1  *ImSig:* **A resource for the identification and quantification of immune signatures in**

2  **blood and tissue transcriptomics data**

3  Ajit Johnson Nirmal[†], Tim Regan[†], Barbara Bo-Ju Shih[†], David Arthur Hume[†], Andrew

4  Harvey Sims[‡], Tom Charles Freeman[†]

5  [†]The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of

6  Edinburgh, Easter Bush, Edinburgh, EH5 9RG, UK.

7  [‡]Applied Bioinformatics of Cancer, Edinburgh Cancer Research Centre, Institute of Genetics

8  and Molecular Medicine, University of Edinburgh, Crewe Road South, Edinburgh, EH4

9  2XU, UK.

10  **Corresponding Author:**

11  Tom C Freeman
12  Systems Immunology Group
13  The Roslin Institute and Royal (Dick) School of Veterinary Studies
14  University of Edinburgh
15  Easter Bush
16  EH25 9RG
17  T: +44 (0)131 651 9203
18  F: +44 (0)131 651 9105
19  tom.freeman@roslin.ed.ac.uk
20
21
22
23
24
25
26
27
28
29
30

**Abstract**

The outcome of many diseases is commonly correlated with the immune response at the site of pathology. The ability to monitor the status of the immune system in situ provides a mechanistic understanding of disease progression, a prognostic assessment and a guide for therapeutic intervention. Global transcriptomic data can be deconvoluted to provide an indication of the cell types present and their activation state, but the gene signatures proposed to date are either disease-specific or have been derived from data generated from isolated cell populations. Here we describe an improved set of immune gene signatures, *ImSig*, derived based on their co-expression in blood and tissue datasets. *ImSig* includes validated lists of marker genes for the main immune cell types and a number of core pathways. When used in combination with network analysis, *ImSig* is an accurate and easy to use approach for monitoring immune phenotypes in transcriptomic data derived from clinical samples.

**Introduction**

The differentiation and activation of immune cells is associated with changes in the expression of hundreds to thousands of genes (1, 2). Genes specifically expressed by a cell type or cells in a particular state of activation can be used as markers (3) to monitor immune cells in a disease environment, facilitate tailored therapies (4) and clinical stratification of diseases (4, 5). Several studies have identified immune cell markers to 'deconvolute' gene expression data. Some of the widely used methods (Table 1) include LLSR (6, 7), qprog (8), DSA (9), PERT (10), MMAD (11) and CIBERSORT (12). For a more detailed review of deconvolution methods read (13). Whilst the derivations differ, all these methods define their marker gene lists based on the comparison of gene expression data from isolated immune cells. Here we employ the principle of co-expression as the basis to derive cell-type specific immune signatures directly from large clinical transcriptomic datasets. This method exploits the fact that the mRNA abundance of genes expressed by a specific cell type will correlate with the number of those cells in a given sample. Between similar samples in a given dataset there are always subtle differences in their cellular composition due to innate variation (e.g. normal variation between individuals, disease severity or subtype etc.), as well as inconsistencies in sampling. Genes expressed by a particular cell type may therefore be identified based on their distinct co-expression profile without the need to physically isolate the cells. We have used this approach to identify robust immune cell type-specific gene expression signatures, known collectively as *ImSig*, derived from and validated on multiple independent datasets to ensure their wide applicability. We have benchmarked *ImSig* against other methods and shown it to out-perform them. We also provide an easy to use algorithm to identify the presence of different immune cell populations in any transcriptomics dataset.

76    **Materials and Methods**

77    *Selection of datasets*

78    The primary datasets used for deriving *ImSig* were identified from the Gene Expression

79    Omnibus (National Centre for Biotechnology Information) or ArrayExpress (European

80    Bioinformatics Institute) databases. Datasets from isolated immune cells were identified, and

81    restricted to only those based on the Affymetrix Human Genome U133 Plus 2.0 Array with

82    availability of raw data (.CEL) files. These included: B cells (germinal centre B cells, naïve B

83    cells, memory B cells, IgM+IgD+CD27+ B cells, class switched B cells, IgM+IgD-CD27+ B

84    cells); plasma B cells; monocytes; T cells (central memory T cells, effector memory T cells,

85    naïve T cells, gamma-delta T cells, CD4- T cells, CD4+ T cells, CD8- T cells, CD8+ T cells);

86    macrophages (resting and activated), neutrophils, NK cells and platelets (see Table S5 for

87    details).

88    A second group of datasets were identified for the purpose of refining and validating the final

89    *ImSig* gene lists. They consisted of blood and tissue datasets, derived from a broad spectrum

90    of diseases and were restricted to data generated on the Affymetrix U133 Plus 2.0 Array with

91    available raw data (see Table S6 for a list of these data).

92    *Processing of microarray datasets*

93    Quality control (QC) of data from each dataset was performed using the ArrayQualityMetrics

94    package in Bioconductor and scored on the basis of six quality metrics (31). Any array failing

95    more than one metric was removed. Following QC, signal intensity were summarised and

96    normalised using robust multi-array average (RMA) in R using the 'oligo package' (32).

97    Data from isolated immune cell populations were merged and normalised as described above.

98    In order to check that samples clustered according to cell type specific rather than study or

4

99   any other factor, the RMA normalised data was loaded into the network analysis tool Miru

100  (Kajeka Ltd., Edinburgh, UK). Miru calculates a matrix of pairwise Pearson correlation

101  coefficients (*r*) expression values between every pair of genes/samples in a dataset. Graph

102  layout in-tool is performed using a modified Fast Multipole Multilevel Method (FMMM)

103  (33) and the resulting network is rendered in a 3-D environment. Networks are composed of

104  nodes (representing transcripts/samples) connected by weighted edges (representing

105  correlation values). After loading the immune cell data a sample similarity network was then

106  plotted at a correlation threshold of *r* > 0.83. All sample outliers i.e. samples that did not

107  group with other samples of the sample type, were removed. The remaining 329 samples

108  clustered based on cell type rather than study (Figure S4). For blood and tissue datasets, the

109  data was collapsed to one probe-set per gene by choosing the probe-set with the highest

110  variance across samples.

111  *Refinement of signatures (Cluster model algorithm)*

112  The quality-controlled datasets were loaded into the network analysis tool Miru. Within the

113  tool, a correlation network was generated and clustered using the MCL algorithm (inflation

114  value: 2.2). A proportion of the genes in each MCL cluster were replaced with random genes

115  in increments of 2% from 0-100% of total genes. The percentage of genes from the original

116  MCL cluster in this modified cluster was defined as $Percent_{similar}$ (percentage of genes with

117  high similarity). The similarity of each gene to other members of the cluster or annotation is

118  defined by the median value of its Pearson correlation coefficients to every other member

119  within a cluster or annotation, and the median of this value from all genes within a cluster or

120  annotation is referred to as $Pearson_{group}$. The decrease in $Pearson_{group}$ with increasing

121  replacements of MCL cluster by random genes was modelled as a sigmoid function of

122  $Percent_{similar}$ using nonlinear least squares in R.

123     In the situation where the groups of genes were of the same signature (different cell type

124     signatures derived from both blood and tissue) instead of modified MCL clusters, the

125     $Percent_{similar}$ is unknown whilst $Pearson_{group}$ can be calculated. Therefore, inverse estimates of

126     $Percent_{similar}$ using $Pearson_{group}$ were made using the R package "investr". The upper and

127     lower threshold of $Pearson_{group}$, beyond which investr function cannot estimate the

128     $Percent_{similar}$, were noted and used as cut off for determining if genes will be discarded from

129     the refined signature.

