis: CPS. Wrote the manuscript: CPS.

302 The authors declare no conflict of interest.

303 **References**

- 304 1. Chaussabel D: Assessment of immune status using blood transcriptomics and potential implications
305 for global health. *Semin Immunol* 2015, 27:58–66.
- 306 2. Li S, Roupahel N, Duraisingham S, Romero-Steiner S, Presnell S, Davis C, Schmidt DS, Johnson SE,
307 Milton A, Rajam G, Kasturi S, Carlone GM, Quinn C, Chaussabel D, Palucka AK, Mulligan MJ, Ahmed R,
308 Stephens DS, Nakaya HI, Pulendran B: Molecular signatures of antibody responses derived from a
309 systems biology study of five human vaccines. *Nat Immunol* 2013, 15:195–204.
- 310 3. Shen-Orr SS, Gaujoux R: Computational deconvolution: extracting cell type-specific information from
311 heterogeneous samples. *Curr Opin Immunol* 2013, 25:571–578.
- 312 4. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, Marshall KA, Phillippy KH,
313 Sherman PM, Holko M, Yefanov A, Lee H, Zhang N, Robertson CL, Serova N, Davis S, Soboleva A: NCBI
314 GEO: archive for functional genomics data sets—update. *Nucleic Acids Res* 2013, 41(Database
315 issue):D991–D995.
- 316 5. Abbas AR, Wolslegel K, Seshasayee D, Modrusan Z, Clark HF: Deconvolution of Blood Microarray Data
317 Identifies Cellular Activation Patterns in Systemic Lupus Erythematosus. *PLoS ONE* 2009, 4:e6098.
- 318 6. Chikina M, Zaslavsky E, Sealfon SC: CellCODE: a robust latent variable approach to differential
319 expression analysis for heterogeneous cell populations. *Bioinformatics* 2015, 31:1584–1591.
- 320 7. Gaujoux R, Seoighe C: Semi-supervised Nonnegative Matrix Factorization for gene expression
321 deconvolution: A case study. *Infect Genet Evol* 2011.
- 322 8. Gong T, Hartmann N, Kohane IS, Brinkmann V, Staedtler F, Letzkus M, Bongiovanni S, Szustakowski JD:
323 Optimal Deconvolution of Transcriptional Profiling Data Using Quadratic Programming with Application
324 to Complex Clinical Blood Samples. *PLoS ONE* 2011, 6:e27156.
- 325 9. Lu P, Nakorchevskiy A, Marcotte EM: Expression deconvolution: a reinterpretation of DNA microarray
326 data reveals dynamic changes in cell populations. *Proc Natl Acad Sci U S A* 2003, 100:10370.
- 327 10. Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, Hoang CD, Diehn M, Alizadeh AA: Robust
328 enumeration of cell subsets from tissue expression profiles. *Nat Methods* 2015, advance online
329 publication.
- 330 11. Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ, Nelson HH, Wiencke JK,
331 Kelsey KT: DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC*
332 *Bioinformatics* 2012, 13:86.
- 333 12. Jaffe AE, Irizarry RA: Accounting for cellular heterogeneity is critical in epigenome-wide association
334 studies. *Genome Biol* 2014, 15:R31.

