1 # Approaches for integrating heterogeneous RNA-seq data reveals cross-talk
2 # between microbes and genes in asthmatic patients

3

4 Daniel Spakowicz*[1,2,3,4], Shaoke Lou*[1], Brian Barron[1], Tianxiao Li[1], Jose L Gomez[5], Qing Liu[5], Nicole Grant[5],
5 Xiting Yan[5], George Weinstock[2], Geoffrey L Chupp[5], Mark Gerstein[1,6,7,8]

6

7 [1] Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT
8 [2] The Jackson Laboratory for Genomic Medicine, Farmington, CT
9 [3] Division of Medical Oncology, Ohio State University College of Medicine, Columbus, OH
10 [4] Department of Biomedical Informatics, Ohio State University College of Medicine, Columbus, OH
11 [5] Section of Pulmonary, Critical Care, and Sleep Medicine, Department of Internal Medicine, Yale University
12 School of Medicine, New Haven, CT
13 [6] Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT
14 [7] Department of Computer Science, Yale University, New Haven, CT
15 [8] Department of Statistics and Data Science, Yale University, New Haven, CT
16 * These authors contributed equally

17

18 **ABSTRACT (337 words)**

19 Sputum induction is a non-invasive method to evaluate the airway environment, particularly for asthma. RNA
20 sequencing (RNAseq) can be used on sputum, but it can be challenging to interpret because sputum contains
21 a complex and heterogeneous mixture of human cells and exogenous (microbial) material. In this study, we
22 developed a methodology that integrates dimensionality reduction and statistical modeling to grapple with the
23 heterogeneity. We use this to relate bulk RNAseq data from 115 asthmatic patients with clinical information,
24 microscope images, and single-cell profiles. First, we mapped sputum RNAseq to human and exogenous
25 sources. Next, we decomposed the human reads into cell-expression signatures and fractions of these in each
26 sample; we validated the decomposition using targeted single-cell RNAseq and microscopy. We observed
27 enrichment of immune-system cells (neutrophils, eosinophils, and mast cells) in severe asthmatics. Second,
28 we inferred microbial abundances from the exogenous reads and then associated these with clinical variables -
29 - e.g., *Haemophilu*s was associated with increased white blood cell count and *Candid*a, with worse lung
30 function. Third, we applied a generative model, Latent Dirichlet allocation (LDA), to identify patterns of gene
31 expression and microbial abundances and relate them to clinical data. Based on this, we developed a method
32 called LDA-link that connects microbes to genes using reduced-dimensionality LDA topics. We found a number
33 of known connections, e.g. between *Haemophilus* and the gene IL1B, which is highly expressed by mast cells.
34 In addition, we identified novel connections, including *Candida* and the calcium-signaling gene CACNA1E,
35 which is highly expressed by eosinophils. These results speak to the mechanism by which gene-microbe
36 interactions contribute to asthma and define a strategy for making inferences in heterogeneous and noisy
37 RNAseq datasets.

38 **INTRODUCTION**

39 **Linking high-dimensional, heterogeneous datasets**

40 RNA sequencing (RNAseq) has become a standard method of analyzing complex communities. Depending on
41 the sample type, these data can be very heterogeneous. A key problem tackled in this paper is dealing with the
42 heterogeneity and noise in RNAseq data in complex samples such as sputum. This can be appreciated by

43  comparing sputum RNAseq to a more traditional experiment, e.g. blood RNAseq, where the sample can be
44  collected consistently and that contains relatively well-defined cell types (Figure 1). In blood, the vast majority
45  of RNAseq reads align to the human genome, and the goal is often to relate the expression of the genes to a
46  phenotype. By contrast, sputum may be less consistently collected, its cell types are less defined, and it may
47  contain RNA from microbes and other organisms that act as cryptic indicators of the environment. This
48  combination of variables and dimensions often requires researchers to collapse the dimensions to
49  appropriately de-noise the analysis. Here, we present such a strategy that uses a number of supervised and
50  unsupervised techniques such as single-cell signatures and latent Dirichlet allocation (LDA). These techniques
51  can produce a low-dimensional representation of common groups of genes, microbes, or other features that
52  tend to increase or decrease in abundance together. Our approach is useful when the heterogeneity comes
53  from the sample type (e.g., sputum) and especially when the samples derive from a heterogeneous population
54  of individuals, such as patients with asthma.

## Interactions between the host and microbes in the lung

56  Asthma is a disease of the airway that can present with many clinical phenotypes. Much work has focused on
57  identifying subgroups of the disease and how each subgroup responds to treatment. For example, Yan et al.
58  introduced transcriptional endotypes of asthma and the Severe Asthma Respiratory Phenotype consortium
59  defined five subtypes of asthma [1]. Some of these subgroups respond differently to environmental and
60  microbial triggers, such as fungal spores. Some fungi have well-defined effects in asthma, but the role of many
61  microbes remains contentious. A simplified model assigns microbes to one of three categories: pathogenic
62  organisms that cause inflammation, beneficial organisms that reduce inflammation, and those that have no
63  effect on inflammation. The majority of the organisms in the lungs are expected to have no effect, and severe
64  asthmatics are expected to have more pathogenic and fewer beneficial microbes.

## Inferring immune cell fractions from RNAseq data

66  The pathology of microbes is often inferred by the number and type of immune cells observed in samples, such
67  as sputum total leukocyte counts [2, 3]. A standard method for counting immune cells in sputum samples uses
68  microscopy, but the resolution is limited to a few cell types [4]. Other cell-counting methods such as flow-
69  sorting can be challenging because of the viscosity and highly variable cell numbers in sputum. An alternative
70  strategy uses cell-type specific expression patterns to deconvolve RNAseq reads from mixtures of cells into
71  fractions of different immune cells [5]. This deconvolution also effectively de-noises heterogeneous datasets by
72  greatly reducing the number of dimensions. Importantly, the RNA needed for this analysis can be purified
73  without poly-A enrichment– here, we use human ribosomal RNA knockdown – which allows for the
74  simultaneous analysis of microbial and human transcripts.