130     In the second stage of the filtering process, signatures with a $Pearson_{group}$ 1) higher than upper

131     threshold were left unchanged; 2) between upper and lower thresholds were reduced in size,

132     using the model above to determine the number of genes to discard; 3) less than the lower

133     threshold were considered to be absent from the dataset. This method of filtering would allow

134     greater stability cross datasets, whilst retaining more flexibility with a more comprehensive

135     list of genes with informative signature. We used the cluster model algorithm on eight blood

136     and eight tissue datasets to refine the *ImSig* signature lists (Figure 2A)

137     *Derivation of ImSig$_{blood}$*

138     The most differentially expressed genes (DEGs) for each isolated immune cell type was

139     determined by calculating the average fold change for a particular cell type relative to the

140     rest. The top 100 DEGs for each cell type were refined across eight blood datasets (Table S6)

141     using the cluster model algorithm. The resultant sets of genes derived from each dataset were

142     then compared, and the most overlapping set of genes were defined as the blood signature set,

143     *ImSig$_{blood}$*.

144     *Derivation of ImSig$_{tissue}$*

6

145     The same approach as of $ImSig_{blood}$ was not successful in defining a tissue specific $ImSig$,

146     since the top 100 DEGs for each cell type were poorly co-expressed in complex tissue

147     datasets. This is likely due to the fact that these 100 DEGs were derived from isolated cells,

148     some of which were cultured *in vitro*, where their phenotype more closely resembled that

149     their counterparts in blood. We therefore decided that the best approach was to use a

150     correlation based approach. The expression data of isolated immune cells was loaded into the

151     network analysis tool Miru. A large and highly structured network graph was constructed

152     using a correlation threshold of $r > 0.8$. The network was then clustered into groups of genes

153     sharing similar profiles using the Markov Clustering (MCL) algorithm with an MCL inflation

154     value set to 2.2 (34). These clusters were then extensively explored to find genes that were

155     distinctively expressed in only one cell type in contrast to the rest. These genes were then

156     explored in the context of four tissue datasets as a class set and network graphs constructed

157     and clustered as described earlier. For each dataset, clusters identified as being specific

158     (based on the added class set) to a particular cell type were isolated. The resultant set of genes

159     were compared to each other and the most common set of genes were refined in another 8

160     tissue datasets (Table S6) using the cluster model algorithm to define the $ImSig_{tissue}$.

161     *Derivation of pathway signatures*

162     Whilst analysing the clusters and refining them to be cell type-specific, we also identified a

163     number of other clusters that were consistently co-expressed across different datasets. With

164     the help of GO Annotation and known marker genes, we were also able to define these

165     clusters as cell cycle-associated, interferon stimulated and protein translational activity. These

166     clusters were further refined in blood and tissue datasets as describe above using the cluster

167     model algorithm.

168     *Validation of ImSig in mixed cell population datasets*

169    *ImSig* was validated using additional independent datasets, including two blood (heart attack

170    blood samples: GSE48060 and type I diabetes mellitus blood samples: GSE55098), two

171    tissue (breast tumour tissue samples: GSE58812 and primary CNS tumour tissue samples)

172    and an infection dataset (*Chlamydia trachomatis* infection tissue sample: GSE20436). All

173    datasets were pre-processed as described above. A number of transcriptomic profiles derived

174    from RNA-seq technology were also analysed by *ImSig* to ensure its wide applicability and

175    lack of platform dependency, in particular RNA-seq data were downloaded from TCGA

176    database.

177    *ImSig cluster scoring algorithm*

178    In order to facilitate the use of *ImSig* a scoring system was devised that supports the

179    identification of any given signature without the need to perform network analysis. For any

180    given transcriptomic dataset, the calculation of the *ImSig* scores is a two-step process where

181    an initial score is first computed based on the following formula:

$$
\text{Intial score}^{(r,i)} = \frac{\text{Median correlation}}{\text{Standard deviation}} * \frac{\text{Observed number of nodes (Genes)}^{(r,i)}}{\text{Maximum possible nodes (Genes)}^{(r,i)}} * \frac{\text{Observed number of edges}^{(r,i)}}{\text{Maximum possible edges}^{(r,i)}}
$$

182

183    Where *r* is the correlation cut-off and *i* is the cell type/pathway signature.

184    Median correlation is calculated by computing the correlation values across samples for all

185    possible pairs of genes within any given signature and then taking the median value. The

186    standard deviation is calculated by computing the mean expression value of all genes within a

187    signature and then calculating its standard deviation across samples. The maximum possible

188    edges is calculated with [n*(n-1)]/2, where *n* is the number of genes in any given signature.

189    The maximum possible nodes is the number of genes defining a particular *ImSig* signature.

8

190   The initial score is computed for all eight cell type clusters (B cells, T cells, monocytes,

191   macrophages, NK cells, neutrophils, plasma cells, platelets) and three pathway clusters (cell

192   division, protein translational and interferon response) using a range of Pearson correlation

193   coefficient thresholds, from 0.50 to 0.99 at 0.01 intervals. The resulting matrix contains 50

194   scores for each of the signature. At this point we set an 'initial score threshold' of 20 and 10

195   for microarray and RNA-seq datasets, respectively (these were determined empirically,

196   Figure 2B&C). Any value below this threshold is not regarded to be a genuine cluster due to

197   a poor correlation between genes within the signature at the set *r*-value. We recommend these

198   thresholds as they are based on observations from numerous datasets. Following this the final

199   *ImSig* score is calculated for each cluster using the following formula.

$$ImSig\ score^i = 1 - \left[ \frac{number\ of\ times\ (intial\ score^i) < 20}{50} \right]$$

200

201   All data should be in log scale for calculating *ImSig* score. An R script is available for

202   running *ImSig* scoring algorithm. The script can be downloaded here:

203   www.github.com/systems-immunology-roslin-institute/ImSig. The final score (*ImSig* score)

204   is a value between 0 and 1. After extensive evaluation, any value above 0.3 is regarded as

205   evidence that the cell type/pathway signature is present in the dataset.

206   *Comparison with CIBERSORT:*

207   A blood and a tissue dataset were used for this purpose where there was some prior

208   knowledge about the cell populations present and their relative abundance in sample sub-

209   groups.

210   A blood dataset (SLE patients: GSE49454) was downloaded from GEO. The authors of this

211   study had provided the cell counts along with the transcriptomics data in this file. Initially,

212   the patients were ordered based on the cell count for each of the different cell types

213    independently (B cells, T cells, NK cells and neutrophils). They were then equally divided

214    into three groups and the top and bottom groups were used for analysis. Two-tailed unequal

215    variance T-test showed a significant alteration in cell counts between these two group of

216    patients for all four cell types ($p<0.05$). Using CIBERSORT and *ImSig* the relative proportion

217    of immune cells were then computed. For CIBERSORT the data was loaded into

218    (https://cibersort.stanford.edu/) as per the authors instructions and the computed relative

219    proportions were downloaded. The relative proportions of immune subtypes were all summed

220    to make up the parent cell type (T cells, B cells, neutrophils, NK cells). Then, each cell type

221    was normalised independently to be represented as a fraction of 1 across samples (i.e., the

222    sum of normalised cell proportion for any cell type is equal to 1). Similarly, for *ImSig* the

223    relative abundance of immune cells were calculated by averaging the expression of signature

224    genes for each sample and then normalised to represent them as a fraction of 1. Two-tailed

225    unequal variance T-test was then used to test for significant change in cell proportions

226    between the two groups of patients in all four cell types.