- 335 13. Shannon CP, Balshaw R, Ng RT, Wilson-McManus JE, Keown P, McMaster R, McManus BM,
336 Landsberg D, Isbel NM, Knoll G, Tebbutt SJ: Two-Stage, In Silico Deconvolution of the Lymphocyte
337 Compartment of the Peripheral Whole Blood Transcriptome in the Context of Acute Kidney Allograft
338 Rejection. *PLoS ONE* 2014, 9:e95224.
- 339 14. Houseman EA, Molitor J, Marsit CJ: Reference-free cell mixture adjustments in analysis of DNA
340 methylation data. *Bioinformatics* 2014, 30:1431–1439.
- 341 15. Abbas AR, Baldwin D, Ma Y, Ouyang W, Gurney A, Martin F, Fong S, van Lookeren Campagne M,
342 Godowski P, Williams PM, Chan AC, Clark HF: Immune response in silico (IRIS): immune-specific genes
343 identified from a compendium of microarray expression data. *Genes Immun* 2005, 6:319–331.
- 344 16. Allantaz F, Cheng DT, Bergauer T, Ravindran P, Rossier MF, Ebeling M, Badi L, Reis B, Bitter H,
345 D’Asaro M, Chiappe A, Sridhar S, Pacheco GD, Burczynski ME, Hochstrasser D, Vonderscher J, Matthes T:
346 Expression Profiling of Human Immune Cell Subsets Identifies miRNA-mRNA Regulatory Relationships
347 Correlated with Cell Type Specific Expression. *PLoS ONE* 2012, 7:e29979.
- 348 17. Carvalho BS, Irizarry RA: A framework for oligonucleotide microarray preprocessing. *Bioinformatics*
349 2010, 26:2363–2367.
- 350 18. Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen KD, Irizarry RA: Minfi: a
351 flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation
352 microarrays. *Bioinformatics* 2014, 30:1363–1369.
- 353 19. Johnson WE, Li C, Rabinovic A: Adjusting batch effects in microarray expression data using empirical
354 Bayes methods. *Biostatistics* 2007, 8:118–127.
- 355 20. Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD: The sva package for removing batch effects and
356 other unwanted variation in high-throughput experiments. *Bioinformatics* 2012, 28:882–883.
- 357 21. Zou H, Hastie T: Regularization and variable selection via the elastic net. *J R Stat Soc Ser B Stat*
358 *Methodol* 2005, 67:301–320.
- 359 22. Langfelder P, Horvath S: WGCNA: an R package for weighted correlation network analysis. *BMC*
360 *Bioinformatics* 2008, 9:559.
- 361 23. Jones MJ, Islam SA, Edgar RD, Kobor MS: Adjusting for Cell Type Composition in DNA Methylation
362 Data Using a Regression-Based Approach. Totowa, NJ: Humana Press; 2015.
- 363 24. Shannon CP, Hollander Z, Wilson-McManus J, Balshaw R, Ng R, McMaster R, McManus BM, Keown P,
364 Tebbutt SJ: White Blood Cell Differentials Enrich Whole Blood Expression Data in the Context of Acute
365 Cardiac Allograft Rejection. *Bioinforma Biol Insights* 2012:49.
- 366

367 **Tables**

368 **Table 1: Description of predicted leukocytes**

Cell name	Abbreviation Used	Description
Granulocytes	Gran	CD15+ granulocytes
Monocytes	Mono	CD14+ monocytes
B cells	Bcell	CD19+ B-lymphocytes
T cells (CD4+)	CD4T	CD3+CD4+ T-lymphocytes
T cells (CD8+)	CD8T	CD3+CD8+ T-lymphocytes
NK	NK	CD56+ Natural Killer (NK) cells

369

370 **Figures**

371 **Figure 1: Schematic representation of the experiment**

372 The model was trained using 172 blood gene expression profiles from the Rapid Transition Program
373 cohort (RTP). For all training samples, cellular composition was first estimated from their DNA
374 methylation profiles (using minfi's 'estimateCellCounts') and then used as the response matrix to fit a
375 multi-response Gaussian model (using glmnet) on the blood gene expression profiles (1). The
376 performance of this model on new data was estimated using cross-validation (2) and confirmed using
377 192 blood expression profiles from the Chronic Heart Failure Program (CHFP), an independent test
378 cohort (3).

379 **Figure 2: Assessing model fit**

380 Predicted proportions from the model are plotted against the DNA methylation-derived cell proportions
381 for each sample in the training data (**A**) or that obtained from CBC/Diffs (**B**). For **A**, linear best-fit line to
382 the data is plotted (blue line) with 95% point-wise confidence interval for fit (grey band) and compared
383 with perfect agreement (red dashed line). For **B**, the sum of the predicted B, CD4+ T, CD8+ T and NK cell
384 proportions is compared to the total lymphocyte proportions from the CBC/Diffs. For each cell type,
385 Spearman's rank correlation (ρ) and the root mean squared error (rmse) are reported.