## Supervised deconvolution and the microbiome

76  While deconvolution to cell fractions effectively de-noises human RNAseq data, an equivalent method does not
77  exist for microbes. Although we can map microbe reads onto their genomes, this approach is imperfect
78  because the genome databases are incomplete and assigning a read to a single genome can be complicated if
79  it matches more than one equally well. One can reduce the dimensions by collapsing microbial strains to
80  different taxonomic ranks (e.g., genus or family); however, taxonomy is notoriously imprecise at defining
81  behavior. For example, many bacteria in the genus *Escherichia* are human commensals, whereas *Escherichia*
82  *coli* OH157:H7 causes hemorrhagic colitis. Alternatively, one can group sequences by the metabolic pathways
83  observed, although this requires high-depth sequencing. Here, we propose a method to reduce the

2

84  dimensionality of microbes by first linking the microbes to human genes, and then applying the relatively well-
85  defined gene dimensionality-reduction methods (e.g., deconvolution to cell types).
86
87  In this paper, we use RNAseq of sputum samples from asthmatic patients to demonstrate dimensionality-
88  reduction strategies and identify microbe-host relationships. We map RNAseq reads onto human or microbial
89  genomes and relate the resulting abundance matrices to each other and to clinical data. Further, we
90  deconvolve the human reads into fractions of the various cell types that make up sputum. Finally, we relate the
91  human genes and microbes using a method we call LDA-link, which identifies relationships between genes,
92  microbes, and cell types. These methods represent a general strategy for dealing with heterogeneous RNAseq
93  data that is applicable to other sample types beyond sputum.
94


95  **RESULTS**


96  **Sequencing and processing with the extracellular RNA processing toolkit (exceRpt) pipeline**

97  We collected induced sputum samples from 115 patients with heterogeneous asthma phenotypes and
98  sequenced these sample using RNAseq. The median read depth per sample was 47.5 million, which meets
99  depth recommendations for analyses of this type [6]. We processed these reads through the exceRpt pipeline
100 [7], which conservatively matches reads to genomes in a sequential order designed to reduce experimental
101 artifacts. In brief, we first aligned the quality filtered reads to the UniVec database of common laboratory
102 contaminants [2], and then aligned the remaining reads to human ribosomal sequences before aligning them to
103 the human genome. We excluded samples with a low ratio of transcript alignments to intergenic sequence
104 alignments, and then aligned the remaining reads to the comparably large sequence space of non-human
105 genomes. We first aligned reads to the relatively well-curated ribosomal databases of bacteria, fungi, and
106 archaea (e.g., Ribosomal Database Project[3]) and then to curated genomes of bacteria, fungi, viruses, plants,
107 and animals. The percent of reads mapping to different biotypes was highly heterogeneous; a median of 60%
108 of the reads aligned to the human reference genome and 50% to annotated transcripts (Figure 1, green bars).
109 A median of 0.7% of the input reads aligned to exogenous sources, with some samples containing as much as
110 28.1% exogenous reads. As a control, we applied the same protocol to blood samples, which demonstrated
111 more homogeneity than sputum (Figure 1, top, "blood").


112 **Overview of the analysis approach**

113 The goal of the analysis was to infer meaningful relationships between the numbers and origins of the RNAseq
114 reads and relate them to clinical phenotypes. We conceptualized the clinical information and RNAseq
115 alignments as a series of tables (**Figure 1**). The clinical table includes patient data collected at the clinic, **C,**
116 including age, weight, lung function tests, etc, with rows indexed by patient ($p$) and roughly 200 clinical
117 variables ($N_c$). Alignments to human protein-coding regions created the gene table, **G**, with $N_p$ rows, as above,
118 and roughly 20,000 genes ($N_g$). Alignments to exogenous genomes created the microbe table (**M**) with $N_p$
119 rows and roughly 1,000 microbes ($N_m$). Given these three tables (**C, G,** and **M**), the basic analysis framework
120 is to correlate columns or rows within or between tables. We represent this by a matrix of correlations, $\mathbf{R}(X_{\cdot,i,}$
121 $Y_{\cdot,j})$, where $X_{\cdot,i}$ is the $i$th column of table **X** and $X_{\cdot,j}$ is the $j$th column of table **Y**. This correlation is summed over
122 the other index, usually $p$. For example, we test the relationship between age and the abundance of each
123 microbe $\mathbf{R}(C_{\cdot,age}, M_{\cdot,m})$ across all patients. Similarly, we correlate the expression of a gene (e.g., *TLR4)* with
124 microbe *Candida* $\mathbf{R}(G_{\cdot,TLR4}, M_{\cdot,Candida})$.
125

126  Individual correlations can be difficult to interpret, particularly in heterogeneous, sparse, or noisy datasets.

127  Organizing the genes into relevant pathways or cell types can reduce the dimensionality and de-noise the

128  analysis. To this end, we deconvolved $G$ ($N_p \times N_g$) into a cell-type fraction table, $F$ ($N_p \times N_f$), and a cell-type

129  signatures table, $S$ ($N_f \times N_g$). However, an analogous supervised method does not exist for the microbes.

130  Therefore, we applied an unsupervised dimensionality-reduction approach, latent dirichlet allocation (LDA),

131  which provides a topic distributions in patients ($\theta^G$, $N_p \times N_k$) across a smaller number ($N_k$=10) of topics and

132  gene topic (and $\varphi^G$, $N_k \times N_g$). This can also be done to the microbe table M and get $\theta^M$ and $\varphi^M$, and the gene

133  and microbe topic can be correlated (e.g. $R(\theta^G_{:,g}, \theta^M_{:,m})$ over all patients).

134

135  The framework described above is useful for identifying linear relationships, but non-linear relationships are

136  also possible. For example, a microbe sensed by a human immune cell could lead to the activation of a

137  transcription factor and the expression of several genes, each of which would have a non-linear relationship to

138  microbe abundance. To identify such relationships, we applied a non-linear ensemble learning algorithm [8, 9],

139  using the de-noised inputs for each gene and microbe ($\varphi^G$ and $\varphi^M$). We call this method LDA-link. Further, we

140  relate the gene and microbe links identified to cell fractions and thereby relate how the host is responding to

141  microbes with regards to immune cell type response with a particular gene.

## Analysis of human-aligned reads

143  Working toward the hypothesis that we can conceptualize human-aligned sputum RNAseq reads as a mixture

144  of immune cell types, each with a distinct expression profile, we deconvolved the Gene table ($G$) into a table of

145  fractions of component cells type ($F$) and cognate cell-type signatures ($S$) by solving the formula $G \sim F * S$.

146  This method relies on knowing the signature gene-set in each cell type, which derived from the blood immune

147  cell high quality profiles. To validate that we could apply these cell expression profiles to sputum, we generated

148  several additional datasets including single-cell RNAseq (scRNAseq), microscopy, and unsupervised

149  decomposition, and then compared the results to the deconvolution table $F$. (Figure 2A, schema).

## Evaluation of deconvolution results by scRNAseq

151  First, we performed scRNAseq on a cohort of similar sputum samples (five control and five asthmatic patients).

152  The single-cell sequences clustered into four groups (Figure 2B, first and second panels). To determine

153  whether the reference profiles that we used to deconvolve the bulk RNAseq recapitulate those found in the

154  single-cell clusters, we co-clustered the reference profiles with the scRNAseq data (Figure 2B, third panel).