227    Similarly, a tissue dataset (trachoma: GSE20436) was downloaded. The patients were divided

228    into three groups as per the level of infectivity according to its authors (controls, symptom

229    +ve/*C. trachomatis* -ve patients, and symptom +ve/*C. trachomatis* +ve patients). As

230    described earlier, the relative proportion of immune cells (T cells, B cells, neutrophils,

231    monocytes, macrophages, NK cells and plasma cells) were computed and normalised using

232    CIBERSORT and *ImSig*. This was followed by a one-way analysis of variance (ANOVA) to

233    test for significant changes in cell numbers between the three groups of patients.

234

235

236 **Results**

237 *Blood and tissue immune signatures (ImSig$_{blood/tissue}$)*

238 *ImSig* was derived as described in the experimental procedures, and as shown in (Figure 1).

239 Briefly, an initial meta-analysis was carried out on 330 samples of isolated human immune

240 cell populations and the top 100 differentially expressed genes were determined for each

241 immune cell type. Using a network-based approach to identify sets of robustly co-expressed

242 (correlated) genes in a variety of blood datasets, the lists were further refined (Figure 2A).

243 The resulting cell-specific marker gene lists were collectively named *ImSig*$_{blood}$. However, the

244 limitations of this approach become evident from network analysis of clinical tissue

245 transcriptomic datasets, where the cell-based marker genes showed independent expression.

246 To overcome this issue we identified the most conserved cell type-specific groups of genes

247 based on their co-expression across four tissue datasets, and further refined them by

248 examining a eight other tissue datasets (Figure 2A). This resulted in our *ImSig*$_{tissue}$ gene

249 signatures. *ImSig*$_{blood}$ contains 491 marker genes and *ImSig*$_{tissue}$ contains 569 marker genes for

250 B cells, monocytes, macrophages (tissue only), neutrophils, NK cells, T cells, plasma cells,

251 platelets (blood only), cell cycle, protein translation and interferon signalling. For a full list of

252 the genes comprising these signatures and numbers for each cell type or pathway see Table

253 S1. GO term analysis confirmed that the cell marker lists for both *ImSig* signatures were

254 highly enriched in genes related to immune function (Table S2, S3). The overlap between

255 blood and tissue signature varied depending on cell type/pathway (Figure S3).

256 *Genes that make up the signatures*

257 Table S1 highlights the sets of genes that distinguish *ImSig*$_{blood}$ and *ImSig*$_{tissue}$. The process

258 signatures; cell cycle, interferon response and protein synthesis (translational activity) are

259 relatively robust in both blood and tissue. The T cell clusters in both cases are anchored and

11

260    validated by the subunits of CD3, but otherwise, there is very little overlap. The implication

261    is that the T cells that enter tissues in a pathological situation are radically different in their

262    gene expression profiles from the bulk of naïve T cells in peripheral blood. Note that there is

263    no evidence of a cluster of genes associated with specific T cell polarisation states. The key

264    transcription factors *FOXP3* (Treg), *RORC* (Th17) and *GATA3* (Th2) do not form part of

265    clusters, since they are expressed by other cell types. However T-BET (*TBX21*), considered

266    to be a Th1 specific transcription factor is in the NK cell cluster a cell type in which it is also

267    strongly expressed. The NK cell cluster also shows considerable divergence between blood

268    and tissue, in particular the NK cell receptor family being much more robustly co-expressed

269    in tissue RNA. In blood, many of these receptors are also detectable in gamma-delta T cells

270    (14). The various myeloid clusters are rather more difficult to be associated with specific cell

271    types. The macrophage cluster, specific to the tissue data set, contains the CSF1R, which is

272    known to be macrophage-specific and essential for differentiation and survival (15), and also

273    contains many of the genes that are up-regulated in monocyte-derived macrophages derived

274    by cultivation in CSF1 (16). An unexpected member of this cluster is CD4. In blood, CD4 is

275    expressed at similar levels in CD4+ T cells and monocytes, and so does not form part of a T

276    cell cluster. In tissue, CD4 is highly-expressed by macrophages, and correlates more highly

277    with their presence than with the presence of T cells. The clusters annotated provisionally as

278    monocyte and neutrophil have very little overlap between the blood and tissue profiles.

279    Archetypal markers, CD14 for monocytes and the G-CSF receptor and chemokine receptor

280    CXCR2 (the receptor for IL8) are co-expressed with very different gene sets in blood and

281    tissues. Hence, it may be more appropriate to consider distinct separate myelomonocytic

282    regulons, reflecting the rapid differentiation of these cells following extravasation. For

283    example, S100A8/A9, which encode the most abundant neutrophil proteins (17), are also

284    expressed by monocytes, but rapidly down-regulated as they differentiate to macrophages.

285    The mRNAs encoding many neutrophil-specific granule proteins (MPO, lactoferrin etc) are

286    expressed most highly in progenitor cells (16), and do not contribute to a signature in either

287    blood or tissue.

288    *ImSig scoring algorithm*

289    The *ImSig* scoring algorithm was developed to reflect the correlation and expression level of

290    the marker genes in any given dataset. The algorithm generates a numerical likelihood score

291    that a given signature is present in a dataset. Based upon empirical evaluation of a wide range

292    of data, an *ImSig* score >0.3 indicates positive identification of the signature in a given

293    dataset (Figure 2B&C). *ImSig* scores for all the validation datasets along with six other RNA-

294    seq datasets are provided in Table S4. Consistent with its derivation, the *ImSig* macrophage

295    signature is absent from any blood datasets, irrespective of platform. Conversely, the platelet

296    signature was not scored positive in any tissue dataset examined. As with other deconvolution

297    methods, *ImSig* works best when the majority of signature genes are present in the dataset to

298    be analysed. Based upon a permutation analysis of the effect of random removal of genes on

299    the *ImSig* score (Figure 2D) a minimum of 75% of the genes from each individual signatures

300    is required for an accurate representation analysis. Being correlation-based, a dataset

301    generally needs to comprise of at least 20 distinct samples is needed to provide sufficient

302    diversity before the *ImSig* algorithm can be applied.

303    *Validation of blood and tissue marker genes*

304    To test its universality, we applied $ImSig_{blood}$ to deconvolution of a range transcriptomics data

305    derived from whole blood or peripheral blood mononuclear cells (PBMC). Examples of these

306    analyses are given here. Data from the blood of 21 control and 31 heart attack patients

307    (GSE48060) identified the presence of B cells, T cells, NK cells, plasma cells, platelets,

308    monocytes and neutrophils (Figure S1A). In terms of the average expression of marker genes,

13

309    no consistent difference was observed between the control and heart attack samples

310    suggesting that relative blood cell numbers were not altered. The macrophage and cell cycle

311    signatures were not detected. In contrast, *ImSig* analysis of PBMC's from control and patients

312    with type 1 diabetes mellitus (GSE55098) identified increased proliferation in a number of

313    samples (Figure 3A) and the analysis also clearly identified the presence of T cells, B cells,

314    along with plasma cells, monocytes, neutrophils, NK cells and platelets (Table S4). Notably

315    there was also significantly lower expression (p=1E-10) of the NK cells markers genes in

316    samples derived type 1 diabetes (Figure 3A) where these cells are known to be dysregulated

317    (18, 19).