386 **Figure 3: Cross-validation performance**

387 Distribution of root mean square error (rmse; **A**), scaled to the expected frequency for each cell, and
388 Spearman's rank correlation (ρ ; **B**) for out-of-sample predictions across a 20 x 10 fold cross-validation
389 are visualized using box-and-whisker plots. The mean and 95% CI are shown as a point and range in the
390 center of each boxplot and represent the expected out-of-sample performance.

391 **Figure 4: Our model accurately predicts the cellular composition of blood**

392 **samples and outperforms existing approaches in Affymetrix Gene ST data**

393 Predicted cell proportions are plotted against the cell proportions obtained from CBC/Diffs in an
394 independent dataset (CHFP cohort) for the model (**A**) or using the method from Abbas *et al.* (**B**). The
395 sum of the predicted B, CD4+ T, CD8+ T and NK cell proportions is compared to the total lymphocyte
396 proportions from the CBC/Diffs. For each cell type, Spearman's rank correlation (ρ) and the root mean
397 squared error (rmse) are reported.

398 **Figure 5: Our model identifies better performing marker genes for use with**
399 **reference-free approaches in Affymetrix Gene ST data**

400 Surrogate proportion variables obtained from CellCODE are plotted against the cell proportions
401 obtained from CBC/Diffs in an independent dataset (CHFP cohort). The sum of the surrogate proportion
402 variables obtained for B, CD4+ T, CD8+ T and NK cells is compared to the total lymphocyte proportions
403 from the CBC/Diffs. Marker genes used by CellCODE were derived from the coefficients of the model (A)
404 or using the recommended set of marker genes (B) derived from the IRIS reference dataset. For each cell
405 type, Spearman's rank correlation (ρ) is reported.

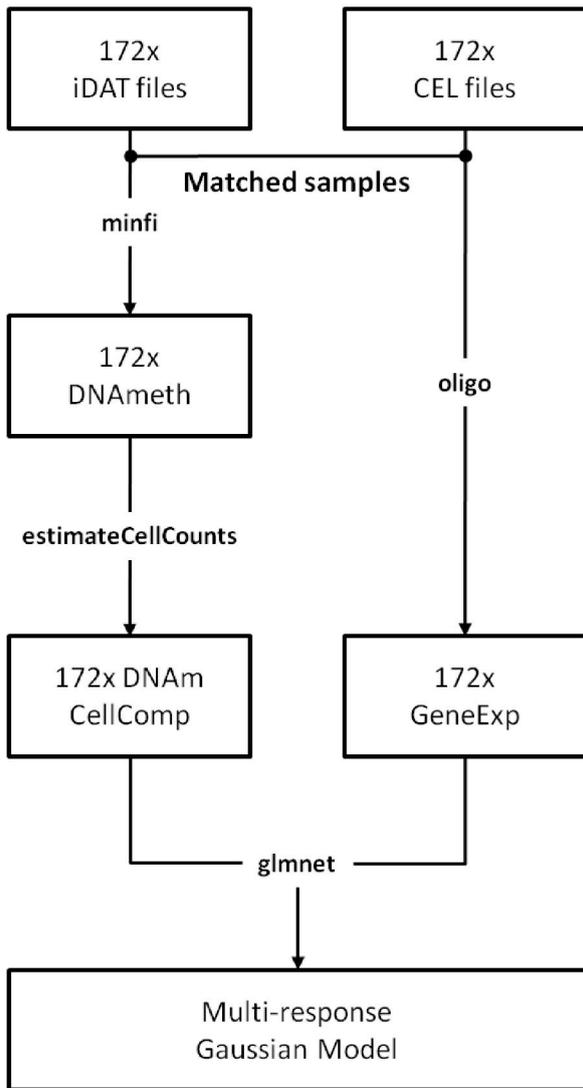
406 **Figure 6: Model predicted cell proportions highlight prednisone-dependent**
407 **changes in peripheral blood composition**

408 Treatment of acute exacerbations (AE) in COPD with prednisone results in important changes in the
409 cellular composition of peripheral blood. The distributions of granulocyte, monocyte, B, CD4+ T, CD8+ T
410 and NK cell proportions are visualized for patients from the Rapid Transition Program (RTP) cohort that
411 were given prednisone or not (p-value is for the unpaired Student's *t*-test comparing the two groups in
412 each case).