155  The reference profiles split into the groups by lineage; for example, those in the lymphoid progenitor line co-

156  clustered with cluster 2, and the myeloblast progenitor line co-clustered with cluster 4. This result suggests that

157  the reference profiles accurately represent the cell types in sputum. The myeloid lineage cluster showed a

158  significant difference in the number of cells between asthmatics and controls (Figure 2C). From this analysis,

159  we concluded that (1) the blood-derived cell profiles appropriately fit the sputum cell types and (2) no additional

160  cell types are needed to deconvolve the sputum bulk RNAseq data.

## Evaluation of deconvolution results by microscopy

162  Second, we evaluated a subset of the samples by microscopy and manually counted the number of

163  neutrophils, eosinophils, lymphocytes, and macrophages. We found good agreement with $F$, when cell counts

164  could be directly compared, i.e. neutrophils and eosinophils were both present in $F$ and counted by

165  microscopy. In cases where the deconvolution method gave higher resolution, (e.g., M0, M1, and M2

4

166  macrophages versus one type of macrophage by microscopy), the aggregation of the relevant columns in $F_f$
167  correlated well with the microscopy counts (Figure 2D).
168

### Association of cell fractions with clinical features

170  Having validated the deconvolution of sputum samples (table $F$), we then correlated the cell fractions with
171  clinical features ($R(F_{.,f}, C_{.,c})$ for all patients). We found that the changes in fractions of several cell types were
172  highly correlated with clinical features (Figure 2E). For example, the fraction of T-regulatory cells negatively
173  correlated with the number of hospitalizations per year, suggesting a beneficial role of these cells in the
174  management of asthma.
175

### Evaluation of deconvolution results by unsupervised decomposition

177  We compared the signal captured by cell-type deconvolution to an unsupervised decomposition method: LDA.
178  Using LDA, we factored the gene expression table into ten topics that conceptually represent gene expression
179  programs. This resulted in a gene-topic-fraction-in-patients table, $\theta^G$ ( $N_p \times N_k$) with $N_k$=10 topics, as well as
180  corresponding gene-topic table, $\varphi^G (N_k \times N_g)$, that are analogous to the supervised deconvolution tables $F$ and
181  $S$. We correlated the cell-type fractions table with the gene topics fraction table ($R(F_{.,f}, \theta_{.,k})$ for all patients, and
182  found agreement between LDA and the cell-signature-based deconvolution for only the most prominent cell
183  type, neutrophils (Figure 2D, topic 4). The top genes associated with topic 4 were enriched in the neutrophil
184  chemotaxis pathway (Figure S8 B).
185

186  However, the remaining topics were comprised of multiple cell types. This suggests that LDA can identify
187  distinct but partially overlapping features in $G$. According to the clustering of $\theta^G$ , a subgroup of severely
188  asthmatic patients was highly correlated with topic four (**Figure S8A**). The top-weighted genes in topic 4 were
189  enriched for the pathways "neutrophil chemotaxis" and "asthma-related genes" (**Figure S8B**). These pathways
190  were not enriched in the analogous cell-type-signatures table $S$, suggesting that LDA topics are distinct from
191  the cell-type signatures, but are also clinically relevant. Moreover, the top-weighted genes in topic 1 of the
192  gene topic components table were mitochondrial genes, and topic 1 was strongly correlated with age. This link
193  shows strong support in the literature, as reactive oxygen species produced by the mitochondria reduce their
194  function over time [10]; however, we did not observe this relationship for any cells in the cell-type-fractions
195  table ($F$). Another method using a very different algorithm than LDA, non-negative matrix factorization (NMF),
196  showed strong agreement with LDA (Figure S2, Nmf.1). This supports the use of supervised deconvolution
197  methods as picking out interpretable signals that are different than those identified by unsupervised methods.
198  Unsupervised decomposition should be considered a set of features distinct from those found through
199  deconvolution.
200

### Analysis of exogenous reads

202  After filtering out contaminants and human reads, we assembled the set of reads that aligned to exogenous
203  genomes into a Microbe table (**M**). The exogenous sequences aligned to mostly bacteria and fungi, although
204  we also observed a few arthropod and helminth reads (**Supplemental Table X**). The dominant phyla observed
205  were from the bacterial kingdom: Proteobacteria, Firmicutes, and then Bacteroidetes. The abundance of
206  Proteobacteria is in contrast to observations from the gut where Bacterioidetes predominate [11]. Also notable

207  was the presence of two phyla of fungi among the eight most abundant overall, although this was in lower
208  abundance than many of the bacterial phyla.


209  **Microbes correlations with clinical information and cell fractions**

210  We correlated the microbe abundances to clinical information ($R(M_{\cdot,m}, C_{\cdot,c})$ for all patients) (**Figure 3A**).
211  *Haemophilus* was associated with increased total white blood cell numbers, as has been described previously
212  [12]. *Candida* was associated with worse lung function test results (e.g., forced expiratory volume and forced
213  vital capacity), which supports the association with a severe form of asthma characterized by eosinophilia [13].
214

215  We next correlated microbe abundances to human immune cell fractions ($\mathbf{R}(M_{\cdot,m}, F_{\cdot,f})$ for all patients) (**Figure
216  3B**). Several correlations demonstrated results with strong literature precedence. For example, studies have
217  previously shown that *Haemophilus* associates with eosinophilia [14], and we observed a significant correlation
218  between *Haemophilus* and the fraction of eosinophils. We also observed a significant correlation between
219  *Haemophilus* and activated mast cells, suggesting an alternative route to *Haemophilus*-induced inflammation
220  [15]. Moreover, the fungal genus *Candida* was also significantly correlated with eosinophils, even more
221  strongly than *Haemophilus*. Pulmonary candidiasis has long been associated with allergic bronchial asthma
222  and inflammation [16], however few lung microbiome studies have examined both bacterial and fungal signals.
223  This highlights the need for a more comprehensive search of the lung microbiome and demonstrates the power
224  of an RNAseq-based method that can report on all kingdoms with the same sample preparation.