318    To validate *ImSig*tissue, we first examined a dataset of triple-negative breast cancers derived

319    from 107 patients (GSE58812). As expected, and in keeping with our previous network

320    analysis of multiple tumour datasets (20), the cell cycle cluster was readily detected,

321    reflecting the heterogeneity in proliferative index between tumours. The analysis revealed

322    macrophages, T cells, B cells, plasma cells, interferon but there was no evidence of platelets,

323    neutrophils and NK cells present in these samples (Figure 3B). The levels of all immune cells

324    (as judged by the average expression of the marker genes) varied greatly between samples.

325    By contrast, a relatively small brain tumour dataset comprising 23 samples of primitive

326    neuroectodermal tumors and medulloblastomas lacked evidence of immune cell infiltration,

327    other than an NK signature (Figure S1B). Being behind the blood-brain barrier, lymphocyte

328    populations in these tumours are likely absent or at very low levels (21) but infiltration of T

329    cells was evident in other brain tumour datasets that we have analysed (Table S4). Neutrophil

330    signatures were absent from tumour datasets. However, as expected, a dataset of eye swabs

331    taken from eyes of controls or children with the symptoms of trachoma (GSE20436) (22) was

332    positive for all signatures of immune cells (Figure S2). Previous studies have shown that in

333    certain chlamydial infections, neutrophils recruit T cells to the site of infection (23), other

14

334    studies report the involvement of NK cells, monocytes and macrophages (24-26). Finally, we

335    demonstrate the explorative power of *ImSig* when coupled with network analysis. The genes

336    comprising the signatures were selected as being core 'invariant' markers of a particular cell

337    type. When used in the context of a correlation analysis of a complete dataset, if the relevant

338    cells are present within the samples, surrounding the signature genes will be other genes

339    expressed in these populations. In this manner one can better evaluate the activation state of

340    immune cells *in situ*. Using the trachoma dataset as an example we highlight known immune

341    related genes that were co-expressed with *ImSig* core signature genes (Figure 4). The

342    associated *ImSig* scores for all the validation datasets can be found in Table S4.

343    *Comparison with CIBERSORT*

344    The ability of *ImSig* and CIBERSORT to identify changes in relative proportions of cells

345    between sample groups was compared using a blood (GSE49454: Systemic lupus

346    erythematosus patients) and a tissue dataset (GSE20436: trachoma). For the blood dataset,

347    cell counts were available for B cells, neutrophils, T cells and NK cells. Both methods

348    generally performed well, $ImSig_{blood}$ demonstrated a significant difference (p<0.05) in all four

349    cell types, although CIBERSORT failed to show a significant difference in B cells (p=0.389)

350    (Figure 5A, Table S7). Samples from the trachoma dataset were divided into three groups of

351    20, based on the level of infection as originally described (for more detail see Methods).

352    Although actual cell counts are not available for these data, it is known that the immune

353    infiltrate increases with the level of infection (27). $ImSig_{tissue}$ showed there to be a significant

354    increase (p<0.05) in all seven cell immune types (B cells, neutrophils, T cells, NK cells,

355    plasma cells, monocytes and macrophages) during an active infection, while significant

356    differences were only reported for T cells and macrophages using CIBERSORT (Figure 5B,

357    Table S7). Moreover, the pattern observed using CIBERSORT did not seem to correlate with

15

358    the infection status of C. *trachomatis* (Figure 5B). CIBERSORT was also used in its native

359    form, i.e. the subtypes were not summed to represent the parent population. A significant

360    change in cell number was observed only for M2 macrophages (p=0.001), activated mast

361    cells (p=0.022) and resting dendritic cells (p=0.0007). The 19 other immune cell groups

362    defined by CIBERSORT showed no significant difference in cell proportion across patient

363    groups (Table S8).

364

365

366

367

368

369

370

371

372

373

374

375

376 **Discussion**

377 In the last few years a number of immune marker gene signatures have been proposed (6-12).

378 The current work is based on the observation that when correlation (co-expression) network

379 analysis is employed to explore large transcriptomics datasets derived from normal or

380 diseased tissues, clusters of genes associated with specific immune cell populations, or

381 specific transcriptional regulons such as protein synthesis, interferon response or cell cycle,

382 are frequently observed clustered together (20, 22, 28, 29). This is because the abundance of

383 mRNAs derived from cell-specific, or process-specific genes is correlated with relative

384 number of those cells expressing those genes within a sample, resulting in their observed co-

385 expression across a sample set. The most important conclusion from our analysis is that

386 signatures based upon cells isolated from blood cannot be applied with any confidence to

387 tissue data.

388 The utility of the blood and tissue *ImSig* gene lists has been demonstrated through

389 applications to a number of datasets. Other approaches to deconvolution include LLSR (7),

390 qprog (8), DSA (9), PERT (10), MMAD (11) and CIBERSORT (12). Each is based on a

391 signature derived by a different data mining approach ranging from simple matrix

392 decomposition to complex iterative procedure. Of these methods CIBERSORT was shown to

393 out-perform others (12) in terms of analysis of tissue data with noise or unknown content and

394 was reported to be able to differentiate closely related cell types. CIBERSORT includes

395 profiles for 22 distinct cell types, including various states of T cell activation and macrophage

396 differentiation. The network analysis of disease datasets herein does not support robust

397 clusters that distinguish macrophage activation states, in keeping with previous analysis (20).

398 In essence, the best one can do is define three myeloid states (neutrophil, monocyte,

399 macrophage), and the inducible genes are disease/lesion specific. Expression QTL analysis of

17

400    inducible gene expression in monocytes suggests that inducible gene expression profiles may

401    also be individual-specific (30).

402    An ideal workflow for employing *ImSig* would involve running the *ImSig* algorithm to

403    identify the different immune cell populations in a dataset and then using the average

404    expression of signature genes to understand the relative proportion of cells between samples

405    and clinical subsets. This can be followed by network analysis which can be used to better

406    understand the wider context of the immune environment. Through observing the genes that

407    closely correlate with the core signature genes, one can better under the type of activation or

408    indeed the level of involvement which these cells play in a given microenvironment of a

409    disease state. As an example we have highlighted a few immune related genes that are co-

410    expressed with our core signature genes in the trachoma dataset (Figure 4). The expression

411    profiles of known immune modulatory genes such as *IFNG, LAG3, CD44, FOX03, FOXP3,*

412    *CD80, IL20, STAT4, IL17A* etc are correlated with the core macrophage and T cell signature

413    genes, suggesting that the macrophages are undergoing classical activation, and the T cells

414    include Th17, TReg and Th1 states. Thus such explorative analysis can be employed using

415    *ImSig* to understand the differentiation state of immune cells between patient groups.

416    The *ImSig* algorithm has been tested on data derived microarray and RNA-seq platforms. We

417    have also tested its applicability across a wide range of datasets derived from blood, tissue,

418    sputum and faecal samples (data not shown). As long as immune cells are present, *ImSig*

419    efficiently identifies the cell types present. We therefore anticipate that *ImSig* and the

420    methodological approaches described here will prove valuable for studying immune cell

421    variation in human transcriptomics data derived from a wide variety of conditions clinical

422    samples.