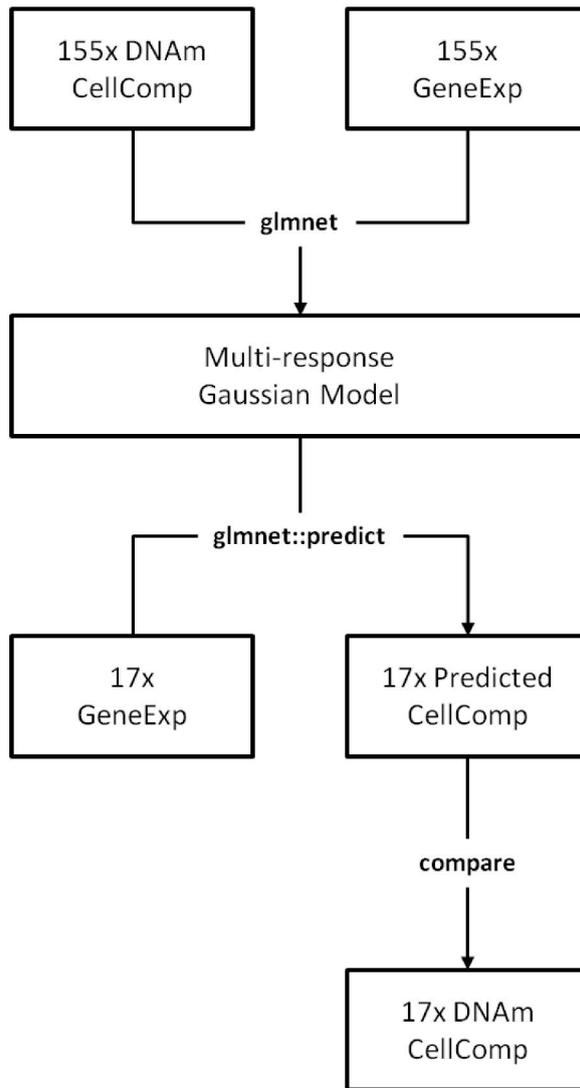
413 **Supplementary Figure S1: DNA methylation-derived composition vs. CBC/Diffs**

414 Predicted proportions were obtained by applying the 'estimateCellCounts' function from the 'minif R
415 package to peripheral blood derived DNA methylation profiles in the Rapid Transition Program (RTP)
416 cohort and plotted against cell proportions obtained from CBC/Diffs. The sum of the predicted B, CD4+
417 T, CD8+ T and NK cell proportions is compared to the total lymphocyte proportions from the CBC/Diffs.
418 For each cell type, Spearman's rank correlation (ρ) and the root mean squared error (rmse) are
419 reported.

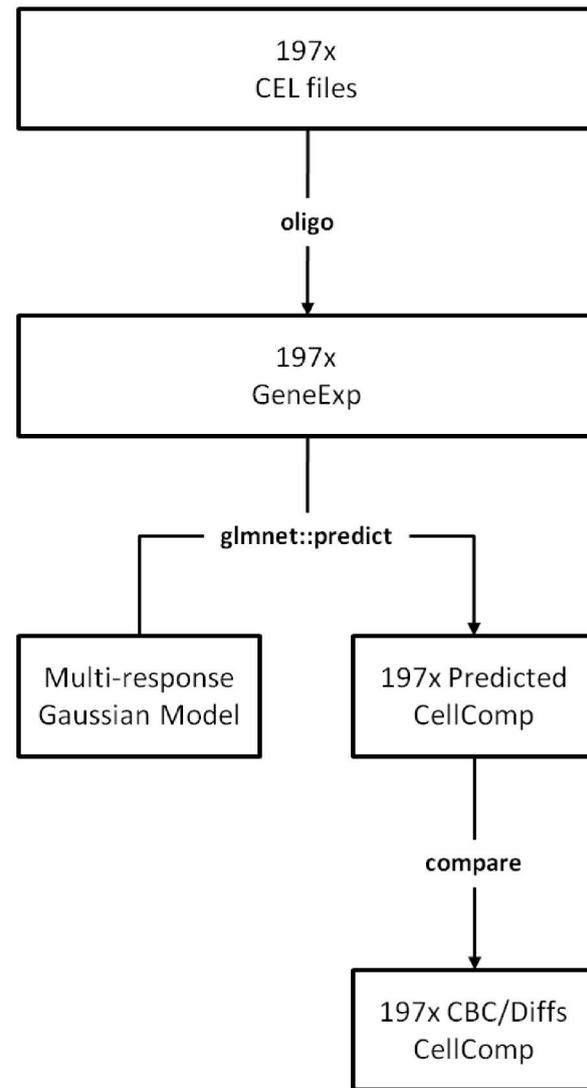
1. Training (RTP Cohort)

