

1 Identifying robust biomarkers of infection
2 through an omics-based meta-analysis

3
4 Ashleigh C Myall^{1,2}, Simon Perkins¹, David Rushton⁴, Jonathan David⁵, Phillippa Spencer⁶,
5 Andrew R Jones^{1,&}, & Philipp Antczak^{1,3,&*}

6 ¹Institute of Systems, Molecular and Integrative Biology, University of Liverpool, Liverpool,
7 United Kingdom

8 ²Department of Mathematics, Imperial College London, London, United Kingdom

9 ³Center for Molecular Medicine, University of Cologne, Cologne, Germany

10 ⁴Defence and Security Analysis Division, Defence Science and Technology laboratory
11 (DSTL), Porton Down Salisbury, United Kingdom

12 ⁵Chemical, Biological and Radiological Division, Defence Science and Technology
13 laboratory (DSTL), Porton Down Salisbury, United Kingdom

14 ⁶Cyber and Information Systems Division, Defence Science and Technology laboratory
15 (DSTL), Porton Down, Salisbury, United Kingdom

16

17 *Corresponding author

18 E-mail: P.Antczak@liverpool.ac.uk (PA)

19 &Joint Senior Authors

20 **Abstract**

21 A fundamental problem for disease treatment is that while antibiotics are a powerful counter
22 to bacteria, they are ineffective against viruses. To ensure a given individual receives optimal
23 treatment given their disease state and to reduce over-prescription of antibiotics leading to
24 antimicrobial resistance, the host response can be measured to distinguish between the two
25 states. To establish a predictive biomarker panel of disease state we conducted a meta-
26 analysis of human blood infection studies using Machine Learning (ML). We focused on
27 publicly available gene expression data from two widely used platforms, Affymetrix and
28 Illumina microarrays, and integrated over 2000 samples for each platform to develop optimal
29 gene panels. On average our models predicted 80% of bacterial and 85% viral samples
30 correctly by class of infection type. For our best performing model, identified with an
31 evolutionary algorithm, 93% of bacterial and 89% of viral samples were classified correctly.
32 To enable comparison between the two differing microarray platforms, we reverse engineered
33 the underlying molecular regulatory network and overlay the identified models. This revealed
34 that although the exact gene-level overlap between models generated from the two
35 technologies was relatively low, both models contained genes in the same areas of the
36 network, indicating that the same functional changes in host biology were being detected,
37 providing further confidence in the robustness of our models. Specifically, this convergence
38 was to pathways including the Type I interferon Signalling Pathway, Chemotaxis, Apoptotic
39 Processes, and Inflammatory / Innate Response. Amongst and related to these pathways we
40 found three genes, *IFI27*, *LY6E*, and *CDI77*, particularly prevalent throughout our analysis.

41 **Author summary**

42 Bacterial and viral disease require specific treatments, and whilst there are various treatment
43 options for specific infection types, rapid diagnosis and identification of the optimal
44 treatment remains challenging. Even in wealthier countries with developed healthcare

45 systems, unnecessary prescription of antibiotics to patients with viral infections is causing
46 phenomena such as multi-drug resistant bacteria. One way to distinguish a viral from
47 bacterial infection is to measure an individual's responses, for example by measuring the
48 expression of particular genes in a blood sample, as different types of infections trigger
49 different types of responses. In our study we analysed thousands of previously collected data
50 sets from human blood, where individuals had either viral, bacterial or no infection (control).
51 We used machine learning to identify "signatures" – small sets of genes that are indicative of
52 the type of infection (if any) carried by an individual. Within data sets we used two different
53 technology platforms had been used to collect data. We demonstrated that their gene-level
54 signatures do not overlap perfectly when derived from the different platforms, the biological
55 networks from which those genes were derived, however, had a high overlap – giving
56 confidence that our models are robust against technology artefacts or bias. We have identified
57 a small set of genes that serve as strong biomarkers of infection status in humans.

58 **Introduction**

59 The varying differences within both classes of bacterial and viral infections cause the body to
60 respond in a distinct way (1). Bacteria can be countered by pathways such as complement-
61 mediated lysis, and the cell-mediated response for those that survive phagocytosis and live
62 within the cell (intracellular bacteria). In this response, cells present bacterial peptides
63 (antigens) on their surface, which are identifiable by Helper T cells that mediate bacterial
64 destruction (2). There are a large variety of viruses and bacteria that affect the host's immune
65 system in various ways. Whilst some response pathways may overlap for bacterial and viral
66 infections, there are however a number key differences (3, 4). In fact, these different response
67 pathways cause varied transcription (expression) of key genes and are the medium for
68 distinguishing disease state based on the host's transcriptional response (5).

69 Differential expression of certain genes related to immunological responses can be indicative
70 of both (i) disease state and (ii) individual pathogens (6). Such knowledge can be exploited in
71 differentiating between viral, bacterial and control biological states. Previous studies
72 demonstrated this by developing a small set of only seven genes that can accurately
73 discriminate bacterial from viral infections across a range of clinical conditions, whilst
74 simultaneously succeeding to determine with high accuracy which patients do not require
75 antibiotics (7). Simultaneously, there have been numerous other studies looking at diagnosing
76 infection based on the host's transcriptional response (8-12).

77 Previous work failed to generalise as the data contains a far smaller set of pathogens that
78 would be encountered in 'real world' scenarios, or studies focussed on single technology
79 platforms, specific pathogens, or geographical regions (which contain populations with
80 different HLA alleles, and different local pathogen groups). To address this lack of
81 generalisation this work aims to utilise a larger scale analysis over a more representative
82 sample set to improve biomarker generalisability. To gain statistical power and develop more
83 robust panels, meta analyses of publicly available data have proven to be an effective
84 technique (13). However, analysis integrating several cohorts together face inherent
85 limitations from systematic variations otherwise known as "batch effects". Without proper
86 handling, these batch effects have been shown detrimental in population level gene
87 expression analysis (14). Computational techniques exist to reduce batch to batch variation
88 (15). ComBat (16), used in our study, is a well-known batch correction algorithm, and has
89 been shown successful at removing batch effects between studies whilst retaining relatively
90 high amounts of the biological variation.

91 Data-driven identification of robust biomarkers is a much-debated subject in the biological
92 field. Several machine learning (ML) approaches have been proposed, with typically good
93 performance on data sets used in a given study, but poorer performance when biomarkers are

94 taken forward for validation. Important is the distinction between uni- and multi-variate
95 approaches to biomarker discovery. While identifying a single predictive marker might be
96 preferred in theory, multi-variate approaches have enabled the discovery of more complex
97 relationships that can provide performance (accuracy; sensitivity) far exceeding univariate
98 predictive models (17). One particular aspect in multi-variate predictive approaches is the
99 optimisation of the representative model, which rarely can be achieved through brute force
100 testing and relies on feature selection algorithms. In addition, models developed by ML
101 approaches provide a more complete understanding of the underlying biological mechanisms,
102 adding to our understanding of these systems. In this publication we focus on the use of the
103 Random Forest (RF) (18) classifier, which has been demonstrated to perform well in real-
104 world classification problems with high dimensionality and biased data (19). RFs are bagged
105 decision tree models, which classify data points on a subset of features and have been praised
106 for their ability to avoid overfitting (20). Unlike Support Vector Machines or Neural
107 Networks (two frequently used models with high predictive capabilities) RFs forego much of
108 the model selection step using an ensemble approach which builds many weak classifiers into
109 a single strong self-averaging, interpolating model (21). Whilst RFs consist of many weaker
110 models, they have been shown highly effective at capturing non-linear relationships between
111 model predictors and outputs in a number of genomic studies (22, 23).

112 In recent years bioinformatician seeking predictive models have been faced with increasingly
113 greater dimensionality to their data. With the needs of interpretable models many have
114 responded and used feature selection procedures, which aims to remove redundant and
115 irrelevant model features (24). The results of a smaller feature set not only offers improving
116 model performance, faster computational implementation, and greater interpretation of the
117 underlying generative process (25); but moreover lines up with the original pattern
118 recognition theory, that RFs, like many other ML models were not designed to cope with

119 large amounts of irrelevant features, often referred to as *the curse of dimensionality* (26). This
120 high dimensionality is especially pronounced in the case of gene expression data with the
121 total human gene set being ~20,000.

122 Various feature selection procedure exist and have been demonstrated in biological problems
123 (24). For this study we focused on Backwards Elimination (BW) (27) forming a well-
124 established benchmark, and an evolutionary algorithm, a more explorative and
125 parameterizable search approach, to obtain reduced model feature sets (17). BW essentially
126 searches for the optimal feature set by progressively eliminating the least important features
127 from a given dataset and testing whether the new model is significantly more accurate than
128 the previous. Whereas evolutionary algorithms are based on evolving population(s) of models,
129 which are repetitively intermixed, and subject to random point mutations. This evolutionary
130 process is assumed to produce converging model populations in terms of performance and
131 their associated feature sets (28).

132 The application of different computational pipelines often leads to different outcomes in
133 disease prediction (29). We believe, it is thus important not only to present performance
134 statistics for one given model generated by an ML pipeline, but to explore the underlying
135 biological response of a set of plausible models. By doing so, it is possible to develop a more
136 robust biomarker panel (mitigating overfitting which would generally produce models hard to
137 interpret biologically), and to understand why a given model, or set of similar models, are
138 valid.