225  **Dimensionality reduction for microbes: clustering and networks**

226  We attempted to de-noise the microbe table ($M^{phylum}$) with a variety of dimensionality-reduction techniques.
227  First, we collapsed the microbes by taxonomy, grouping them to the rank of phylum, and then hierarchically
228  cluster the patients based on their phylum abundance (Figure 3C **HierClust**($M^{phylum}$)). The hierarchical
229  clustering showed that the phylum distributions formed three clusters of patients. We related these clusters to
230  the clinical variable "asthma severity" and observed that cluster 2 was enriched for patients identified as having
231  moderate or severe asthma. This cluster was characterized by the highest relative abundance of the phylum
232  Proteobacteria (Figure 3C). Notably, the genus Haemophilus belongs to this phylum, consistent with the
233  correlations observed at the genus rank (Figures 3A, 3B).
234

235  Similarly, we could de-noise the microbe table using a co-abundance network, by correlating the genus-level
236  abundances ($\mathbf{R}(M_{\cdot,m}, M_{\cdot,m})$ and identifying significant modules (**Supplemental Figure Z**). An interpretation of
237  these modules is that they define metabolic niches, where microbes either directly compete for metabolites or
238  there is interdependency in metabolite production. Such networks could be created from other tables, such as
239  the topic distribution of microbes ($\mathbf{R}(\varphi_{\cdot,m}^{M}, \varphi_{\cdot,m}^{M})$ for all the topics) (**Figure 3D**). These modules represent
240  another unit that could be related to the clinical information (**C**) and the cell-type fractions (**F**).


241  **LDA-link for the identification of links between genes and microbes**

242  How much cross-talk exists between microbes and human cells in the airway remains contentious [17]. We feel
243  this is partly due to the heterogeneous and noisy data from airway samples, where it is often difficult to find
244  strong correlations using standard algorithms. We therefore sought to link genes to microbes via a new method
245  called LDA-link.
246

247  LDA-link connects genes to microbes using a combination of linear correlation, unsupervised decomposition
248  and an ensemble learning classifier. We hypothesized that the only strongest gene and microbe correlations

249 would be observable through the noise in the RNAseq data. Therefore, we used these strong links as a training
250 set to find other links, after taking steps to reduce the noise in the data. We reduced the noise using LDA and
251 then identified links using a random forest classifier, described in more detail below and in the methods
252 section.
253
254 To define the training, set we first related columns between the gene and microbe tables ($\mathbf{R}(G_{\cdot,g}, M_{\cdot,m})$),
255 yielding many low-scoring correlations. However, a relatively small number were strong ($\mathbf{R} > 0.4$) and highly
256 significant (p < 1E-5 after FDR correction) (**Figure 4A**). We selected the very strong correlations as true-
257 positive links between genes and microbes in the training set, and non-correlated pairs ($-0.05 < \mathbf{R} < 0.05$) as
258 true-negative links. The genes involved in these strong correlations were enriched for pathways related to
259 microbial interactions in the airway, including "Asthma & Bronchial Hypersensitivity" and "Respiratory Syncytial
260 Virus Bronchiolitis" (**Figure 4B**), suggesting that the small set of strong linear correlations were relevant to
261 asthma.
262
263 Next, we trained a random forest classifier on the linear correlations described above. To reduce the noise in
264 the data, the features used as inputs to the classifier were the LDA topics for each gene and microbe
265 ($\varphi_{\cdot,g}^{G}, \varphi_{\cdot,m}^{M}$). That is, for each gene-microbe pair, we concatenated the gene and microbe topics into a single
266 vector (length 20). The Gini index showed the most important features in defining links between genes and
267 microbes were gene topics #7 and #8, and microbe topic #1 (**Figure 4 C-F**). The genes that comprise the most
268 influential gene topic #8, are enriched for the pathway "Inflammatory Response", and specifically the cytokines
269 IL2 and IL6. It is tempting to speculate that these genes are strong predictors of a link between genes and
270 microbes because they indicate when the presence of a microbe has triggered an inflammatory response.

271 **Cross-talk between genes and microbes defined by LDA-link**

272 LDA-link identified connections between genes and microbes reported elsewhere in the literature as well as
273 novel observations. A bipartite graph summarizes a subset of the connections, showing in most cases several
274 genes linked to each microbe (**Figure 5A**, for a complete list see **Supplemental Table X**). Notably, both fungi
275 and bacteria showed these links, further highlighting the need to evaluate more than bacteria when performing
276 microbiome experiments in the airway. The gene lactotransferrin was linked to *Aeromonas,* which has been
277 associated with gastroenteritis and skin infections and has been previously reported to bind lactoferrin [18].
278 *Burkholderia*, a gram-negative bacterial genus, is recognized as an important pathogen in the mucus-filled
279 lungs of patients with cystic fibrosis; it was linked to gene MUC6, which encodes a secreted protein
280 responsible for the production of mucin [19]. *Haemophilus* was observed to be linked to NFKB Inhibitor Zeta,
281 which is induced by the bacterial cell wall component lipopolysaccharide [20]. In addition, *Haemophilus* was
282 linked to the cytokine interleukin 1 beta (IL1B), an important mediator of the inflammatory response. IL1B
283 hypersensitivity is a hallmark of the asthma phenotype. *Pasteurella* was also linked to IL1B, and its toxin has
284 been shown to induce expression of IL1B [21]. In addition to single gene-microbe pairs, we layered on pathway
285 and cell deconvolution data to identify larger-scale effects of microbes.
286
287 Microbes were linked to genes that are enriched in pathways relating to auto-immunity and inflammation as
288 well as cytokine receptors and their interactions (**Figure 5B**). The microbes associated with cytokine pathways
289 included *Synechococcus*, *Lactococcus, Dialister*, *Psychrobacter*, *Moraxella*, *Brenneria*, *Proteus*, *Haemophilus*,
290 and *Pasteurella*. In addition, we related the cell-type signatures table ($\mathbf{S}_{f,g}$) to identify the immune cell types that
291 are related to each microbe (**Figure 5C**). We observed the *Haemophilus*-IL1B linkage in monocytes and mast
292 cells. Samples containing *Haemophilus* triggered more activated mast cells according to its cell fraction
293 (**Figure 5C inset**) [22-25]. Similarly, the fungal genus *Candida* was linked to the gene GCSAML, which was

294 highly expressed by eosinophils. The presence of *Candida* was associated with increased numbers of
295 Eosinophils in the airway.
296

297 **DISCUSSION**

298 Heterogeneity and noise are common problems in biological datasets. Heterogeneity can derive from mixtures
299 of different cell types, such as in sputum, or from sparsity, such as in microbiome or single-cell RNAseq data.
300 Unsupervised methods of dimensionality reduction can effectively eliminate these issues, but suffer from
301 decreased interpretability. That is, variables are collapsed together for reasons that are often opaque.
302 Supervised dimensionality reduction maintains interpretability because variables are collapsed using prior
303 knowledge, such as the genes in a pathway or the expression patterns of a cell type. Here, we combined
304 unsupervised and supervised approaches to de-noise the data while retaining interpretability.
305