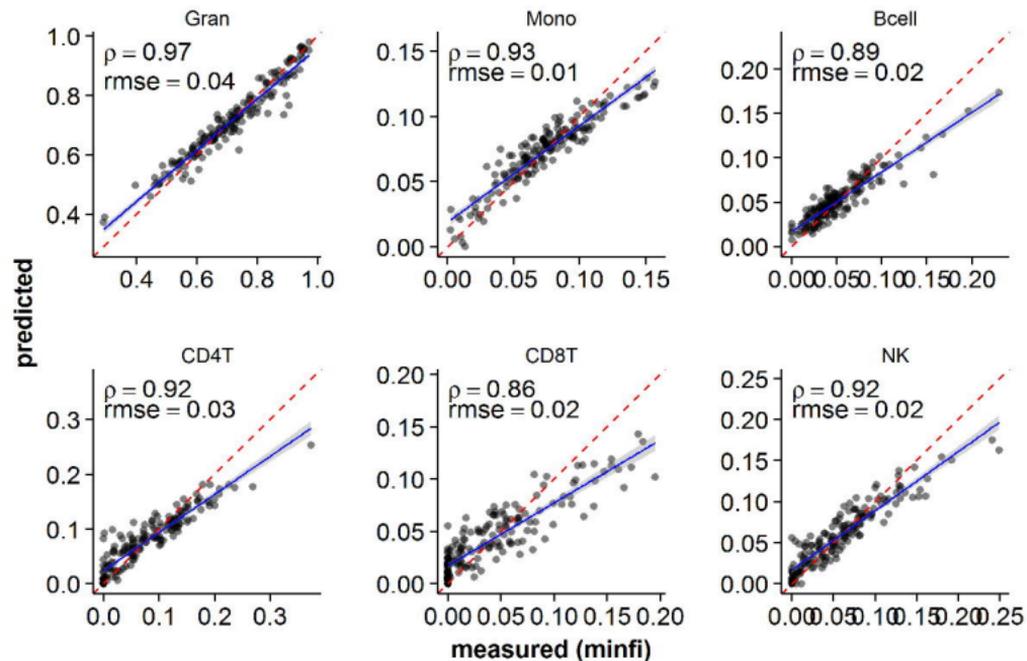
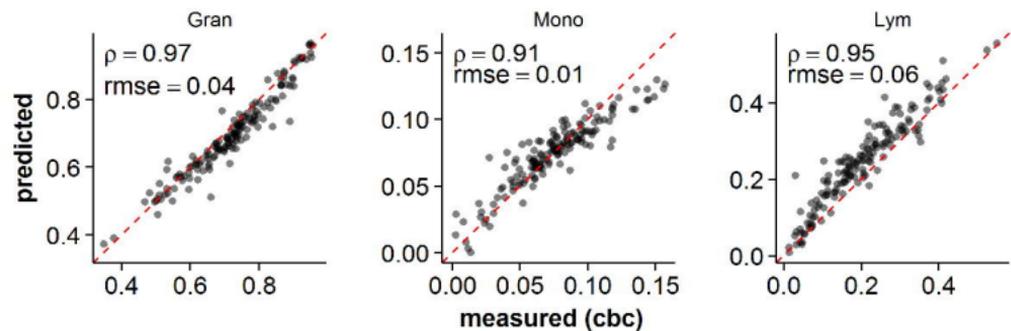