139 In this work, we have performed a meta-analysis over publicly available transcriptomics data
140 (human blood samples where individuals had bacterial, viral or no infection), from two
141 microarray technologies (Affymetrix and Illumina). We applied feature selection and
142 machine learning for biomarker discovery and predictive model generation, and lastly we

143 explored the biological context of the resulting models by reverse engineering the underlying
144 networks. Representing omics data as a network, has several key benefits. One can often
145 better represent many complex systems as connected components, and the genome is no
146 exception (30). Clustering is one popular method to explore these complex networks and
147 many algorithms exist to reveal insight into these complex structure (31). Visualising a
148 clustered network allows us to explore aspects of this generative process, and how feature
149 selection unfolds over it. However, network construction can often be sensitive to the
150 computational approach and parameterization applied (32, 33). In our approach, we validated
151 our findings and mitigate any potential bias in network generation and clustering by
152 illustrating that the biological driven feature selection is consistent across two separate
153 networks, containing different studies, and derived from different technological platforms.

154 Materials and Methods

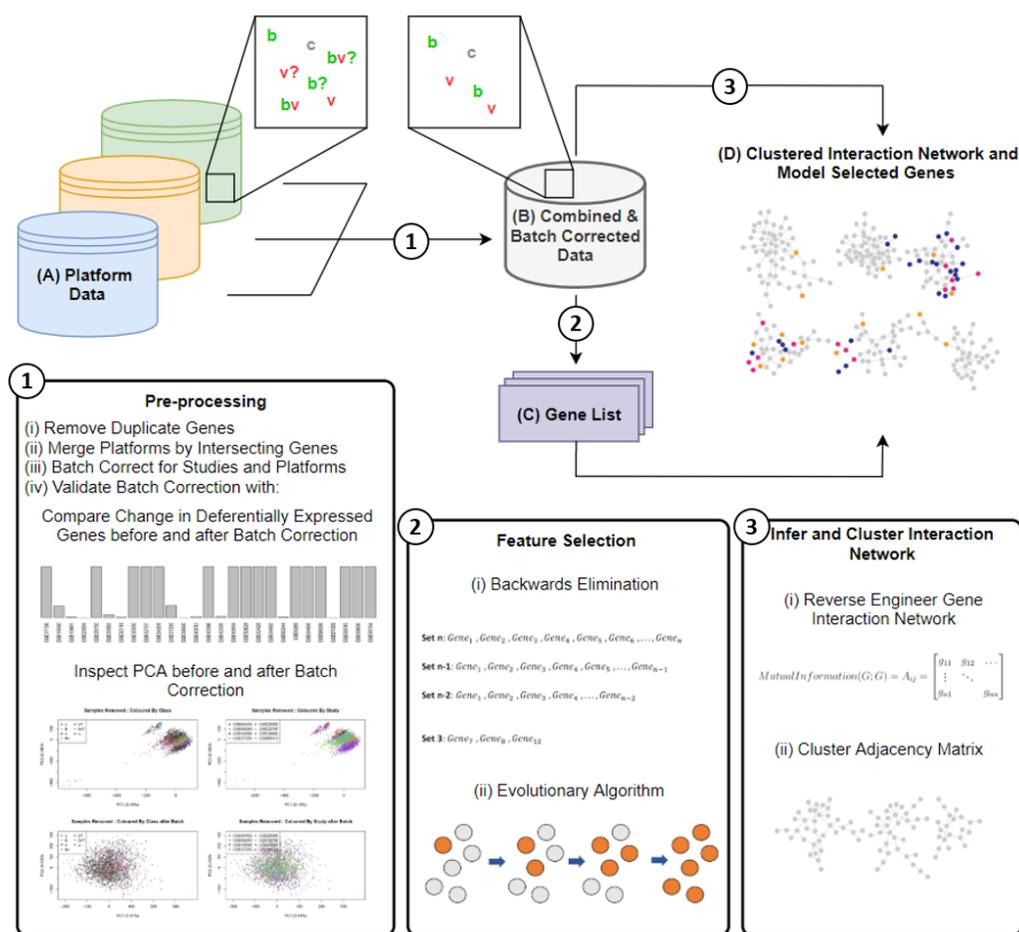


Fig 1. Conceptual overview. Individual data (A), containing bacterial (b), viral (v), control (c), and samples with lower levels of study confidence (?) is merged by common genes and pre-processed by Step 1. Step 1 outputs a combined and batch corrected dataset (B), where only b/v/c samples are present. Two instances of (B) are formed, one where samples of lower levels of support are integrated into b/v classes, and the other completely omitting uncertain samples. Feature selection is performed on data B in Step 2 using (i) Backwards Elimination, and (ii) an Evolutionary algorithm. Step 2's output is a number of Gene Lists (C) obtained in the feature selection. Data B is also used to infer and cluster a gene interaction network, by (i) reverse engineering the gene interaction network, and (ii) clustering the adjacency matrix. (D) is then formed as the clustered interaction network overlaid with genes found in the best performing mode of each dataset and search procedure.

155 To identify and validate a panel of biomarkers able to differentiate bacterial and viral
 156 infections, we performed a meta-analysis of GEO gene expression data, all from open source
 157 microarray human blood infection studies. Our analysis was divided into three major method
 158 steps: i) pre-processing, ii) feature selection, and iii) inferring a gene interaction network, to
 159 discover and our validate gene lists (Fig 1). Following the major steps, we performed and
 160 report the results of a final out of sample test on data not previously used in the training phase
 161 for greater validation.

162 **Pre-processing**

163 **Data.** Datasets from four technological platforms (two from Affymetrix platforms and two
 164 from Illumina platforms), consisting of 3868 samples, from 21 different studies, were
 165 included in the analysis (Table 1). These datasets were selected from a wider pool identified
 166 in an initial scan of online databases, based on a variety of factors including: microarray
 167 platform manufacturer (most prevalent platforms – Affymetrix and Illumina) and study set
 168 size (larger studies with more predictive power), class pathogen strain distribution (aiming
 169 for an equal distribution across the data); and ability to merge with other datasets in our
 170 analysis.

171 **Table 1. Summary of platform level Affymetrix and Illumina datasets prior to pre-processing.**

Manufacturer	Platform (GPL)	Studies (GSE)	Distinct Genes	Sample Count (%)	Bacterial (b) (%)	Uncertain Bacterial (b?) (%)	Viral (v) (%)	Uncertain (v?) (%)	Control (c) (%)	Other (%)
Affymetrix	GPL570	GSE49954, GSE50628, GSE54992, GSE25504, GSE66099, GSE69606, GSE6269, GSE18090, GSE28750, GSE34205	22,213	615 (100)	27 (4.4)	227 (36.9)	164 (26.7)	0 (0)	156 (25.4)	41 (6.7)
Affymetrix	GPL571	GSE52428, GSE95104, GSE17156	13,383	834 (100)	60 (7.2)	0 (0)	358 (42.9)	348 (41.7)	68 (8.2)	0 (0)

Affymetrix	GPL9188	GSE30550	13,383	268 (100)	0 (0)	0 (0)	132 (49.3)	119 (44.4)	17 (6.3)	0 (0)
Illumina	GPL10558	GSE29385, GSE32707, GSE37250, GSE40396, GSE60244, GSE64456, GSE68310	19,947	2,151 (100)	215 (10.0)	141 (6.6)	1069 (49.7)	0 (0)	467 (21.7)	259 (12.0)

172 **Deduplicating genes by probes and merging datasets.** Dataset columns (originally
 173 microarray ProbeIDs) were first deduplicated and substituted by their gene mappings. Where
 174 duplicate ProbeIDs existed for the same gene we selected a representative ProbeID with the
 175 highest average intensity across samples (34). Samples from datasets of the same
 176 manufacturer were then merged by common genes, first at the level of studies within the
 177 same platform, then by platforms in the same manufacturer.

178 **Batch correction and evaluation.** Batch corrections targeted two non-biological sources of
 179 systematic variation: (i) inter-platform study batch effects (differences between platforms),
 180 and (ii) intra-platform batch effects (differences between studies within a batch). Batch
 181 correction was implemented with ‘ComBat’ (16) in a two-step sequential batch correction
 182 pipeline (S1 Appendix.docx). We repeated this process for both Affymetrix and Illumina
 183 datasets separately to form batch corrected Affymetrix data, and batch corrected Illumina
 184 data.

185 Batch correction was verified to retain biological variation and remove technical variation
 186 using two validation steps (Fig 1 Step 1). Firstly, we tested whether pre and post batch
 187 correction significant features overlapped significantly. Secondly, we performed Principle
 188 Component Analysis (PCA) (35) visualising the data in two dimensions and comparing the
 189 PCA plots of before and after batch correction. For a successful batch correction, pre-batch
 190 correction sample clustering in the PCA would be visually removed in the PCA plot of post
 191 batch corrected data.

192 **Dealing with study sample ambiguity: forming a confirmed and integrated dataset**

193 **instance.** To include more data, including some class ambiguity in the original studies, we
194 formed a modelling dataset which integrated bacterial and viral samples with lower levels of
195 confidence (b?, and v?)(Table 1). This integrated dataset contained only classes labeled b/v/c
196 (Fig 1). For Affymetrix this formed (Affy_I) and similarly, for Illumina this formed
197 Illumina_I. Two additional datasets of confirmed sample classes only, were also generated
198 and included in the study but presented only in the Appendix.

199 **Feature Selection.** To search for optimal panels of genes we implemented two search feature
200 selection procedures: (i) the well-known Backward Elimination process (27), and (ii) a
201 genetically inspired search algorithm (GALGO) (17). Both search procedures operated using
202 the RF Classifier, implemented in the R Ranger package (36) a fast and parallelisable
203 implementation of RFs for high dimensional data.

204 **Dataset Preparation.** For dealing with un-even class distributions present in our data (Table
205 1) we employed two strategies. Firstly, we used a study aware data split which insured
206 relatively equal class proportions across both training, test and evaluation data splits.
207 Secondly, we ensured that classification accuracy bias due to larger class proportion of
208 disease states was minimized by weighing smaller classes correspondingly higher (18). This
209 ensures that our model will not be biased to classifying samples with a larger proportion in
210 the dataset.