423

424  **Author contributions**

425  A.J.N performed the majority of work described here. A.J.N, T.R, B.J.S, D.A.H, A.H.S and

426  T.C.F wrote and edited the manuscript. A.H.S and T.C.F supervised the project.

427  **Acknowledgements**

433

434

435

436

437

438

439

440

441

442

443

**References**

444

445    1.    Mabbott, N. A., J. K. Baillie, H. Brown, T. C. Freeman, and D. A. Hume. 2013. An
446            expression atlas of human primary cells: inference of gene function from
447            coexpression networks. *BMC Genomics* 14: 632-632.
448    2.    Robinette, M. L., A. Fuchs, V. S. Cortez, J. S. Lee, Y. Wang, S. K. Durum, S.
449            Gilfillan, M. Colonna, and C. the Immunological Genome. 2015. Transcriptional
450            programs define molecular characteristics of innate lymphoid cell classes and subsets.
451            *Nat Immunol* 16: 306-317.
452    3.    Pui, C.-H., and W. E. Evans. 1998. Acute Lymphoblastic Leukemia. *New England*
453            *Journal of Medicine* 339: 605-615.
454    4.    van 't Veer, L. J., and R. Bernards. 2008. Enabling personalized cancer medicine
455            through analysis of gene-expression patterns. *Nature* 452: 564-570.
456    5.    Golub, T. R., D. K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. P. Mesirov, H.
457            Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, C. D. Bloomfield, and E. S.
458            Lander. 1999. Molecular Classification of Cancer: Class Discovery and Class
459            Prediction by Gene Expression Monitoring. *Science* 286: 531-537.
460    6.    Abbas, A. R., D. Baldwin, Y. Ma, W. Ouyang, A. Gurney, F. Martin, S. Fong, M. van
461            Lookeren Campagne, P. Godowski, P. M. Williams, A. C. Chan, and H. F. Clark.
462            2005. Immune response in silico (IRIS): immune-specific genes identified from a
463            compendium of microarray expression data. *Genes Immun* 6: 319-331.
464    7.    Abbas, A. R., K. Wolslegel, D. Seshasayee, Z. Modrusan, and H. F. Clark. 2009.
465            Deconvolution of Blood Microarray Data Identifies Cellular Activation Patterns in
466            Systemic Lupus Erythematosus. *PLoS ONE* 4: e6098.
467    8.    Gong, T., N. Hartmann, I. S. Kohane, V. Brinkmann, F. Staedtler, M. Letzkus, S.
468            Bongiovanni, and J. D. Szustakowski. 2011. Optimal Deconvolution of
469            Transcriptional Profiling Data Using Quadratic Programming with Application to
470            Complex Clinical Blood Samples. *PLoS ONE* 6: e27156.
471    9.    Zhong, Y., Y.-W. Wan, K. Pang, L. Chow, and Z. Liu. 2013. Digital sorting of
472            complex tissues for cell type-specific gene expression profiles. *BMC Bioinformatics*
473            14: 89.
474    10.   Qiao, W., G. Quon, E. Csaszar, M. Yu, Q. Morris, and P. W. Zandstra. 2012. PERT:
475            A Method for Expression Deconvolution of Human Blood Samples from Varied
476            Microenvironmental and Developmental Conditions. *PLoS Comput Biol* 8: e1002838.
477    11.   Liebner, D. A., K. Huang, and J. D. Parvin. 2014. MMAD: microarray
478            microdissection with analysis of differences is a computational tool for deconvoluting
479            cell type-specific contributions from tissue samples. *Bioinformatics* 30: 682-689.
480    12.   Newman, A. M., C. L. Liu, M. R. Green, A. J. Gentles, W. Feng, Y. Xu, C. D. Hoang,
481            M. Diehn, and A. A. Alizadeh. 2015. Robust enumeration of cell subsets from tissue
482            expression profiles. *Nat Meth* 12: 453-457.
483    13.   Hackl, H., P. Charoentong, F. Finotello, and Z. Trajanoski. 2016. Computational
484            genomics tools for dissecting tumour-immune cell interactions. *Nat Rev Genet* 17:
485            441-458.
486    14.   The Fantom Consortium. 2014. A promoter-level mammalian expression atlas. *Nature*
487            507: 462-470.
488    15.   Hume, D. A., and K. P. A. MacDonald. 2012. Therapeutic applications of
489            macrophage colony-stimulating factor-1 (CSF-1) and antagonists of CSF-1 receptor
490            (CSF-1R) signaling. *Blood* 119: 1810.
491    16.   Joshi, A., C. Pooley, T. C. Freeman, A. Lennartsson, M. Babina, C. Schmidl, T.
492            Geijtenbeek, F. C. the, T. Michoel, J. Severin, M. Itoh, T. Lassmann, H. Kawaji, Y.

493         Hayashizaki, P. Carninci, A. R. R. Forrest, M. Rehli, and D. A. Hume. 2015.
494         Technical Advance: Transcription factor, promoter, and enhancer utilization in human
495         myeloid cells. *Journal of Leukocyte Biology* 97: 985-995.
496   17.   Perera, C., H. P. McNeil, and C. L. Geczy. 2009. S100 Calgranulins in inflammatory
497         arthritis. *Immunol Cell Biol* 88: 41-49.
498   18.   Rodacki, M., B. Svoren, V. Butty, W. Besse, L. Laffel, C. Benoist, and D. Mathis.
499         2007. Altered Natural Killer Cells in Type 1 Diabetic Patients. *Diabetes* 56: 177-185.
500   19.   Qin, H., I.-F. Lee, C. Panagiotopoulos, X. Wang, A. D. Chu, P. J. Utz, J. J. Priatel,
501         and R. Tan. 2011. Natural Killer Cells From Children With Type 1 Diabetes Have
502         Defects in NKG2D-Dependent Function and Signaling. *Diabetes* 60: 857-866.
503   20.   Doig, T. N., D. A. Hume, T. Theocharidis, J. R. Goodlad, C. D. Gregory, and T. C.
504         Freeman. 2013. Coexpression analysis of large cancer datasets provides insight into
505         the cellular phenotypes of the tumour microenvironment. *BMC Genomics* 14: 1-16.
506   21.   Louveau, A., T. H. Harris, and J. Kipnis. 2015. Revisiting the Mechanisms of CNS
507         Immune Privilege. *Trends in Immunology* 36: 569-577.
508   22.   Natividad, A., T. C. Freeman, D. Jeffries, M. J. Burton, D. C. Mabey, R. L. Bailey,
509         and M. J. Holland. 2010. Human conjunctival transcriptome analysis reveals the
510         prominence of innate defense in Chlamydia trachomatis infection. *Infect Immun* 78.
511   23.   de Oca, R. M., A. J. Buendía, L. Del Río, J. Sánchez, J. Salinas, and J. A. Navarro.
512         2000. Polymorphonuclear Neutrophils Are Necessary for the Recruitment of CD8+ T
513         Cells in the Liver in a Pregnant Mouse Model of Chlamydophila abortus (Chlamydia
514         psittaci Serotype 1) Infection. *Infection and Immunity* 68: 1746-1751.
515   24.   Belay, T., F. O. Eko, G. A. Ananaba, S. Bowers, T. Moore, D. Lyn, and J. U.
516         Igietseme. 2002. Chemokine and Chemokine Receptor Dynamics during Genital
517         Chlamydial Infection. *Infection and Immunity* 70: 844-850.
518   25.   Liu, W., and K. A. Kelly. 2008. Prostaglandin E2 modulates dendritic cell function
519         during chlamydial genital infection. *Immunology* 123: 290-303.
520   26.   Ren, Q. U. N., S. J. Robertson, D. Howe, L. F. Barrows, and R. A. Heinzen. 2003.
521         Comparative DNA Microarray Analysis of Host Cell Transcriptional Responses to
522         Infection by Coxiella burnetii or Chlamydia trachomatis. *Annals of the New York*
523         *Academy of Sciences* 990: 701-713.
524   27.   Hu, V. H., M. J. Holland, and M. J. Burton. 2013. Trachoma: Protective and
525         Pathogenic Ocular Immune Responses to Chlamydia trachomatis. *PLoS Neglected*
526         *Tropical Diseases* 7: e2020.
527   28.   Freeman, T. C., A. Ivens, J. K. Baillie, D. Beraldi, M. W. Barnett, D. Dorward, A.
528         Downing, L. Fairbairn, R. Kapetanovic, S. Raza, A. Tomoiu, R. Alberio, C. Wu, A. I.
529         Su, K. M. Summers, C. K. Tuggle, A. L. Archibald, and D. A. Hume. 2012. A gene
530         expression atlas of the domestic pig. *BMC Biology* 10: 1-22.
531   29.   Sharp, G. C., J. L. Hutchinson, N. Hibbert, T. C. Freeman, P. T. K. Saunders, and J.
532         E. Norman. 2016. Transcription Analysis of the Myometrium of Labouring and Non-
533         Labouring Women. *PLoS ONE* 11: e0155413.
534   30.   Fairfax, B. P., P. Humburg, S. Makino, V. Naranbhai, D. Wong, E. Lau, L. Jostins, K.
535         Plant, R. Andrews, C. McGee, and J. C. Knight. 2014. Innate Immune Activity
536         Conditions the Effect of Regulatory Variants upon Monocyte Gene Expression.
537         *Science (New York, N.Y.)* 343: 1246949-1246949.
538   31.   Kauffmann, A., R. Gentleman, and W. Huber. 2009. arrayQualityMetrics—a
539         bioconductor package for quality assessment of microarray data. *Bioinformatics* 25:
540         415-416.
541   32.   Carvalho, B. S., and R. A. Irizarry. 2010. A framework for oligonucleotide microarray
542         preprocessing. *Bioinformatics* 26: 2363-2367.