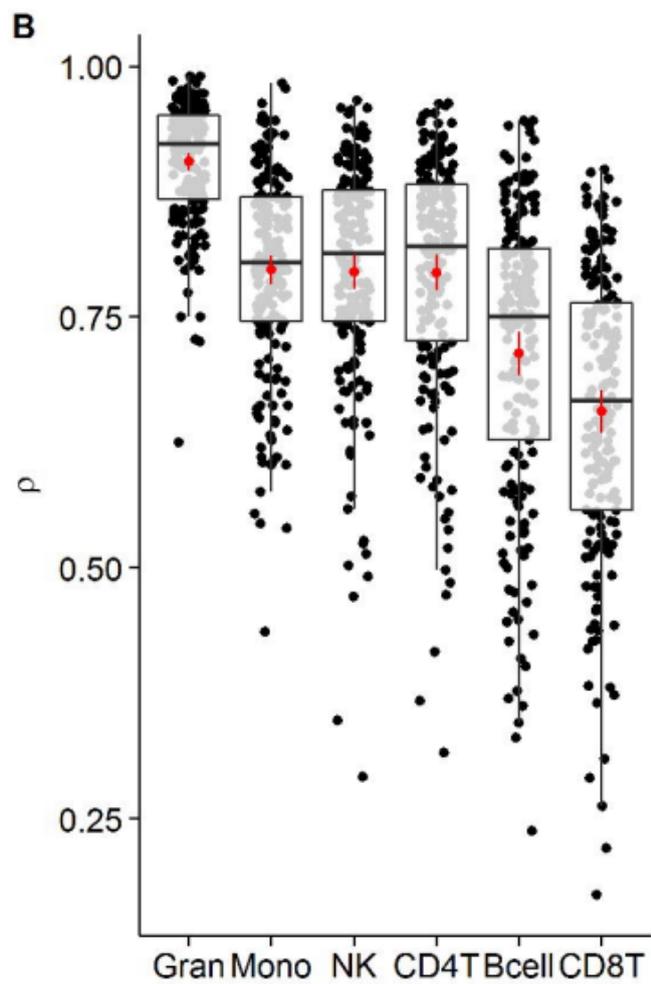
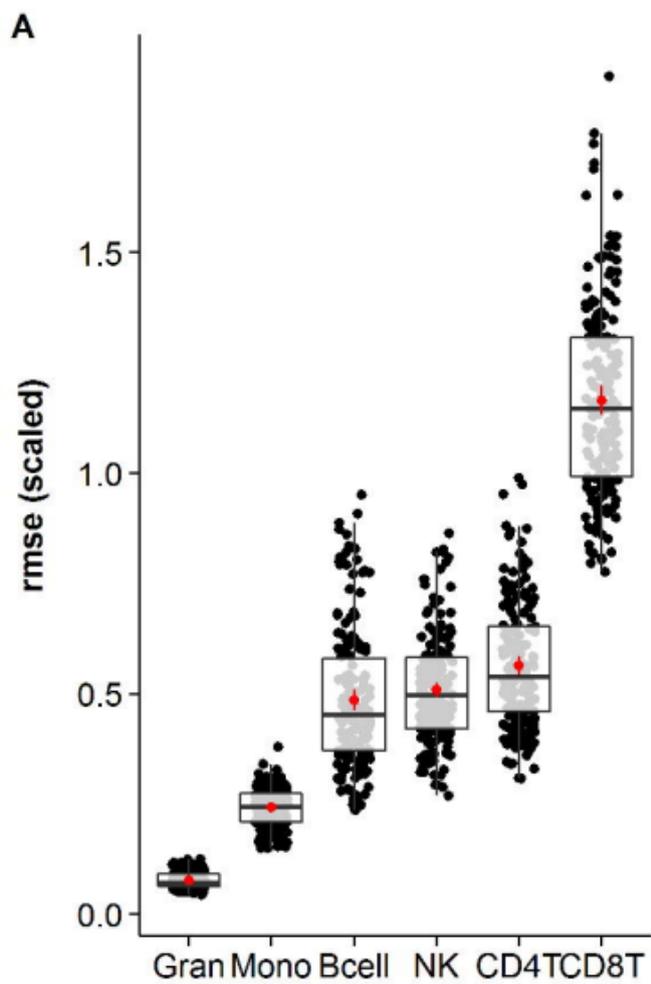
2. Cross-Validation (RTP Cohort)