211 **Backward Elimination.** We operated on a 60/20/20 training/test/evaluation data split for
212 each dataset processed in BW (37). On each training set we ran 240 BW search procedures,
213 using Out-of-bag (OOB) error as the minimisation criterion and implementation using the
214 VarSelRF R package (38). Each run generated a single optimal model which minimised OOB.

215 For each dataset a single representative model was selected from the 240 runs which
216 maximised accuracy on test data.

217 **Genetic-algorithm.** The Genetic-Algorithm (GA) optimized approach is an efficient method
218 for creating suitable multivariate models. We used the R library GALGO (17) to identify a
219 small feature model by continuously crossing a number of small feature models
220 (chromosomes of features) with each other, hypothetically identifying better models with
221 successive generations. We used an initialised fitness goal of 0.95, model size (chromosome
222 size) of 15 genes, and k-fold cross-validation to counter overtraining. In the RF, larger classes,
223 namely viral, were also penalized, as to ensure equal predictions across classes. After 250
224 models, we generated a representative model through a frequency based forward selection
225 strategy which ensures only genes that contributed to predictions are included in the final
226 model (S2 Appendix).

227 **Inferring underlying interaction network**

228 We reverse engineered gene regulatory networks using ARACNe (39) which builds an
229 adjacency matrix of genes with their mutual information from expression data (Fig 1). These
230 networks allow identification of functional relationships between genes and their
231 corresponding products (40, 41). In addition, they can provide insight into the functionally
232 relevant groups of genes for distinguishing disease state, by examining locations of RF
233 selected genes.

234 To select significant interactions within our dataset we used a p-value threshold < 0.05 in the
235 ARACNe procedure. The approach can then estimate a mutual information threshold that is
236 relevant for the provided dataset and a specified p-value. With our data this resulted in a
237 threshold of $MI > 0.0176$ to be retained. From the gene pairs of mutual information, we
238 formed an edge table which was the basis for our interaction network. Nodes are genes and

239 edge weights are the mutual information between two genes, where greater mutual
240 information would suggest a stronger relationship. We then loaded our networks in
241 Cytoscape (42) which visualises molecular interaction networks and has support for a number
242 of clustering algorithms.

243 To identify highly interconnected sub-networks within our reconstructed regulatory network
244 we utilised the Cytoscape clustering plugin GLay (32). GLay uses an implementation of the
245 Girvan-Newman Edge-betweenness algorithm (43) which we used to split our networks it into
246 clusters of connected genes. This resulted in a number of smaller sub-networks and allowed
247 us to inspect their functional roles within the larger network. We then mapped higher level
248 ontologies, such as pathways and gene ontology from gene symbols and used the DAVID (44)
249 tool to provide enrichment analysis. The enrichment analysis looked at several different
250 ontologies, providing an indication of overrepresentation, which we used to infer the likely
251 biological function of a given cluster. Each cluster analysis generated an enrichment table
252 detailing enriched ontology terms along with enrichment ratios and (adjusted) p-values. From
253 the enrichment table we then produced a dotplot which depicted enrichment ratio, p-value
254 and gene count, along with a colour scheme denoting different ontologies, for visual
255 interpretation.

256 For clusters of genes with enriched and significant terms related to the immune response, we
257 labelled them manually as Functionally Relevant (FR) clusters. These FR clusters allowed us
258 to make inferences about which biological functions hold predictive power, by overlaying
259 model selected genes onto our labelled gene regulatory network.

260 **Out of sample testing**

261 Out of sample testing usually refers to testing a model on data not previously seen in model
262 training and selection (37). Whilst a validation set was held back for both Affymetrix and

263 Illumina data, the validation data still contained samples from the same manufacturer and
264 group of studies used in training. Hence, within the original ‘discovery dataset’, gene lists
265 could still be overfit to some non-biological effect persisting in either the manufacturer
266 technology or set of studies present, which was not removed by batch correction.

267 To properly test generalisability and investigate any discovery data bias, we evaluated the
268 best performing models discovered on both Affymetrix and Illumina data by retraining and
269 testing them on non-discovery data (Affymetrix Gene Lists to Illumina Data, and Illumina
270 Gene lists to Affymetrix Data). These non-discovery datasets contained samples from
271 different studies and technology and therefore represented the ideal validation datasets. With
272 similar error between discovery and non-discovery data one can be confident that models
273 have not overfitted to a given dataset and are suggested to be generalisable.

274 **Results**

275 **Pre-processing**

276 Gene de-duplication and data merging was successful for both Affymetrix and Illumina. In
277 the final Illumina datasets 19,947 distinct genes were found intersecting all studies, whereas
278 for Affymetrix Data we found 13,383 (Table 2). This lower Affymetrix count was due to
279 platforms GPL571 and GPL9188 having only 13,383 distinct genes (Table 1). This gene loss
280 from intersection resulted in the omittance of 8,830 gene columns, which were present for the
281 615 samples in GPL570.

282 Affymetrix platforms were successfully merged and combined via our batch correction
283 pipeline, indicated by non-significant changes in differentially expressed (DE) genes and
284 removal of clustering in our PCA analysis between both study and platform batch corrections
285 (S1 Appendix). Illumina based datasets were represented by a single platform, GPL10558.

286 Batch correction did not result in significant changes to DE genes and removed the
287 previously observed clustering by study (S Fig 2).

288 This resulting two datasets Affy_I and Illumina_I contained 1676 and 1892 samples
289 respectively (**Error! Reference source not found.**). It is evident there is an uneven class
290 distribution present in both datasets. Both Affy_I and Illumina_I are made up of more than 50%
291 viral samples (66.89% and 56.50%, Table 2). The most underrepresented class is bacterial
292 samples, with both datasets comprising fewer than 20% samples labelled as bacterial (Table
293 2).

294 **Table 2. Merged and batch corrected modelling dataset description.** Merged and batch corrected
295 Affymetrix and Illumina (ambiguous classes integrated) dataset breakdown by distinct genes,
296 platforms, class make up, and sample count.

Dataset	Distinct Genes	Platforms	Bacterial Samples	Viral Samples	Control Samples	Total Samples
Affy_I	13,383	GPL570 GPL571 GPL9188	314 (18.74%)	1121 (66.89%)	241 (14.38%)	1676
Illumina_I	19,947	GPL10558	356 (18.82%)	1069 (56.50%)	467 (24.68%)	1892

297 **Biomarker lists**

298 Running GA and BW on both Affymetrix and Illumina generated an ensemble of models for
299 each method-datasets pair. For BW this was an ensemble of optimal models, one per run of
300 the algorithm. For GA this was the evolved chromosomes obtained by repeats of the search
301 procedure. From this ensemble of models, we computed relative gene selection frequencies
302 (top 16 genes displayed Table 3).

303 **Table 3. Top 16 Gene selection for Affymetrix and Illumina models and their relative selection**
304 **frequencies.** Frequency provided in brackets is based on the model selection frequency in each
305 optimisation run (the number of times a gene was selected across the number of optimised models).

306 Bold genes are included amongst 3 of models top 16 selection, and underlined genes are included in
 307 all four.

Affymetrix Genes (relative frequency)		Illumina (relative frequency)	
BW	GA	BW	GA
<i>MS4A4A</i> (1.00)	<i>PCOLCE2</i> (1.00)	<i>IFI44</i> (1.00)	<i>IFI27</i> (1.00)
<i>MTHFD2</i> (1.00)	<i>CEP55</i> (0.97)	<i>MCEMP1</i> (1.00)	<i>EPSTI1</i> (0.41)
<i>RSL24D1</i> (1.00)	<i>HBA1.HBA2</i> (0.88)	<i>CD177</i> (1.00)	<u><i>LY6E</i></u> (0.39)
<i>TSPO</i> (1.00)	<i>CDC27</i> (0.66)	<i>GPR84</i> (1.00)	<i>SPATS2L</i> (0.34)
<u><i>LY6E</i></u> (1.00)	<i>TSPO</i> (0.56)	<i>EIF1</i> (1.00)	<i>RSAD2</i> (0.26)
<i>MMP8</i> (1.00)	<u><i>LY6E</i></u> (0.50)	<i>IFI27</i> (1.00)	<i>IFIT5</i> (0.24)
<i>NSUN7</i> (1.00)	<i>MMP8</i> (0.47)	<i>EPSTI1</i> (1.00)	<i>IFI44</i> (0.24)
<i>IFI27</i> (1.00)	<i>PGD</i> (0.47)	<i>REPINI</i> (1.00)	<i>ZDHH19</i> (0.22)
<i>CXCL10</i> (1.00)	<i>RSL24D1</i> (0.47)	<u><i>LY6E</i></u> (1.00)	<i>FCGR1A;FCGR1CP</i> (0.21)
<i>ITGAM</i> (1.00)	<i>SIGLEC1</i> (0.47)	<i>ALKBH5</i> (1.00)	<i>IFI44L</i> (0.19)
<i>PSMA6;KIAA0391</i> (1.00)	<i>IFI44</i> (0.44)	<i>EEF2</i> (1.00)	<i>MCEMP1</i> (0.19)
<i>GRB10</i> (1.00)	<i>OAS3</i> (0.44)	<i>RBM33</i> (1.00)	<i>PRC1</i> (0.18)
<i>GYG1</i> (1.00)	<i>WNT10B</i> (0.44)	<i>ARRB1</i> (0.99)	<i>HPGD</i> (0.17)
<i>PGD</i> (1.00)	<i>ADAMTS3</i> (0.41)	<i>DSCR3</i> (0.99)	<i>OAS2</i> (0.17)
<i>CD177</i> (0.99)	<i>HPR.HP</i> (0.38)	<i>TSPAN18</i> (0.99)	<i>HERC5</i> (0.17)
<i>OLAH</i> (0.99)	<i>OLAH</i> (0.38)	<i>FCGR1A;FCGR1CP</i> (0.96)	<i>IFITM3</i> (0.15)

308

309 BW search procedures in both technologies converged to a small set of genes, indicated by
 310 high relative selection rate calculated by the number of times a gene was selected across the
 311 multiple runs performed in each optimisation procedure. For Affymetrix 14 were included at
 312 a rate of 1.0, whereas for Illumina BW results contain 12 genes at a rate of 1.0 (Table 3).
 313 GA's on the other hand contained a much wider gene selection in the evolved chromosome,
 314 in both search procedures only a single gene was included at a relative rate of 1.0. which
 315 reflects the more varied selection in GA search procedures.