306 The field is increasingly appreciating the role of the airway microbiome in the development of disease.
307 Commensal microbiota have been shown in other contexts to be strong regulators of host immune system
308 development and homeostasis [26]. Disturbances in the composition of commensal bacteria can result in
309 imbalanced immune responses and affect an individual's susceptibility to various diseases, including those that
310 are inflammatory (e.g., inflammatory bowel disease and colon cancer), autoimmune (e.g., celiac disease and
311 arthritis), allergic (e.g., asthma and atopy), and metabolic (e.g., diabetes, obesity, and metabolic syndrome)
312 (reviewed in [27]). Investigating the microbiota in the lower respiratory tract is a relatively new field in
313 comparison to the extensive work on the intestinal tract. In fact, the lung was excluded from the original Human
314 Microbiome Project because it was not thought to have a stable resident microbiome [11]. A limited number of
315 reports have investigated the changes in the lung microbiota between healthy, non-smoking and smoking
316 individuals as well as in patients suffering from cystic fibrosis, chronic obstructive pulmonary disease, or
317 asthma [2, 28-30]. Despite emerging data on the airway microbiota, little is known about the role of the lung
318 microbiome in modulating pulmonary mucosal immune responses. LDA-link can find relationships between
319 microbes and genes and link them to immune cells and their responses.
320

321 The linkages identified here suggest major processes by which lung immune cells respond to microbes. We
322 found that mast cells respond to *Haemophilus* and *Pasteurella* via IL1B and that eosinophils respond to
323 *Candida* via GCSAML. While experimental validation of these linkages is needed, these results represent
324 observations that would be missed by analyses that do not deconvolve RNAseq data into cell fractions, or that
325 analyze only human RNAseq reads. We expect LDA-link to be broadly useful in relating heterogeneous or
326 noisy RNAseq data.

327 **METHODS**

328 **Sample collection and sequencing**

329 Sputum induction was performed with hypertonic saline, the mucus plugs were dissected away from the saliva,
330 the cellular fraction was separated, and the RNA was purified as described previously [1]. Briefly, RNA was
331 purified using the All-in-One purification kit (Norgen Biotek) and its integrity was assayed by an Agilent
332 bioanalyzer (Agilent Technologies, Santa Clara, CA). Ribosomal depletion was performed with the RiboGone-
333 Mammalian kit (Clontech Cat. Nos. 634846 & 634847 ) and cDNA was created with the SMARTer Stranded
334 RNAseq Kit (Cat. Nos. 634836). Samples were sequenced using an Illumina HiSeq 4000 with 2x125 bp reads,
335 with an average of 47.5 million reads per sample.

**RNAseq processing by exceRpt**

An adapted version of the software package exceRpt [7] was used to process the sputum RNAseq data. Briefly, RNAseq reads were subjected to quality assessment using FastQC software v.0.10.1 (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/) both prior to and following 3' adapter clipping. Adapters were removed using FastX v.0.0.13 (http://hannonlab.cshl.edu/fastx_toolkit/). Identical reads were counted and collapsed to a single entry and reads containing N's were removed. Clipped, collapsed reads were mapped directly to the human reference genome (hg19) and pre-miRNA sequences using STAR [31]. Reads that did not align were mapped against a ribosomal reference library of bacteria, fungi, and archaea, compiled by the Ribosome Database Project [32], and then to genomes of bacteria, fungi, plants, and viruses, retrieved from GenBank [32]. In cases where RNAseq reads aligned equally well to more than one microbe, a "last common ancestor" approach was used, and the read was assigned to the next node up the phylogenetic tree, as performed by similar algorithms [7, 33].


**Data tables notation**

We use the following notation to define matrices associated with $p$ patients (115) (Figure 1):

**C**: Clinical table ($N_p \times N_c$), c is the clinical index

**G**: Gene table ($N_p \times N_g$), bulk-RNA seq table before deconvolution,

**M**: Microbe abundance table ($N_p \times N_c$)

**F**: Cell fractions table ($N_p \times N_f$), resulting from the deconvolution of $\boldsymbol{G}_{p,g}$

**S**: Cell signatures table ($N_f \times N_g$), resulting from the deconvolution of $\boldsymbol{G}_{p,g}$

$\boldsymbol{\theta}^G$ : Patient topic table ($N_p \times N_k$) after LDA inference based on gene table $\boldsymbol{G}_{p,g}$

$\boldsymbol{\varphi}^G$ : Gene topic table ($N_k \times N_g$) after LDA inference based on gene table $\boldsymbol{G}_{p,g}$

$\boldsymbol{\theta}^M$ : Patient topic table ($N_p \times N_k$) after LDA inference based on microbe table $\boldsymbol{M}_{p,m}$

$\boldsymbol{\varphi}^M$ : Microbe topic table ($N_k \times N_m$) after LDA inference based on table $\boldsymbol{M}_{p,m}$

**L**: gene microbe linkage table ($N_g \times N_m$) predicted by LDA-link


**Dimensionality Reduction**


Supervised, deconvolution

The gene table (**G**) was deconvolved using the transcriptomes from 22 flow cytometry-sorted and sequenced immune cell types (lm22) using the CIBERSORT tool [5]. Briefly, a pre-defined set of characteristic gene expression patterns for each cell type was used to identify the fraction of each cell type given a mixture of expression by solving for the equation:

$$G = F * E$$

Where **G** is the Gene table of human protein-coding gene expression from the exceRpt pipeline, **F** is the Cell Fraction table, and **E** is the characteristic gene expression calculated within CIBERSORT. Support Vector Regression was used to perform variable selection, reducing the number of characteristic genes used to distinguish cell types and thereby reducing overfitting. The above equation was then solved to provide an estimate of **F**. P-values for the fit of **E** and **F** to **G** demonstrated that all samples were significant at $\alpha = 0.05$.