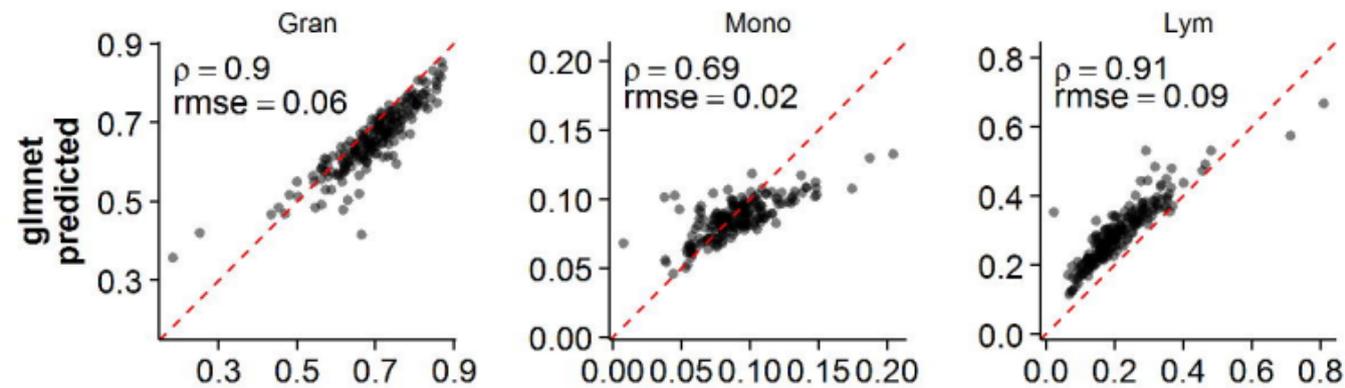
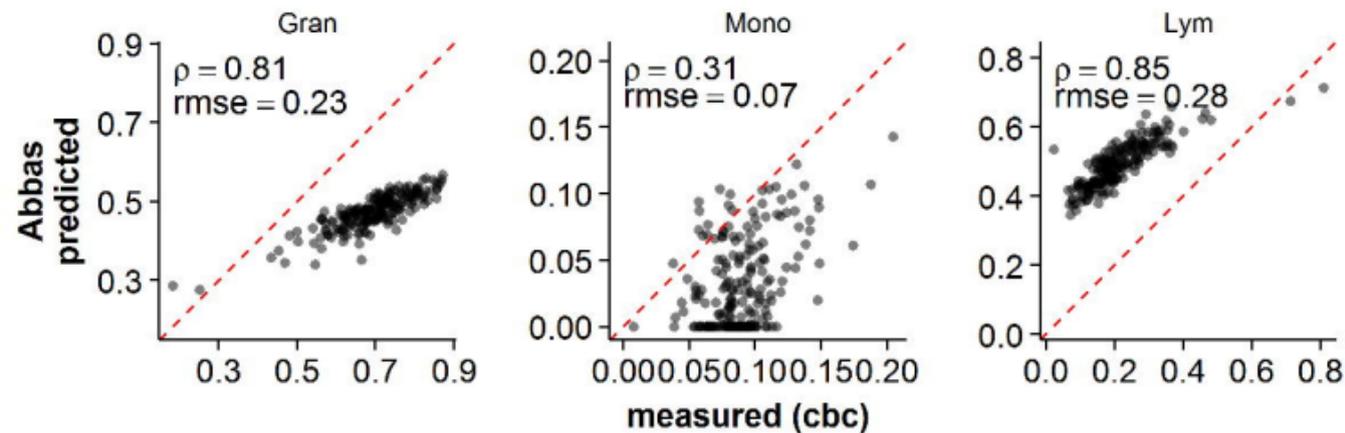