316 Overall search results (aggregated between runs by frequency) from BW and GA in both
 317 Affymetrix and Illumina all contained *LY6E* (Lymphocyte antigen 6E, UniProt: Q16553)
 318 amongst their 9 most frequently selected genes (Table 3). Amongst the next widely selected

319 genes were *IFI27* (Interferon alpha-inducible protein 27, mitochondrial, UniProt: P40305)
320 and *IFI44* (Interferon-induced protein 44, UniProt: Q8TCB0), both in the top 16 by gene
321 selection frequency for three of the four search procedures (Table 3). These 3 genes (*LY6E*,
322 *IFI27*, and *IFI44*) are all type-I interferon-inducible genes (ISGs), demonstrated to have
323 altered expressions in disease states, and known to be highly effective at countering infection
324 (45-48). Furthermore, an additional number of other ISGs were also found amongst the
325 frequently selected model genes (*MS4A4A*, *IFI44L*, *OAS2*, and *IFIT5*).

326 Additionally, several other most frequently selected genes have been linked to certain disease
327 states in the literature. Particularly increased levels of *MMP8* have been observed in HIV-
328 infected patients, which cross-references well as a high proportion of samples in our
329 modelling data coming from HIV viral studies (49). *SIGLEC1* is a Type I transmembrane
330 protein expressed by a subpopulation of macrophages and was one of fifteen genes found
331 upregulated during *in vivo* respiratory syncytial virus infections (50), whilst also said to
332 initiate the formation of the virus-containing compartment (51).

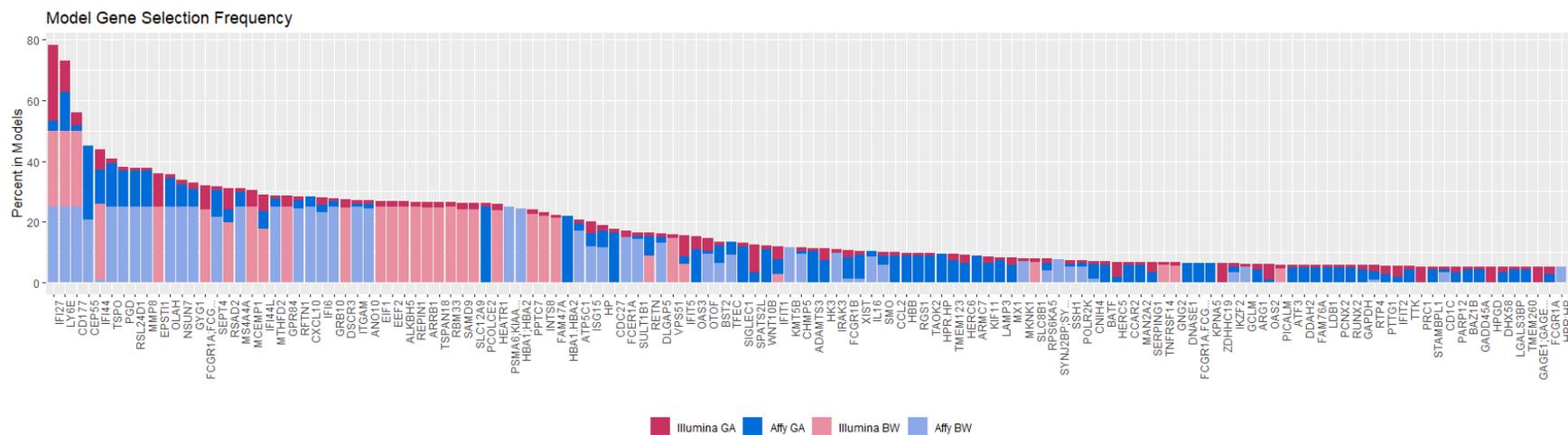


Fig 2. Gene frequency in Affymetrix and Illumina models. Each Model frequency is scaled between 1 and 25. Model overlapping gene frequencies are then stacked and coloured by model-dataset combination. Affymetrix Models by shades of blue and Illumina models by shades of red.

334 To further investigate gene convergence, we compared the relative model gene inclusion
335 rates for all search procedures together. We scaled each model gene frequency (between 1
336 and 25), then plotted them together as a stacked bar plot. Fig 2 shows the resulting stacked
337 frequency, where genes are visualised for greater than 5% aggregated inclusion across all
338 search procedures. Similarly, to our top 16 gene comparison, *LY6E* is indicated as important,
339 being represented in all search procedures. However, interestingly *IFI27* is also included
340 amongst all search procedures. Furthermore *CD177*, a neutrophil-specific receptor and
341 known to be at increased expression for patients in septic shock (52, 53), was selected
342 relatively frequently and present in all search procedures.

343 One interesting aspect to look at is the intersection of this between Genes frequently selected
344 between Affymetrix and Illumina generated models. We identified 88 genes intersecting
345 between Affymetrix and Illumina (S1 Table) and performed functional enrichment analysis
346 of them using DAVID. We found both highly enriched and significant terms relating to the
347 immune response. Included in the list of significant pathways was, in order of significance

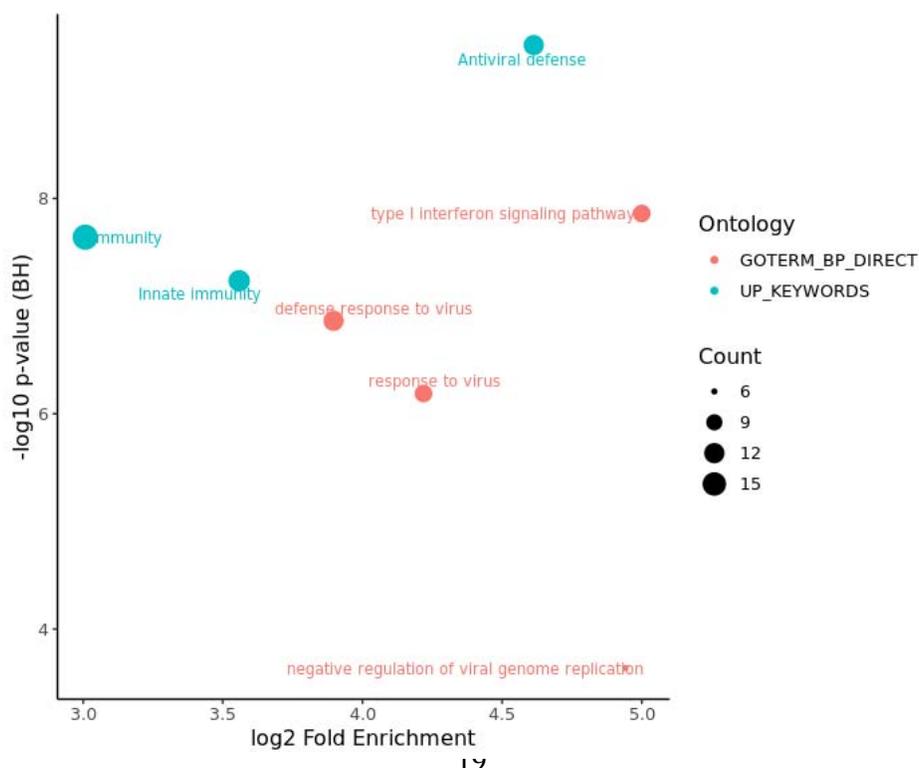


Fig 3. Functional enrichment analysis of the identified 88 genes intersecting between Affymetrix and Illumina search procedures. ‘Antiviral defense’ is the most significant term, whilst ‘type I

348 ‘Antiviral defense’ comprising of 12 genes, the ‘type I interferon signalling pathway’ which
 349 included 10 genes, and ‘Immunity’ encompassing 17 of the 88 genes intersecting between
 350 Affymetrix and Illumina search procedures (**Fig 3**).

351 For each search procedure we obtained a final representative model (Affy_BW, Affy_GA,
 352 Illumina_BW, and Illumina_GA) and evaluated its performance on a held-out data split.
 353 Model performance was recorded as the size of the gene list and its class-based performance
 354 in terms of: Balanced Accuracy, Sensitivity, Specificity, and McNemar’s Test p-value which
 355 tests for consistency in responses and can reveal bias to classifying a certain class (all metrics
 356 derived from the evaluation data split) (54).