376  Following the solution of **F**, a Cell Signature table **S** was calculated to estimate the expression of $g$ genes, as
377  opposed to the reduced set appropriate for the characteristic expression evaluation, by solving the equation:
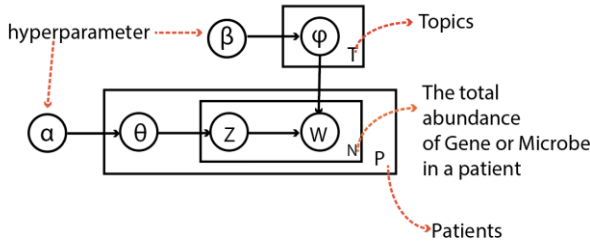378
379                                  **S** = **G** * **F**

380  <u>Decomposition though a generative model</u>
381  The Gene table **G** was decomposed using LDA and Non-negative Matrix Factorization (NMF).
382  For LDA, the abundance values for bulk RNAseq and exogenous RNA were scaled down to reduce
383  computation intensity during sampling. More simply, the RPM expression values were converted to integers,
384  and then divided by 10. The max value was set to 1,000.
385
386  Given each patient ($p$), all of the genes and microbes were treated like corpus of words in the traditional LDA
387  application. The word ($w$) was gene or microbe, and the word count was gene expression or microbe
388  abundances. We built LDA models for genes and microbes, respectively.
389



$$\varphi_k \sim Dirichlet\ (\beta)$$
$$\theta_p \sim Dirichlet\ (a)$$
$$Z_{p,n} \sim Multinomial\ (\theta_p)$$
$$w_{p,n} \sim Multinomial\ (\varphi_{z_{p,n}})$$

394
395  Given $p, w, k, v, N_p, N_w, N_k, N_v, \alpha, \beta, Z, \theta, \varphi, W,$ where $p, w, k, v$ denote a patient, a word in a document, a topic
396  and a word in the corpus respectively; $N_p$ is the number of documents(patients ), $N_w$ is the number of words
397  (gene or microbe) in a document, $N_k$ is the number of topics (set as 10), $N_v$ is the corpus for all the documents;
398  $\alpha$ ($N_k$ dimensional vector)and $\beta$ ($N_v$-dimensional vector) are the hyper parameters for $\theta$ ($N_p \times$
399  $N_k$, the distributoin of topics in documents ) and $\varphi$ ($N_k \times N_v$, the distribution of word for topics)  W is an $N_w$-
400  dimensional vector that denotes the word (gene or microbe expression) in a document (patients). Z is the $N_w$-
401  dimensional vector of integers between 1 and $N_k$ for the topic of word in a document.
402
403  The joint distribution of the LDA model is $\mathcal{P}(Z, W;\ \alpha, \beta)$ and $\varphi$ and $\theta$ are integrated out as:
404

$$\mathcal{P}(Z, W;\ \alpha, \beta) = \int_\varphi \prod_{i=1}^{N_k} \mathcal{P}(\varphi_i; \beta) \prod_{j=1}^{N_p} \prod_{t=1}^{N_w} \mathcal{P}\left(W_{j,t}\middle|\varphi_{Z_{j,t}}\right) d\varphi \int_\theta \prod_{i=1}^{N_p} \mathcal{P}(\theta_i; \alpha) \prod_{j=1}^{N_w} \mathcal{P}(Z_{i,j}|\theta_i) d\theta$$

406

$$= \prod_{k=1}^{K} \frac{\Delta(n_{\cdot,k} + \beta)}{\Delta(\beta)} \prod_{s=1}^{S} \frac{\Delta(n_{s,\cdot} + \alpha)}{\Delta(\alpha)}$$

408  Where $\Delta(\alpha) = \frac{\prod_{k=1}^{K} \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^{K} \alpha_k)}$
409
410  Gibbs sampling equation can be derived from $\mathcal{P}(Z, W;\ \alpha, \beta)$ to approximate the distribution of $\mathcal{P}(Z|W;\ \alpha, \beta)$
411  because $\mathcal{P}(W;\ \alpha, \beta)$ is invariant to $Z$.  Given $Z_{m,n}$ denotes the  topic of the  $n$th word token in the $m$th document,

412 and also assume that its word symbol is the $v$th word in the vocabulary, the conditional probability can be inferred
413 as follows:

414

415

416
$$\mathcal{P}(Z_{m,n} = k | Z_{\neg(m,n)}, W; \alpha, \beta) = \frac{\mathcal{P}(Z,W;\alpha,\beta)}{\mathcal{P}(Z_{\neg(m,n)},W;\alpha,\beta)} = \frac{\mathcal{P}(w,z)}{\mathcal{P}(w_{m,n}, w_{\neg(m,n)}, z_{\neg(m,n)})} = \frac{\mathcal{P}(w,z)}{\mathcal{P}(w_{\neg(m,n)}, z_{\neg(m,n)})} \cdot \frac{1}{\mathcal{P}(w_{m,n}=t)}$$

417
$$\propto \frac{\mathcal{P}(w,z)}{\mathcal{P}(w_{\neg(m,n)}, z_{\neg(m,n)})}$$

418
$$= \frac{\prod_{k=1}^{K} \frac{\Delta(n_{\cdot,k} + \beta)}{\Delta(\beta)} \prod_{p=1}^{P} \frac{\Delta(n_{p,\cdot} + \alpha)}{\Delta(\alpha)}}{\prod_{k=1}^{K} \frac{\Delta(n_{\neg(m,n),k} + \beta)}{\Delta(\beta)} \prod_{p=1}^{P} \frac{\Delta(n_{p,\neg(m,n)} + \alpha)}{\Delta(\alpha)}}$$

419
$$= \frac{\Delta(n_{\cdot,k} + \beta)}{\Delta(n_{\neg(m,n),k} + \beta)} \cdot \frac{\Delta(n_{p,\cdot} + \alpha)}{\Delta(n_{p,\neg(m,n)} + \alpha)}$$

420

421

422 After sampling, the expectation of the θ (doc → topic) and φ(topic → word) matrix can be inferred as follows
423 given the symmetric hyper-parameters $\alpha$ and $\beta$ were used:

424

425
$$\theta_{p,k} = \frac{n_{p,k} + \alpha}{\sum_{i=1}^{K} n_{p,i} + N_k \alpha}$$

426
$$\varphi_{k,v} = \frac{n_{k,v} + \beta}{\sum_{i=1}^{V} n_{k,i} + N_v \beta}$$

427

428 We instantiated the variables θ and φ to $\theta_{p,t}^{G}$, $\theta_{p,t}^{M}$, and $\varphi_{k,g}^{G}$, $\varphi_{k,m}^{M}$, where $\theta_{p,t}^{G}$, $\theta_{p,t}^{M}$ denotes the gene and
429 microbe topic fraction in patient ; $\varphi_{k,g}^{G}$, $\varphi_{k,m}^{M}$ denotes the gene and microbe topic.

430 **Single-cell RNAseq**

431 Sputum cells were separated on a Fluidigm C1 medium-sized channel. The mRNA was purified from
432 approximately 500pg-1ng of total RNA using the Clontech SMARTer Ultra Low RNA Kit and poly-dA-selected
433 using SPRI beads and dT primers. Full-length cDNA was sheared into 200-500bp DNA fragments by
434 sonication (Covaris, Massachusetts, USA), and then indexed and size validated by LabChip GX. Two nM
435 libraries were loaded onto Illumina version 3 flow cells and sequenced using 75bp single-end sequencing on
436 an Illumina HiSeq 2000 according to Illumina protocols. Data were cleaned, processed, aligned, and quantified
437 following the SINCERA pipeline [34].