543    33.    Stefan Hachul, M. J. 2007. Large-Graph Layout Algorithms at Work: An
544           Experimental Study. *Journal of Graph Algorithms and Applications* 11: 345--369.
545    34.    Enright, A. J., S. Van Dongen, and C. A. Ouzounis. 2002. An efficient algorithm for
546           large-scale detection of protein families. *Nucleic Acids Research* 30: 1575-1584.

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572 **Table 1**

573 Summary of the most widely used immune signatures and deconvolution methods

| Authors | Year | Signature derived from | Deconvolution method | No. of Cell types | Total no of markers | Cell types (unique genes) |
|---|---|---|---|---|---|---|
| Abbas et al. | 2005 | Isolated immune cells | No deconvolution algorithm | 6 | 959 unique genes | B Cell (91), Dendritic Cell (70), Lymphoid (234), Monocyte (82), Myeloid (344), Neutrophil (45), NK Cell (17), T Cell (76) |
| Palmer et al. | 2006 | Isolated immune cells | No deconvolution algorithm | 4 | 1146 unique genes | B cells (427), T cells (241), Granulocytes (411), Lymphocytes (67) |
| Abbas et al. | 2009 | Isolated immune cells | Linear least-squares fits | 17 | 359 Affy u133a probes | Resting helper T cells, Activated helper T cells, Resting cytotoxic T cells, Activated cytotoxic T cells, Resting B cells, Activated B cells, BCR-ligated B cells, IgA/IgG memory B cells, IgM memory B cells, Plasma cells, Resting NK cells, Activated NK cells, Monocytes, Resting dendritic cells, Activated Monocytes, Activated dendritic cells, Neutrophils |
| Nicholas et al. | 2009 | Isolated immune cells | No deconvolution algorithm | 8 | 1842 unique genes | T cells (48), Monocytes (186), B cells (218), NK cells (75), Granulocytes (757), Erythroblast (299), Megakaryocyte (262) |

| Gong et al. | 2011 | Isolated immune cells | Quadratic Programming | Uses signature from other studies | | |
|---|---|---|---|---|---|---|
| Zhong et al. | 2013 | Isolated immune cells | Linear model & Quadratic Programming | Uses signature from other studies | | |
| Newman et al. | 2015 | Isolated immune cells | Support vector machine | 22 | 547 unique genes | B cells naïve, B cells memory, Plasma cells, T cells CD8, T cells CD4 naïve, T cells CD4 memory resting, T cells CD4 memory activated, T cells follicular helper, T cells regulatory (Tregs), T cells gamma delta, NK cells resting, NK cells activated, Monocytes, Macrophages M0, Macrophages M1, Macrophages M2, Dendritic cells resting, Dendritic cells activated, Mast cells resting, Mast cells activated, Neutrophils, Eosinophils |

574

575

576

577

578

579

580     **Figure Legends**

581     **Figure 1: Derivation and application of blood and tissue *ImSig*. A)** Flow chart depicts the

582     systematic derivation of *ImSig*. The transcriptome of isolated immune cells was subjected to

583     differential gene expression analysis or correlation analysis to derive a preliminary list. This

584     was further refined using the cluster model algorithm to define the blood and tissue-specific

585     immune signatures (*ImSig*). **B)** Application of signatures involves running *ImSig* scoring

586     algorithm on any transcriptomic data to identify the different immune cells present within the

587     samples followed by network analysis to study the genes that are correlated best with the core

588     signature genes.

589     **Figure 2: Cluster model algorithm refinement and *ImSig* algorithm. A)** The plots

590     represents the outcome of running the cluster model algorithm over a blood and a tissue

591     dataset. Each node represents a unique gene and plotted as a function of its median

592     correlation value within the signature. Blue colour represents the genes that were kept and red

593     represents the genes that were discarded after running the algorithm. The algorithm was

594     applied to eight blood and eight tissue datasets (only 2 shown above). All the blue nodes were

595     then pooled to identify the most commonly occurring genes across datasets, which then

596     formed the basis of defining *ImSig*. **B and C**, Line plots showing 'initial score' calculated for

597     every correlation cut-off between 0.50 and 0.99 while calculating the *ImSig* score. For **B)**

598     microarray dataset (heart attack, GSE48060), the threshold line is drawn at 20 and for **C)**

599     RNA-seq dataset (Brucellosis; E-GEOD-69597), the threshold line is drawn at 10. **D)** Plots

600     showing the effect of loss of signature genes on *ImSig* score. These were calculated by

601     performing a permutation analysis of removing signature genes randomly.

602     **Figure 3: Deconvolution of blood and tissue datasets**. **A)** Correlation network of gene

603     expression data from blood samples of patients with type I diabetes mellitus represented and

604     **B)** samples from breast cancer patients. Each cluster represents a unique cell type. Nodes

605     derived from other signatures which were included in the graph but did not cluster are

606     reduced in size. Histogram plots represent the average expression profile of the *ImSig*

607     signatures across samples.