357 **Table 4. Overall optimal model performance.** Model performance break down by Affymetrix and
 358 Illumina data sets on the held out test dataset in terms of final model gene size, Balanced Accuracy,
 359 Sensitivity, Specificity, and McNemar’s Test p-value

	Model-Dataset Combination	Gene-set Size	Balanced Accuracy (B/C/V)	Sensitivity (B/C/V)	Specificity (B/C/V)	McNemar’s Test p-value
Affymetrix Identified Models	BW	33	0.94 / 0.78 / 0.86	0.90 / 0.57 / 0.97	0.93 / 0.96 / 0.76	3.57e-3
	GA	36	0.93 / 0.82 / 0.89	0.88 / 0.66 / 0.97	0.99 / 0.97 / 0.81	4.90e-10
Illumina Identified Models	BW	30	0.86 / 0.70 / 0.78	0.80 / 0.47 / 0.87	0.93 / 0.92 / 0.87	2.36e-3
	GA	37	0.82 / 0.58 / 0.89	0.83 / 0.58 / 0.89	0.93 / 0.94 / 0.77	4.33e-15
Average	Average	34	0.89 / 0.72 / 0.89	0.85 / 0.57 / 0.89	0.95 / 0.95 / 0.81	5.93e-3

			0.86	0.93	0.80	
--	--	--	------	------	------	--

360

361 Average model size was similar between both Affymetrix and Illumina models (30-37 genes)
362 (Table 4). On average models classified 0.89 of Bacterial, 0.72 of Control and 0.86 of Viral
363 classes correctly across all datasets. In particular, the Affymetrix models, BW and GA,
364 performed particularly well in terms of balanced accuracy on bacterial samples (0.94 and
365 0.93 respectively). In terms of sensitivity all models performed well for bacterial and viral
366 classes (on average 0.85, and 0.93 respectively), however control sample performance was
367 worse when compared to the viral and bacterial classes (0.57). Evaluating model specificity,
368 bacterial performance was particularly high over all models (averaging 0.95) which would
369 suggest we can determine what a bacterial sample is particularly well regardless of the model
370 used (Table 4).

371 **Inferred interaction networks**

372 We inferred the underlying gene regulatory networks for both Affymetrix and Illumina

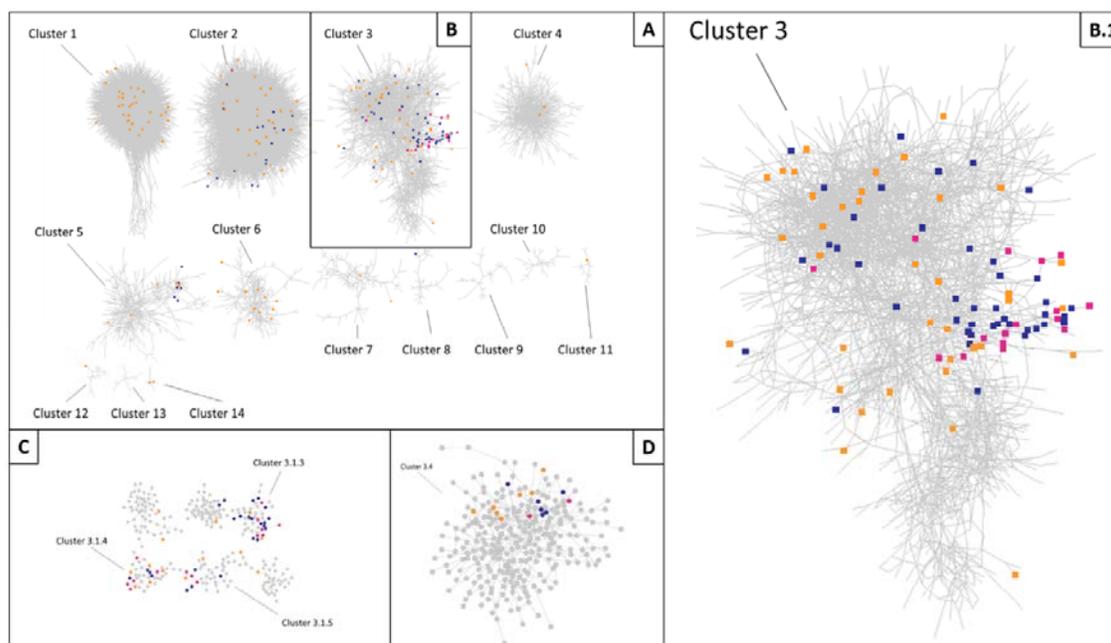


Fig 4. Clustered Illumina interaction network. Illumina models' selected genes are blue, Affymetrix selected genes are orange, and those intersecting both technologies are pink. (A) Illumina Interaction network after initial clustering (visualising clusters > 10 Genes). (B) Cluster 3, containing the most selected genes which intersected between Affymetrix and Illumina models. (B.1) Cluster 3 Enlarged. (C) Highly selected sub clusters of Cluster 3. (D) Cluster 3.4, a sub cluster of Cluster 3 containing two genes which were selected by both Affymetrix and Illumina models.

373 datasets, but present here the analysis on the larger Illumina network (Affymetrix analysis in
374 S3 Appendix). GLayer clustering of the gene interaction network initially revealed 14 clusters
375 containing more than 10 genes (Fig 4). To enable a more granular analysis of specific
376 network sections (those indicated to be functionally relevant in the immune response (FR) as
377 indicated by enrichment analysis, or containing genes selected by our models) we further
378 partitioned several of the initial clusters, forming a network hierarchy (limited to a depth of 3).
379 This resulted in 110 distinct groups of genes which we analysed (Table 5).

380 In Illumina 24 of the 110 clusters had enriched and significant terms related to functions of
 381 the immune system in our DAVID analysis (Table 5). Of these 24 FR clusters, 10 had at been
 382 selected by at least one Illumina model. These 10 clusters contained 55 genes in the union of
 383 Illumina models (68% of all 81 Illumina selected genes in the network). Additionally, a small
 384 number of clusters (four) were selected by every model.

385 **Table 5. Illumina interpreted inferred interaction network properties.** Clusters have been labelled
 386 either functionally related to the immune response (FR). For a cluster to be labelled as FR, functional
 387 enrichment analysis of their gene list will have revealed terms both enriched and significant
 388 implicated in the host response to disease.

Manufac turer	Nodes (Genes)	Sub-clusters of more than 4 Genes	FR Clusters (% of all)	FR Clusters selected by > 1 Model (% of all)	FR Clusters selected by all four Models (% of all)
Illumina	19839	110 (1.00)	24 (21)	10 (9)	4 (4)

389 **Affymetrix – Illumina cluster comparison.** We found a similar number of clusters
 390 converged to between both Affymetrix and Illumina gene lists in their respective networks
 391 (S3 Appendix). It is clear the RF models are selecting genes from multiple of these
 392 uncorrelated clusters, to build a stronger, less correlated model feature sets able to define
 393 disease state. For greater biological understanding we compared the most selected clusters
 394 from both the Affymetrix and Illumina Interaction Network. In Illumina this was Cluster
 395 3.1.3 (S3 Appendix). Whilst the size between both clusters was not comparable (Affymetrix –
 396 Cluster 5 being 435 Genes and Illumina Cluster 3.1.3 being only 47) we found an intersection
 397 of 16 Genes (*DDX60, IFI35, IFI44, IFI44L, IFIH1, IFIT1, IFIT2, IRF7, ISG15, MX1, OAS2,*
 398 *SCO2, TIMM10, TRAFD1, TRIM22 and ZBP1*) which was statistically significant (p-value <
 399 3.18e-12), 10 of which known to be ISGs (*IFI35, IFI44, IFI44L, IFIH1, IFIT1, IFIT2, IRF7,*
 400 *ISG15, MX1, OAS2*) (47). Performing DAVID enrichment analysis on both clusters, we find
 401 in Illumina Cluster 3.1.3 one highly enriched term ‘type I interferon signalling pathway’

402 albeit with a non-significant p-value (S3 Appendix). We do not see the same term in the
403 Affymetrix cluster; however, it does contain numerous ISGs, which we saw commonly
404 amongst gene lists. This convergence between independent feature selection across separate
405 manufacturers and different studies reinforces the high predictive power of ISGs for
406 discriminating disease state across infection studies.

407 **Independent cluster convergence between Affymetrix and Illumina models.** To examine
408 whether convergence between Affymetrix and Illumina was also to the same clusters
409 containing the same genes we looked at where in the Illumina interaction network Affymetrix
410 gene lists selected from (Fig 4, full break down in S3 Appendix). Although selected genes
411 varied between Affymetrix and Illumina, we indeed found that both converged around the
412 same clusters of genes. Moreover, we found that 19 clusters (including lower level sub
413 clusters) were selected by both Affymetrix and Illumina models in the Illumina interaction
414 network. Interestingly amongst this set, the four sub clusters intersecting across all Illumina
415 gene lists (all from within the larger Illumina-Cluster 3: Fig 4) were also selected by
416 Affymetrix gene lists: Illumina-Cluster 3.1.3, Illumina-Cluster 3.1.4, Illumina-Cluster 3.1.5,
417 and Illumina-Cluster 3.4. All of these clusters contained genes revealed by selection
418 frequency analysis in previous section 4.2.