438 **Pathogen-to-host linkage identification**

439 Microbe relative abundances and gene TPM values were correlated as follows, with $G_{\cdot,i}$ for $i$ gene and $M_{\cdot,j}$ for $j$
440 microbe:

441

442
$$R(i,j) = \frac{\sum_{p=1}^{N_p} (G_{p,i} - \overline{G_{\cdot,i}})(M_{p,j} - \overline{M_{\cdot,j}})}{\sqrt{\sum_{p=1}^{N_p} (G_{p,i} - \overline{G_{\cdot,i}})^2} \sqrt{\sum_{p=1}^{N_p} (M_{p,j} - \overline{M_{\cdot,j}})^2}}$$

443

444 Gene-microbe correlations with $p$-values less than 1e-5 (absolute correlation greater than 0.4) were chosen as
445 the positive links in a training set. Negative links in the training set were defined as an absolute correlation of
446 less than 0.05. This approach resulted in 302 positive and 650,398 negative links. A random forest algorithm
447 was trained on this set, which can accommodate the highly unbalanced dataset as well as potentially identify
448 non-linear links between genes and microbes. Down-sampling and up-sampling techniques were tested but did
449 not significantly improve the model. In the final model, we adopted the upscaling technique and tested it using
450 cross-validation. The positive dataset was upscaled to very high levels. We use 2-fold cross validation to
451 validate the performance. Simply, we randomly select half training data to train the model, and use the
452 remaining records to test the performance and repeat this for ten times. The AUC and AUPR were 0.994 and
453 0.996 on average, respectively.
454

455 **Microbe co-abundance network**
456 The raw abundance $M$ and LDA microbe topic matrices $\varphi^M$ , which represent the microbe's weight to each
457 topic, were generated.
458

459 The correlation network between different microbes was calculated using Pearson correlation. The cutoff to
460 define a co-abundance edge was 0.8 for $R(\varphi^M_{.,m}, \varphi^M_{.,m})$ and 0.3 for $R(, \varphi^M_{.,m}$. The microbe network modules, which
461 were densely connected themselves but sparsely connected to other modules, were clustered based on
462 between-ness [35] and the other algorithms; we also tested label propagation and fast greedy algorithms [36].
463

464 We also compared the LDA topics with the microbes in the same clusters. If a microbe was the top 10 most
465 highly contributed for a topic, then we labeled the topic number in the bracket. Some microbes may have
466 multiple topic labels because they highly contribute to multiple topics.


467 **DATA ACCESS**
468 Sputum bulk-cell RNAseq data can be found under the bioproject SRRXXXXX and sputum single-cell RNAseq
469 data at SRRYYYYYY.


470 **ACKNOWLEDGMENTS**

474


475 **DISCLOSURE DECLARATION**
476 The authors declare no conflicts of interest. [[[check with GC]]]
477


478 **ABBREVIATIONS**
479 LDA: Latent Dirichlet Allocation; PCA: Principle Component Analysis;
480

**FIGURES:**

**Figure 1. RNAseq alignment summary for control and asthmatic sputum, showing fractions of reads that aligned to different biotypes. Alignments to the protein-coding biotype were used to generate the gene expression matrix (G), which was then deconvolved into a cell fraction matrix (F) and cell expression (E). The exogenous reads were used to generate a microbial profile matrix (M). These matrices were then related to the clinical phenotype matrix (P) for biological insight.**


Figure 2. Deconvolution of RNAseq human reads into cell fractions using cell signature deconvolution. A) Schematic showing the imputation of a cell fraction matrix and cell-specific expression matrix. B) Imputed cell fractions were validated using microscopy; Cell fractions were then correlated with SARP cluster for  two major cell type: (C) Machrophases.M0 and (D) Mast cell activiated.  E) the cell fraction of LM22 gene signature deconvolution are correlated with the topic distribution of samples from LDA analysis. F)  G) tSNE analysis and clustering using single cell RNAseq from Asthmatic patient and control. H) The fraction of single cells for different cell types clusters between Asthmatic patient and Controls.


Figure 3. Exogenous RNAseq analysis. (A).  The correlations between microbes abundance and cell fraction based on LM22 signature (B) The correlation between microbes abundance and clinical information. (C) The microbes abundance shows clear patterns that associated with Asthmatic severity. (D) The co-abundance network and overlay with the associated topics of microbes.


Figure 4. Prediction of cross-talk between microbe and gene. (A) The diagram to combine linear and LDA-based non-linear algorithms to identify gene microbe linkages. (B) simple correlation to identify strong linkages between microbes and genes. Gene set over represent analysis for genes. X-axis is the -log(p-value) . (C) the importance of features (LDA topics for gene and microbes) in the RandomForest model by Gini index. The top 20 associated gene in topics 8 (D) and topic 7 (E) of genes, and topic 1 (F) of microbes.


Figure 5. The linkage between microbes and genes reflects the heterogeneity of different cell types. (A). Linkages between microbes and genes. (B) The linkages indicated by the cell proportion of certain types.


**SUPPLEMENTAL INFORMATION**


**Figure S1:**

The distribution of cell fraction of sample in different group.


**Figure S2**

The heatmap between NMF component and cell fraction from LM22.

515    **Figure S3**

516    Overall view of correlation of all the extracellular organism with clinical features.


517    **Figure S4**

518    Heatmap of the correlation of topics (gene and microbe ) with clinical information.


519    **Figure S5**

520    Co-abundance network based on correlation of abundance.\

521

522    Figure S6

523

524    Top associated microbes in microbe topics

525

526    Figure S7

527

528    Top associated genes in gene topics

529

530    Figure S8

531

532        (A) The topic distribution of patient. (B) The gene enrichment analysis of top genes in topic 4.