608     **Figure 4: Network graph to highlight a few closely correlated immune related genes**

609     **with *ImSig*. A)** Correlation network of gene expression data from trachomatis infection

610     (GSE20436). The nodes represent unique genes and the *ImSig* genes are coloured to highlight

611     the immune cluster. **B)** A close up of the immune cluster. The *ImSig* related genes are

612     coloured to represent different immune cell types, while the remaining genes are reduced in

613     node size. We highlight a few well known immune modulatory genes with a greater node size

614     and marking their gene symbols alongside. **C)** Bar plots represents the average expression

615     intensity of individual genes across samples. The top panel (Green) plots represents a few

616     marker genes to understand macrophage biology and the bottom panel (dark grey) to
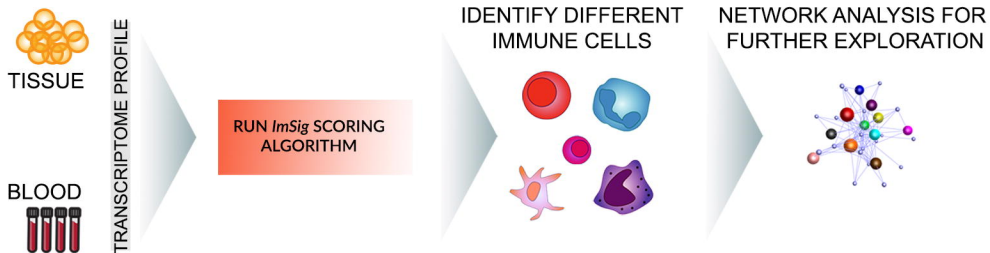
617     understand the T cell biology.

618     **Figure 5: Comparison of *ImSig* with CIBERSORT. A)** Comparison performed using a

619     blood dataset. The boxplots show the relative abundance of immune in cells in the two patient

620     groups computed by CIBERSORT and *ImSig*. The actual median cell count for the four

621     immune cell types were (high, low) Neutrophils (2655, 6160), T cells (617.5, 1988), B cells

622     (35, 293) & NK cells (22.5, 176.5). Significant difference was observed for T cells,

623     Neutrophils and NK cells using CIBERSORT while all differences seen in *ImSig* including B

624     cells are significant (P value <0.05). **B)** Comparison performed using a tissue dataset. The

625     boxplots show the relative abundance of immune in cells in the three different patient groups

626     computed by CIBERSORT and *ImSig*. Significant difference was observed only for

627     macrophages and T cells using CIBERSORT while all differences seen in *ImSig* are

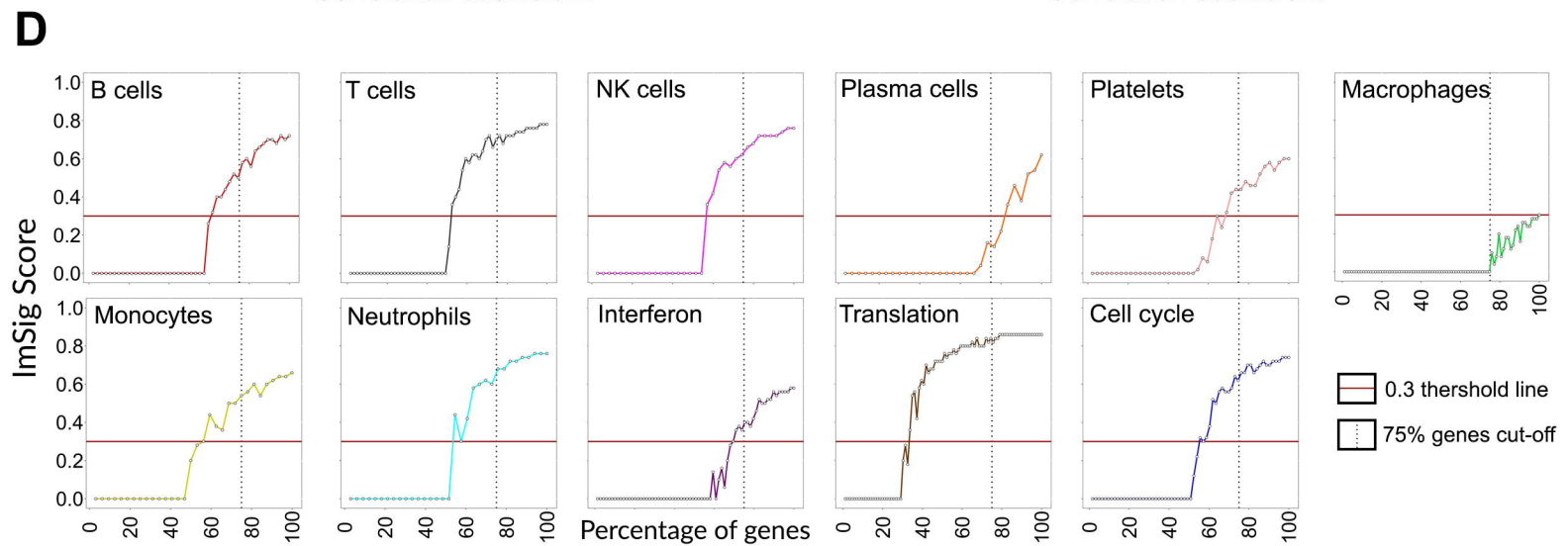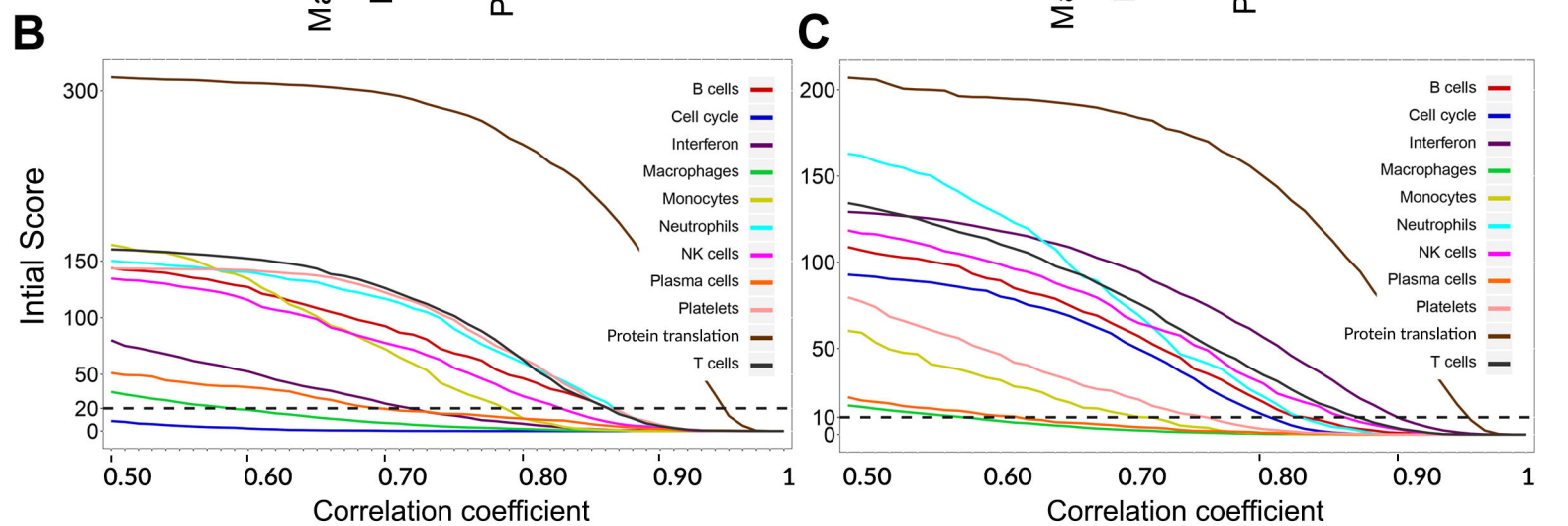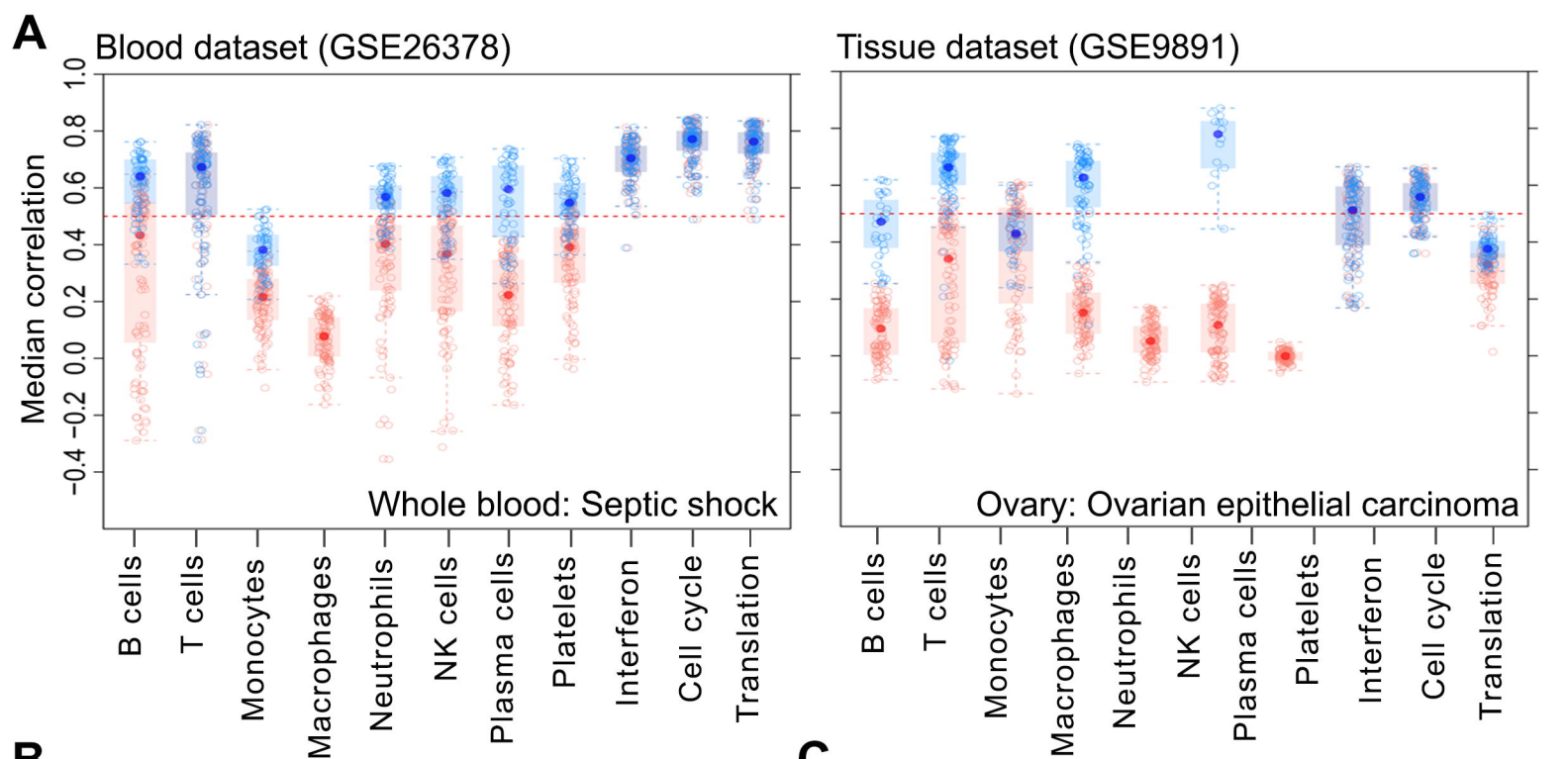628     significant (P value <0.05).