419 We investigated all four clusters selected by all Illumina models (Clusters 3.1.3, 3.1.4, 3.1.5
420 and 3.4) and found they could be separated functionally to different aspects of an immune
421 response. As mentioned, enrichment analysis on Illumina Cluster 3.1.3 revealed the ISGs to
422 be present. However, enrichment analysis also revealed a number of both highly enriched and
423 significant terms related to viral infections ('response to Viruses', 'defense response to
424 virus'), and most prominently 'Antiviral Defense' which is no surprise given the high number
425 of interferon related genes in the cluster (S3 Appendix). Comparing the 47 genes in Clusters
426 3.1.3 to our model frequency analysis revealed 18 overlapping genes (*DHX58*, *EPSTII*,

427 *HERC5, IFI44, IFI44L, IFI6, IFIT1, IFIT2, IFIT5, ISG15, MX1, OAS2, OAS3, RSAD2, RTP4,*
428 *SAMD9, SPATS2L, and TMEM123).*

429 For cluster 3.1.4, in which *LY6E* resides, it bears relation to cell signalling with by far the
430 most significant and enriched term ‘chemotaxis’ (S3 Appendix). Chemotaxis is well known
431 to play critical role in host response to infections, and is specifically involved in recruitment
432 of leukocytes, and movement of lymphocytes around the body (55). The intersect of cluster
433 3.1.4 with our model frequency analysis was also large, being 12 of its 40 genes (*ATF3,*
434 *CCL2, CXCL10, HERC6, LAMP3, LGALS3BP, LY6E, OTOF, PARP12, SEPT4, SERPING1,*
435 *and SIGLEC1).*

436 Cluster 3.1.5 contains genes involved in programmed cell death, containing several
437 significant and enriched terms like ‘Apoptosis’, ‘Regulation of apoptotic process’ and
438 ‘apoptotic process’ (S3 Appendix). A total of 3 of its 37 genes intersected our model
439 frequency analysis (*CHMP5, FCGR1A, and FCGR1B).*

440 Illumina cluster 3.4 contained genes more related to general innate responses with enriched
441 terms containing ‘Inflammatory response’ and ‘innate immune response’ with non-significant
442 p-values (S3 Appendix). Amongst the genes are a number related to the Toll-like receptor
443 family (also an enriched and significant term), which respond to microbial products and
444 viruses, and are key-receptors of the innate immune system (56). Although not visible in the
445 functional enrichment analysis, Illumina Cluster 3.4 also contained a number of Interleukin
446 genes (*IL1B, IL1R1, IL4R, IL18R1, IRAK3*), known to be involved in inflammation and
447 fundamental to innate immunity (57). Out of the 253 genes in cluster 3.4, 15, including
448 *CD177*, intersected with previous frequency analysis (*BATF, CD177, DDAH2, GADD45A,*
449 *GPR84, GRB10, GYG1, HK3, IRAK3, MAN2A2, MKNK1, NSUN7, SULT1B1, TSPO, and*
450 *ZDHHC19).*

451 **Cross manufacturer gene list performance**

452 We evaluated each of the BW & GA representative models from Affymetrix on the Illumina
453 Data and Illumina Models on the Affymetrix data. Contrasting each model's performance
454 between these two discovery and non-discovery datasets we get the performance results
455 depicted in Fig 5. This figure shows the difference between overall accuracy, and class-based
456 accuracy, speciality and sensitivity when generalising our models to data pertaining from a
457 different technology and set of studies.

458 In terms of overall accuracy (Fig 5 A) Affymetrix models, both GA ad BW, performed worse
459 when applying to the Illumina data. However, the drop was less than 0.1 for both Affymetrix

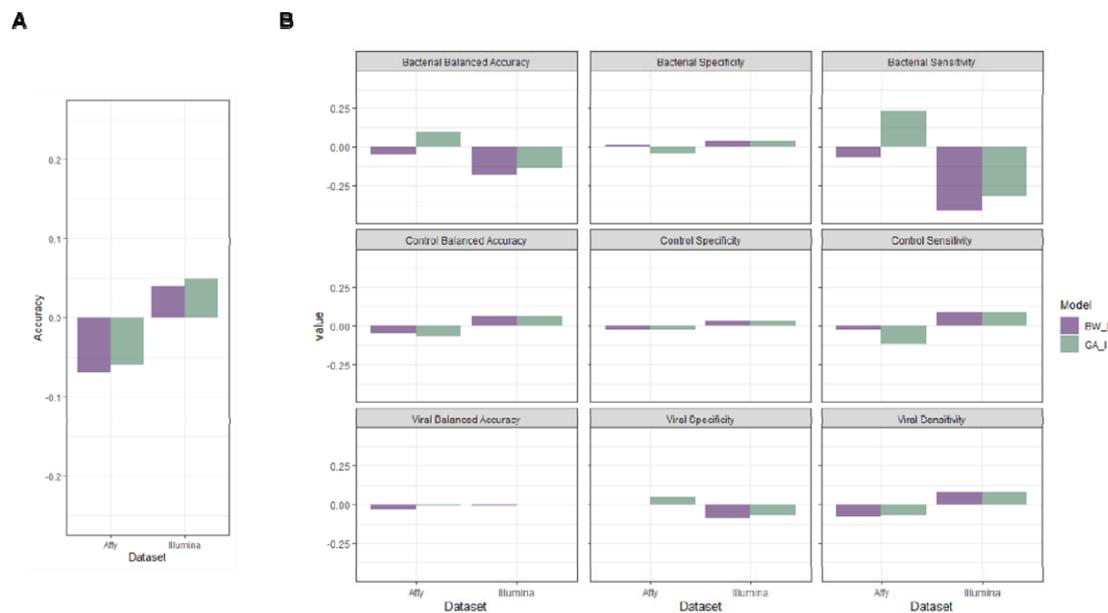


Fig 5. Cross manufacturer model change in performance. Difference in performance when taking Affymetrix derived models and testing on the Illumina data, and the Illumina derived models when testing on the Affymetrix data. (A) Difference in performance in terms of overall Accuracy. (B) Class based performance in terms of Balanced Accuracy, Sensitivity, and Specificity. For each performance measure, bars are grouped by model, and each bar refers to the difference between performance on the original dataset (which each model was discovered on) and the performance on the data it had not been exposed too. For Affymetrix models this would contrast the performance on the Affymetrix data, with the same model's performance on the Illumina data.

460 GA and BW. Whereas for Illumina, both GA and BW models slightly gained accuracy when
461 applied to the Affymetrix data (0.04 and 0.05 respectively).

462 Looking specifically at bacterial performance (Fig 5 B), both Illumina models performed
463 worse on the Affymetrix data in terms of bacterial balanced Accuracy (BW_I 0.71 and GA_I
464 0.73 2dp). Whereas the Affymetrix models performed well on the Illumina data (BW_I 0.89
465 and GA_I 0.89 2dp). In terms of bacterial specificity there was little change for all models,
466 staying within +/- 0.05 2dp of change in performance. However, in terms of bacterial
467 sensitivity, the Illumina models performed particularly worse on the Affymetrix data (BW_I
468 0.44 and GA_I 0.47 2dp).

469 Across viral class specific metrics (Fig 5 B), no model had any large change in Balanced
470 Accuracy (change < 0.05 2dp). The largest metric change was seen in sensitivity, with
471 Affymetrix models slightly decreasing, but with an original score of 0.97 and 0.95 for BW_I
472 and GA_I they are still performing well when ran on the Illumina data.

473 Overall, both Affymetrix and Illumina models performed well given that data was pertaining
474 from different manufacturers and different groups of studies. Particularly stability around
475 viral performance suggests a robustness to the gene lists for classifying viral samples
476 correctly. However, given that bacterial performance change was very comparable to viral, it
477 too suggests a strong ability to classify bacterial samples, even when moving out of the
478 original dataset.

479

480 **Discussion**

481 Due to the amount of relevant data, we focused our analysis on studies from two of the
482 largest microarray platforms, Affymetrix and Illumina. Whilst these both determine the
483 expression levels of genes and are common in large-scale population studies, differences in
484 quantification and normalisation of gene expression values create technical difference (58).
485 Studies within manufacturers were successfully batch corrected, indicated by non-significant
486 changes in differentially expressed genes and removal of sample clustering by studies and
487 platforms in PCA analysis. However, the combination of studies between manufacturers was
488 unsuccessful, leading to two parallel analyses on the combined and batch corrected versions
489 of (i) Affymetrix and (ii) Illumina datasets which minimized biological variation loss.

490 Simpler solutions are more specifically justifiable and allow for greater interpretation, which
491 is the motivation for feature selection amongst models in biological data. We employed two
492 feature selection algorithms using the Random Forest Classifier over our Data: Backwards
493 Elimination and GALGO – both essentially cutting the noise and finding the most significant
494 biological variation responsible for predicting disease state. It is unknown without a brute
495 force search whether a *truly* optimal combination of genes has been found, however both BW
496 and GA approaches converged around a small group of genes located in uncorrelated and
497 functionally separable clusters. Models were found to be strongly enriched for the ISGs. In
498 fact, *IFI27* and *LY6E* (both ISGs) were included in all Affymetrix and Illumina models. *IFI27*
499 is involved in various signalling pathways affecting apoptosis (59-61). Whereas, *LY6E*
500 belongs to a class of interferon-inducible factors that broadly enhance viral infectivity (62).
501 *LY6E* has also been attributed a diverse set of effects, including attenuating T-cell receptor
502 signalling (63) and suppressing responsiveness to *Lipopolysaccharide* which stimulate
503 immune responses (64). Moreover, *IFI27* was shown by Tang et al. to be a *single-gene*
504 biomarker that discriminates between influenza, and other viral and bacterial infections in

505 patients with suspected respiratory infection (65). However, this single-gene biomarker
506 approach lacks generalisability and robustness when predicting a more varied pathogen set.
507 As we have observed, performance in our meta-analysis was greatly improved by including
508 more genes in our models.

509 Our larger set of RF selected genes contained numerous examples confirmed by previous
510 studies to be implicated in disease states. For instance, our results coincide with recent meta-
511 analysis, by Andres-Terre et al., looking at transcriptional signatures of infections,
512 specifically in distinguishing influenza from other viral and bacterial infections, which found
513 127 multi-gene signatures, 27 of which were also present in our representative models (*ATF3*,
514 *BST2*, *CXCL10*, *EIF2AK2*, *HERC5*, *HERC6*, *IFI27*, *IFI44*, *IFI44L*, *IFI6*, *IFIT1*, *IFIT2*, *IFIT5*,
515 *ISG15*, *JUP*, *LGALS3BP*, *LY6E*, *MRPL44*, *MTHFD2*, *MX1*, *OAS1*, *OAS2*, *OAS3*, *OASL*,
516 *RSAD2*, *RTP4*, *SERPING1*, *SPATS2L*) serving to validate our successful data integration and
517 biological findings (66). Notably amongst these coinciding genes are *IFI27* and *LY6E*, again
518 confirming the validity of our converging feature selection.