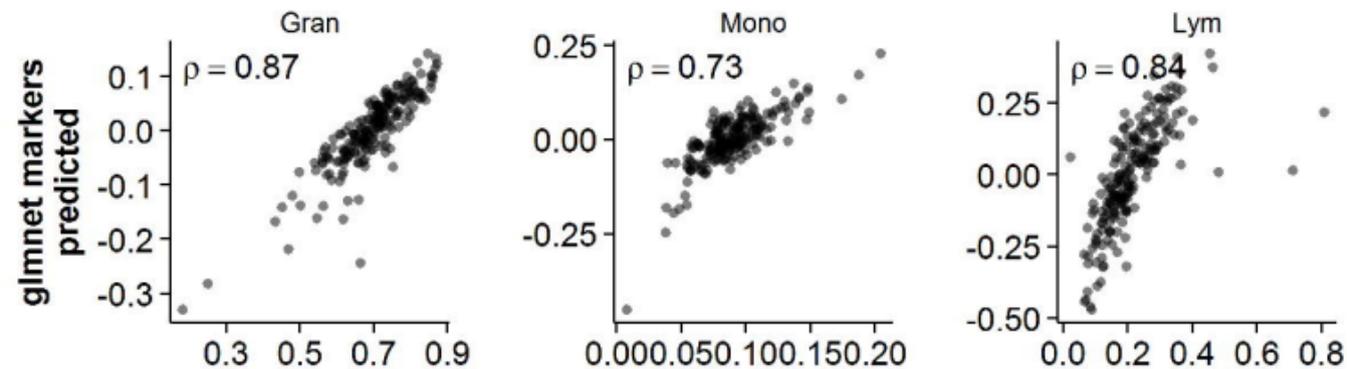
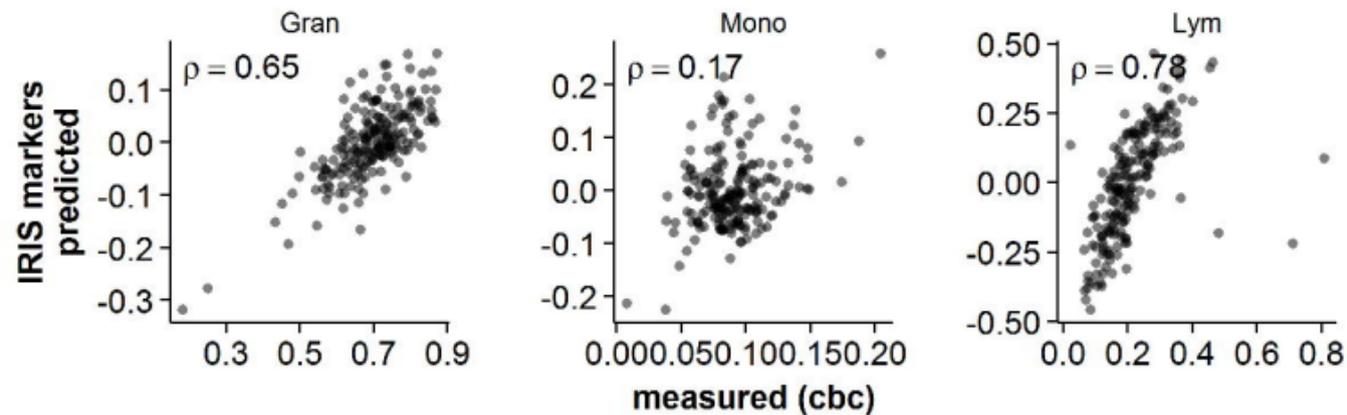
3. Independent Validation (CHFP Cohort)



A**B**



A**B**

A**B**

glmnet predicted