**A** Derivation of *ImSig*

Isolated immune cells → DIFFERENTIAL EXPRESSION → TOP 100 DEG'S → REFINING IN → BLOOD

Tissue samples → CORRELATION ANALYSIS → CELL TYPE SPECIFIC CLUSTERS → TISSUE

CLUSTER MODEL ALGORITHM → BLOOD & TISSUE SPECIFIC *ImSig*

**B** Application of *ImSig*

TISSUE / BLOOD → TRANSCRIPTOME PROFILE → RUN *ImSig* SCORING ALGORITHM → IDENTIFY DIFFERENT IMMUNE CELLS → NETWORK ANALYSIS FOR FURTHER EXPLORATION

**A**

Blood dataset (GSE26378)    Tissue dataset (GSE9891)

Median correlation

Whole blood: Septic shock    Ovary: Ovarian epithelial carcinoma

B cells, T cells, Monocytes, Macrophages, Neutrophils, NK cells, Plasma cells, Platelets, Interferon, Cell cycle, Translation

**B**

Intial Score

Correlation coefficient

B cells
Cell cycle
Interferon
Macrophages
Monocytes
Neutrophils
NK cells
Plasma cells
Platelets
Protein translation
T cells

**C**

Correlation coefficient

B cells
Cell cycle
Interferon
Macrophages
Monocytes
Neutrophils
NK cells
Plasma cells
Platelets
Protein translation
T cells

**D**

ImSig Score

B cells    T cells    NK cells    Plasma cells    Platelets    Macrophages

Monocytes    Neutrophils    Interferon    Translation    Cell cycle

Percentage of genes

0.3 thershold line
75% genes cut-off

**A** Trachoma Dataset

Immune cluster

r= 0.8
Nodes: 7474
Edges: 295,295

**B**

IL22
INTERFERON
CCL20
CXCL11
GZMB
CXCL10
NK CELLS
CELL CYCLE
LAG3
IFNG
IL21
MONOCYTES
IL10
STAT4
CXCR3
ITCH
IL17A
CTLA4
IL21R
IL23R
FOXP3
IL1B
TNF
NOS2
CCR4
IL4R
CD80
PLASMA CELLS
NEUTROPHILS
MACROPHAGES
B CELLS
CCR10

**C**

NOS2 · TNF · CCL17 · IL1RN · IFNG · LAG3 · CD44 · FOXO3

Expression intensity

A- Control   B- Trachoma Patients (*C. trachomatis* -ve)   C- Trachoma Patients (*C. trachomatis* +ve)

**A. Blood dataset (GSE49454)**

B cells | T cells | Neutrophils | NK cells

a- Low cell count
b- High cell count

**B. Tissue dataset (GSE20436)**

B cells | T cells | Neutrophils | NK cells

Plasma cells | Monocytes | Macrophages

CIBERSORT
*ImSig*

a- CONTROL
b- TRACHOMA PATIENTS
(*C. trachomatis* -ve)
c- TRACHOMA PATIENTS
(*C. trachomatis* +ve)