519 By inferring the underlying interaction network, we discovered that convergence was not
520 only happening to a set of genes, but also, and more prominently, convergence was focusing
521 around particular groups of functionally similar genes. This gene-group convergence only
522 emerged as part of an in-depth investigation into the driving forces of feature selection from a
523 biological network perspective. When representative members of these uncorrelated gene
524 clusters are taken together, they can form highly predictive gene lists. With the ability to
525 define the host response to viral and bacterial infections, genes of our identified clusters are
526 likely good at approximating key functions important in disease state prediction. Notably, the
527 four functional groups of genes were indicated to be: Type I interferon-inducible genes
528 (ISGs), Chemotaxis genes, Apoptotic Processes genes, and Inflammatory / Innate Response
529 genes, which were prevalent in every model (both Affymetrix and Illumina). Within this

530 cluster convergence we found a highly selected group of genes to be ISGs (the most frequent
531 between both Affymetrix and Illumina models). This is no surprise, given Type I Interferons
532 serve as a link between the innate and adaptive immune systems (67) and have a broad range
533 of effects on both innate and adaptive immune cells during infection with viruses, bacteria,
534 and parasites (47). Their varying sensitivity to particular forms of pathogens is likely why a
535 number can be used in conjunction for classification with RFs. While ISGs exact function are
536 not fully understood, it appears our RF models have identified their strong connection to
537 disease state (68, 69). Whilst convergence was prominent around four functional groups of
538 genes, we also note that both in Affymetrix and Illumina, a greater more variable set of
539 functional gene groups were used in addition within our gene lists. Hence, there is a degree of
540 variability in gene solutions, and it seems there is an interchangeable portion of our gene lists
541 in which a number of genes from uncorrelated functional groups of genes can be used to
542 achieve high performance in defining disease state.

543 Finally, we verified our gene lists for generalisability by retraining and evaluating on data
544 from a different manufacturer to which they were discovered in (Affymetrix Gene lists to
545 Illumina and Illumina Gene lists to Affymetrix). It is apparent that all gene lists tend to do
546 better on Affymetrix data, regardless of which set they were discovered on, which suggests
547 that the dataset, not the gene lists, is influencing performance. Hence, we have uncovered the
548 differentiating biological signatures underlying able to define bacterial and viral infections.

549

550

551 **Conclusions**

552 Our meta-analysis of Affymetrix and Illumina human blood infection data has revealed
553 several panels of genes which are able to distinguish well between bacterial and viral
554 infections. The difference in technology and gene coverage between Affymetrix and Illumina
555 did not allow for a direct integration in our analysis. However, we were able to confirm that
556 convergence was occurring independent of the technology, to both the same genes and the
557 same functional groups of genes. This technology independent differentiable signal is
558 learnable, and we demonstrated its presence by reconstructing the underlying regulatory gene
559 network and overlaying models from the two datasets.

560

561 **Acknowledgments**

562 We thank all the contributing studies for generating and making publicly available their
563 respective datasets. We also gratefully acknowledge DSTL (www.gov.uk/dstl) for providing
564 support.

565 This work was also supported by the Chem-Bio Diagnostics program contract HDTRA1-12-
566 D-0003-0023 from the Department of Defense Chemical and Biological Defense program
567 through the Defense Threat Reduction Agency (DTRA).

568

569

570
571

ReferencesBibliography

- 572 1. Shi Z, Gewirtz AT. Together Forever: Bacterial-Viral Interactions in Infection and Immunity.
573 *Viruses*. 2018;10(3):122.
- 574 2. Chaplin DD. Overview of the immune response. *The Journal of allergy and clinical*
575 *immunology*. 2010;125(2 Suppl 2):S3-S23.
- 576 3. Rock KL, Reits E, Neefjes J. Present Yourself! By MHC Class I and MHC Class II Molecules.
577 *Trends Immunol*. 2016;37(11):724-37.
- 578 4. Yewdell JW, JR B. Mechanisms of Viral Interference with MHC Class I Antigen Processing and
579 Presentation. In: *Annual Reviews Collection Bethesda (MD): National Center for Biotechnology*
580 *Information (US);* 2002.
- 581 5. Manger ID, Relman DA. How the host 'sees' pathogens: global gene expression responses to
582 infection. *Current Opinion in Immunology*. 2000;12(2):215-8.
- 583 6. Suarez NM, Bunsow E, Falsey AR, Walsh EE, Mejias A, Ramilo O. Superiority of transcriptional
584 profiling over procalcitonin for distinguishing bacterial from viral lower respiratory tract infections in
585 hospitalized adults. *J Infect Dis*. 2015;212(2):213-22.
- 586 7. Sweeney TE, Wong HR, Khatri P. Robust classification of bacterial and viral infections via
587 integrated host gene expression diagnostics. *Sci Transl Med*. 2016;8(346):346ra91.
- 588 8. Ramilo O, Allman W, Chung W, Mejias A, Ardura M, Glaser C, et al. Gene expression patterns
589 in blood leukocytes discriminate patients with acute infections. *Blood*. 2006;109(5):2066-77.
- 590 9. Hu X, Yu J, Crosby SD, Storch GA. Gene expression profiles in febrile children with defined
591 viral and bacterial infection. *Proceedings of the National Academy of Sciences*. 2013;110(31):12792-
592 7.
- 593 10. Nascimento EJM, Braga-Neto U, Calzavara-Silva CE, Gomes ALV, Abath FGC, Brito CAA, et al.
594 Gene expression profiling during early acute febrile stage of dengue infection can predict the disease
595 outcome. *PloS one*. 2009;4(11):e7892-e.
- 596 11. Zaas AK, Chen M, Varkey J, Veldman T, Hero AO, Lucas J, et al. Gene Expression Signatures
597 Diagnose Influenza and Other Symptomatic Respiratory Viral Infections in Humans. *Cell Host &*
598 *Microbe*. 2009;6(3):207-17.
- 599 12. Dawany N, Showe LC, Kossenkov AV, Chang C, Ive P, Conradie F, et al. Identification of a 251
600 gene expression signature that can accurately detect M. tuberculosis in patients with and without
601 HIV co-infection. *PloS one*. 2014;9(2):e89925-e.
- 602 13. Lagani V, Karozou AD, Gomez-Cabrero D, Silberberg G, Tsamardinos I. A comparative
603 evaluation of data-merging and meta-analysis methods for reconstructing gene-gene interactions.
604 *BMC Bioinformatics*. 2016;17(5):S194.
- 605 14. Akey JM, Biswas S, Leek JT, Storey JD. On the design and analysis of gene expression studies
606 in human populations. *Nature Genetics*. 2007;39(7):807-8.
- 607 15. Benito M, Parker J, Du Q, Wu J, Xiang D, Perou CM, et al. Adjustment of systematic
608 microarray data biases. *Bioinformatics*. 2004;20(1):105-14.
- 609 16. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using
610 empirical Bayes methods. *Biostatistics*. 2007;8(1):118-27.
- 611 17. Trevino V, Falciani F. GALGO: an R package for multivariate variable selection using genetic
612 algorithms. *Bioinformatics*. 2006;22(9):1154-6.
- 613 18. Breiman L. Random Forests. *Machine Learning*. 2001;45(1):5-32.
- 614 19. Denil M, Matheson D, Freitas ND. Narrowing the Gap: Random Forests In Theory and In
615 Practice. In: Eric PX, Tony J, editors. *Proceedings of the 31st International Conference on Machine*
616 *Learning; Proceedings of Machine Learning Research: PMLR;* 2014. p. 665--73.
- 617 20. Segal M. *Machine Learning Benchmarks and Random Forest Regression*. Technical Report,
618 Center for Bioinformatics & Molecular Biostatistics, University of California, San Francisco. 2003.
- 619 21. Cawley GC, Talbot NLC. On Over-fitting in Model Selection and Subsequent Selection Bias in
620 Performance Evaluation. *J Mach Learn Res*. 2010;11:2079–107.

- 621 22. Díaz-Uriarte R, Alvarez de Andrés S. Gene selection and classification of microarray data
622 using random forest. *BMC Bioinformatics*. 2006;7(1):3.
- 623 23. Jiang H, Deng Y, Chen H-S, Tao L, Sha Q, Chen J, et al. Joint analysis of two microarray gene-
624 expression data sets to select lung adenocarcinoma marker genes. *BMC Bioinformatics*. 2004;5(1):81.
- 625 24. Saeys Y, Inza I, Larrañaga P. A review of feature selection techniques in bioinformatics.
626 *Bioinformatics*. 2007;23(19):2507-17.
- 627 25. Guyon I, Elisseeff A. An introduction to variable and feature selection. *Journal of machine*
628 *learning research*. 2003;3(Mar):1157-82.
- 629 26. Bellman RE. *Adaptive control processes: a guided tour*: Princeton university press; 2015.
- 630 27. Huang da W, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene lists
631 using DAVID bioinformatics resources. *Nature protocols*. 2009;4(1):44-57.
- 632 28. de la Fraga LG, Coello Coello CA. A Review of Applications of Evolutionary Algorithms in
633 Pattern Recognition. In: Wang PSP, editor. *Pattern Recognition, Machine Intelligence and Biometrics*.
634 Berlin, Heidelberg: Springer Berlin Heidelberg; 2011. p. 3-28.
- 635 29. Uddin S, Khan A, Hossain ME, Moni MA. Comparing different supervised machine learning
636 algorithms for disease prediction. *BMC Medical Informatics and Decision Making*. 2019;19(1):281.
- 637 30. Newman M. *Networks: An Introduction*. Oxford University Press. 2010.
- 638 31. Rui X, Wunsch D. Survey of clustering algorithms. *IEEE Transactions on Neural Networks*.
639 2005;16(3):645-78.
- 640 32. Su G, Kuchinsky A, Morris JH, States DJ, Meng F. GLayer: community structure analysis of
641 biological networks. *Bioinformatics*. 2010;26(24):3135-7.
- 642 33. Maier M, Luxburg Uv, Hein M. Influence of graph construction on graph-based clustering
643 measures. *Proceedings of the 21st International Conference on Neural Information Processing*
644 *Systems*; Vancouver, British Columbia, Canada: Curran Associates Inc.; 2008. p. 1025–32.
- 645 34. Wang X, Lin Y, Song C, Sibille E, Tseng GC. Detecting disease-associated genes with
646 confounding variable adjustment and the impact on genomic meta-analysis: With application to
647 major depressive disorder. *BMC Bioinformatics*. 2012;13(1):52.
- 648 35. Pearson K. LIII. On lines and planes of closest fit to systems of points in space. *The London,*
649 *Edinburgh, and Dublin Philosophical Magazine and Journal of Science*. 1901;2(11):559-72.
- 650 36. Wright MN, Ziegler A. ranger: A Fast Implementation of Random Forests for High
651 Dimensional Data in C++ and R. 2017. 2017;77(1):17.
- 652 37. Hastie T, Tibshirani R, Friedman J. *The elements of statistical learning: data mining, inference,*
653 *and prediction*: Springer Science & Business Media; 2009.
- 654 38. Diaz-Uriarte R. GeneSrF and varSelRF: a web-based tool and R package for gene selection
655 and classification using random forest. *BMC Bioinformatics*. 2007;8(1):328.
- 656 39. Margolin AA, Nemenman I, Basso K, Wiggins C, Stolovitzky G, Dalla Favera R, et al. ARACNE:
657 an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context.
658 *BMC bioinformatics*. 2006;7 Suppl 1(Suppl 1):S7-S.
- 659 40. Boucher B, Jenna S. Genetic interaction networks: better understand to better predict.
660 *Frontiers in genetics*. 2013;4:290.
- 661 41. Mani R, St.Onge RP, Hartman JL, Giaever G, Roth FP. Defining genetic interaction.
662 *Proceedings of the National Academy of Sciences*. 2008;105(9):3461-6.
- 663 42. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: a software
664 environment for integrated models of biomolecular interaction networks. *Genome research*.
665 2003;13(11):2498-504.
- 666 43. Newman ME, Girvan M. Finding and evaluating community structure in networks. *Physical*
667 *review E*. 2004;69(2):026113.
- 668 44. Huang DW, Sherman BT, Tan Q, Collins JR, Alvord WG, Roayaei J, et al. The DAVID Gene
669 Functional Classification Tool: a novel biological module-centric algorithm to functionally analyze
670 large gene lists. *Genome Biol*. 2007;8(9):R183-R.

- 671 45. Ronnblom L, Eloranta ML. The interferon signature in autoimmune diseases. *Curr Opin*
672 *Rheumatol.* 2013;25(2):248-53.
- 673 46. Schneider WM, Chevillotte MD, Rice CM. Interferon-stimulated genes: a complex web of
674 host defenses. *Annu Rev Immunol.* 2014;32:513-45.
- 675 47. McNab F, Mayer-Barber K, Sher A, Wack A, O'Garra A. Type I interferons in infectious disease.
676 *Nat Rev Immunol.* 2015;15(2):87-103.
- 677 48. Kyogoku C, Smiljanovic B, Grun JR, Biesen R, Schulte-Wrede U, Haupl T, et al. Cell-specific
678 type I IFN signatures in autoimmunity and viral infection: what makes the difference? *PLoS One.*
679 2013;8(12):e83776.
- 680 49. Singh H, Samani D, Nambiar N, Ghate MV, Gangakhedkar RR. Prevalence of MMP-8 gene
681 polymorphisms in HIV-infected individuals and its association with HIV-associated neurocognitive
682 disorder. *Gene.* 2018;646:83-90.
- 683 50. Jans J, Unger WWJ, Vissers M, Ahout IML, Schreurs I, Wickenhagen A, et al. Siglec-1 inhibits
684 RSV-induced interferon gamma production by adult T cells in contrast to newborn T cells. *Eur J*
685 *Immunol.* 2018;48(4):621-31.
- 686 51. Hammonds JE, Beeman N, Ding L, Takushi S, Francis AC, Wang JJ, et al. Siglec-1 initiates
687 formation of the virus-containing compartment and enhances macrophage-to-T cell transmission of
688 HIV-1. *PLoS Pathog.* 2017;13(1):e1006181.
- 689 52. Demaret J, Venet F, Plassais J, Cazalis M-A, Vallin H, Friggeri A, et al. Identification of CD177
690 as the most dysregulated parameter in a microarray study of purified neutrophils from septic shock
691 patients. *Immunology Letters.* 2016;178:122-30.
- 692 53. Stroncek DF. Neutrophil-specific antigen HNA-2a, NB1 glycoprotein, and CD177. *Curr Opin*
693 *Hematol.* 2007;14(6):688-93.
- 694 54. Dietterich TG. Approximate Statistical Tests for Comparing Supervised Classification Learning
695 Algorithms. *Neural Computation.* 1998;10(7):1895-923.
- 696 55. Jin T, Xu X, Hereld D. Chemotaxis, chemokine receptors and human disease. *Cytokine.*
697 2008;44(1):1-8.
- 698 56. Das A, Guha P, Sen D, Chaudhuri TK. Role of toll like receptors in bacterial and viral diseases
699 – A systemic approach. *Egyptian Journal of Medical Human Genetics.* 2017;18(4):373-9.
- 700 57. Dinarello CA. Interleukin-1 in the pathogenesis and treatment of inflammatory diseases.
701 *Blood.* 2011;117(14):3720-32.
- 702 58. Barnes M, Freudenberg J, Thompson S, Aronow B, Pavlidis P. Experimental comparison and
703 cross-validation of the Affymetrix and Illumina gene expression analysis platforms. *Nucleic acids*
704 *research.* 2005;33(18):5914-23.
- 705 59. Rosebeck S, Leaman DW. Mitochondrial localization and pro-apoptotic effects of the
706 interferon-inducible protein ISG12a. *Apoptosis : an international journal on programmed cell death.*
707 2008;13(4):562-72.
- 708 60. Liu N, Zuo C, Wang X, Chen T, Yang D, Wang J, et al. miR-942 decreases TRAIL-induced
709 apoptosis through ISG12a downregulation and is regulated by AKT. *Oncotarget.* 2014;5(13):4959-71.
- 710 61. Gytz H, Hansen MF, Skovbjerg S, Kristensen AC, Horlyck S, Jensen MB, et al. Apoptotic
711 properties of the type 1 interferon induced family of human mitochondrial membrane ISG12
712 proteins. *Biology of the cell.* 2017;109(2):94-112.
- 713 62. Mar KB, Rinkenberger NR, Boys IN, Eitson JL, McDougal MB, Richardson RB, et al. LY6E
714 mediates an evolutionarily conserved enhancement of virus infection by targeting a late entry step.
715 *Nat Commun.* 2018;9(1):3603.
- 716 63. Saitoh S, Kosugi A, Noda S, Yamamoto N, Ogata M, Minami Y, et al. Modulation of TCR-
717 mediated signaling pathway by thymic shared antigen-1 (TSA-1)/stem cell antigen-2 (Sca-2). *Journal*
718 *of immunology (Baltimore, Md : 1950).* 1995;155(12):5574-81.
- 719 64. Meng F, Lowell CA. Lipopolysaccharide (LPS)-induced macrophage activation and signal
720 transduction in the absence of Src-family kinases Hck, Fgr, and Lyn. *The Journal of experimental*
721 *medicine.* 1997;185(9):1661-70.

- 722 65. Tang BM, Shojaei M, Parnell GP, Huang S, Nalos M, Teoh S, et al. A novel immune biomarker
723 IFI27 discriminates between influenza and bacteria in patients with suspected respiratory infection.
724 The European respiratory journal. 2017;49(6).
- 725 66. Andres-Terre M, McGuire HM, Pouliot Y, Bongen E, Sweeney TE, Tato CM, et al. Integrated,
726 Multi-cohort Analysis Identifies Conserved Transcriptional Signatures across Multiple Respiratory
727 Viruses. Immunity. 2015;43(6):1199-211.
- 728 67. Tough DF. Type I interferon as a link between innate and adaptive immunity through
729 dendritic cell stimulation. Leuk Lymphoma. 2004;45(2):257-64.
- 730 68. Hertzog PJ, O'Neill LA, Hamilton JA. The interferon in TLR signaling: more than just antiviral.
731 Trends in immunology. 2003;24(10):534-9.
- 732 69. Kovarik P, Castiglia V, Ivin M, Ebner F. Type I Interferons in Bacterial Infections: A Balancing
733 Act. Frontiers in immunology. 2016;7(652).

734

735 **Supporting information**

736 **S1 Appendix. Pre-processing.**

737 **S1 Fig. Affymetrix Interaction Network.** Affymetrix recovered interaction network at first
738 level of clustering. Selected model genes are highlighted.

739 **S1 Table. Model Gene Selection Frequency.** Affymetrix and Illumina model selected genes
740 with relative frequency of selection (genes with greater than 5% aggregated inclusion across
741 all search procedures).

742 **S2 Appendix. Biomarker Search.**

743 **S2 Fig. Illumina Interaction Network.** Illumina recovered interaction network at first level
744 of clustering. Selected model genes are highlighted.

745 **S2 Table. Highly selected gene clusters from Affymetrix and Illumina interaction**
746 **network.** Table containing the genes from the 4 highly model selected Illumina clusters, and
747 5 highly model selected gens from the Affymetrix clusters.

748 **S3 Appendix. Inferred Interaction Networks.**

749 **S3 Table. Out-sample results of gene lists.** The out-sample results from running Affymetrix
750 derived gene lists on the Illumina data, and the Illumina derived gene lists on the Affymetrix
751 data

752