533

534    Figure S9

535    Main pathways get involved by microbe linked genes.

536


537    **REFERENCES**

538    1.    Yan X, Chu JH, Gomez J, Koenigs M, Holm C, He X, et al. Noninvasive analysis of the sputum
539    transcriptome discriminates clinical phenotypes of asthma. Am J Respir Crit Care Med. 2015;191(10):1116-25.
540    Epub 2015/03/13. doi: 10.1164/rccm.201408-1440OC. PubMed PMID: 25763605; PubMed Central PMCID:
541    PMCPMC4451618.
542    2.    Huang YJ, Nariya S, Harris JM, Lynch SV, Choy DF, Arron JR, et al. The airway microbiome in patients
543    with severe asthma: Associations with disease features and severity. J Allergy Clin Immunol. 2015;136(4):874-
544    84. Epub 2015/07/30. doi: 10.1016/j.jaci.2015.05.044. PubMed PMID: 26220531; PubMed Central PMCID:
545    PMCPMC4600429.
546    3.    Gibson PG, Girgis-Gabardo A, Morris MM, Mattoli S, Kay JM, Dolovich J, et al. Cellular characteristics
547    of sputum from patients with asthma and chronic bronchitis. Thorax. 1989;44(9):693-9. Epub 1989/09/01. doi:
548    10.1136/thx.44.9.693. PubMed PMID: 2588203; PubMed Central PMCID: PMCPMC462047.
549    4.    Belda J, Leigh R, Parameswaran K, O'Byrne PM, Sears MR, Hargreave FE. Induced sputum cell
550    counts in healthy adults. Am J Respir Crit Care Med. 2000;161(2 Pt 1):475-8. Epub 2000/02/15. doi:
551    10.1164/ajrccm.161.2.9903097. PubMed PMID: 10673188.
552    5.    Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of cell subsets
553    from tissue expression profiles. Nat Methods. 2015;12(5):453-7. Epub 2015/03/31. doi: 10.1038/nmeth.3337.
554    PubMed PMID: 25822800; PubMed Central PMCID: PMCPMC4739640.
555    6.    Consortium EP. An integrated encyclopedia of DNA elements in the human genome. Nature.
556    2012;489(7414):57-74. Epub 2012/09/08. doi: 10.1038/nature11247. PubMed PMID: 22955616; PubMed
557    Central PMCID: PMCPMC3439153.

7.      Rozowsky J, Kitchen RR, Park JJ, Galeev TR, Diao J, Warrell J, et al. exceRpt: A Comprehensive Analytic Platform for Extracellular RNA Profiling. Cell Syst. 2019;8(4):352-7 e3. Epub 2019/04/09. doi: 10.1016/j.cels.2019.03.004. PubMed PMID: 30956140.

8.      Welling SH, Clemmensen LK, Buckley ST, Hovgaard L, Brockhoff PB, Refsgaard HH. In silico modelling of permeation enhancement potency in Caco-2 monolayers based on molecular descriptors and random forest. Eur J Pharm Biopharm. 2015;94:152-9. Epub 2015/05/26. doi: 10.1016/j.ejpb.2015.05.012. PubMed PMID: 26004819.

9.      Welling SH, Refsgaard HHF, Brockhoff PB, Clemmensen LH. Forest Floor Visualizations of Random Forests. arXiv e-prints [Internet]. 2016 May 01, 2016. Available from: https://ui.adsabs.harvard.edu/abs/2016arXiv160509196W.

10.     Payne BA, Chinnery PF. Mitochondrial dysfunction in aging: Much progress but many unresolved questions. Biochim Biophys Acta. 2015;1847(11):1347-53. Epub 2015/06/09. doi: 10.1016/j.bbabio.2015.05.022. PubMed PMID: 26050973; PubMed Central PMCID: PMCPMC4580208.

11.     Turnbaugh PJ, Ley RE, Hamady M, Fraser-Liggett CM, Knight R, Gordon JI. The human microbiome project. Nature. 2007;449(7164):804-10. Epub 2007/10/19. doi: 10.1038/nature06244. PubMed PMID: 17943116; PubMed Central PMCID: PMCPMC3709439.

12.     Essilfie AT, Simpson JL, Horvat JC, Preston JA, Dunkley ML, Foster PS, et al. Haemophilus influenzae infection drives IL-17-mediated neutrophilic allergic airways disease. PLoS Pathog. 2011;7(10):e1002244. Epub 2011/10/15. doi: 10.1371/journal.ppat.1002244. PubMed PMID: 21998577; PubMed Central PMCID: PMCPMC3188527.

13.     Bousquet J, Chanez P, Lacoste JY, Barneon G, Ghavanian N, Enander I, et al. Eosinophilic inflammation in asthma. N Engl J Med. 1990;323(15):1033-9. Epub 1990/10/11. doi: 10.1056/NEJM199010113231505. PubMed PMID: 2215562.

14.     Ahren IL, Eriksson E, Egesten A, Riesbeck K. Nontypeable Haemophilus influenzae activates human eosinophils through beta-glucan receptors. Am J Respir Cell Mol Biol. 2003;29(5):598-605. Epub 2003/04/12. doi: 10.1165/rcmb.2002-0138OC. PubMed PMID: 12689921.

15.     Galli SJ, Tsai M. Mast cells in allergy and infection: versatile effector and regulatory cells in innate and adaptive immunity. Eur J Immunol. 2010;40(7):1843-51. Epub 2010/06/29. doi: 10.1002/eji.201040559. PubMed PMID: 20583030; PubMed Central PMCID: PMCPMC3581154.

16.     Masur H, Rosen PP, Armstrong D. Pulmonary disease caused by Candida species. Am J Med. 1977;63(6):914-25. Epub 1977/12/01. PubMed PMID: 343588.

17.     Mathieu E, Escribano-Vazquez U, Descamps D, Cherbuy C, Langella P, Riffault S, et al. Paradigms of Lung Microbiota Functions in Health and Disease, Particularly, in Asthma. Front Physiol. 2018;9:1168. Epub 2018/09/25. doi: 10.3389/fphys.2018.01168. PubMed PMID: 30246806; PubMed Central PMCID: PMCPMC6110890.

18.     Ascencio F, Ljungh A, Wadstrom T. Characterization of lactoferrin binding by Aeromonas hydrophila. Appl Environ Microbiol. 1992;58(1):42-7. Epub 1992/01/01. PubMed PMID: 1311545; PubMed Central PMCID: PMCPMC195170.

19.     Sajjan U, Keshavjee S, Forstner J. Responses of well-differentiated airway epithelial cell cultures from healthy donors and patients with cystic fibrosis to Burkholderia cenocepacia infection. Infect Immun. 2004;72(7):4188-99. Epub 2004/06/24. doi: 10.1128/IAI.72.7.4188-4199.2004. PubMed PMID: 15213163; PubMed Central PMCID: PMCPMC427436.

20.     Park CY, Heo JN, Suk K, Lee WH. Sodium azide suppresses LPS-induced expression MCP-1 through regulating IkappaBzeta and STAT1 activities in macrophages. Cell Immunol. 2017;315:64-70. Epub 2017/04/11. doi: 10.1016/j.cellimm.2017.02.007. PubMed PMID: 28391993.

21.     Hildebrand D, Bode KA, Riess D, Cerny D, Waldhuber A, Rommler F, et al. Granzyme A produces bioactive IL-1beta through a nonapoptotic inflammasome-independent pathway. Cell Rep. 2014;9(3):910-7. Epub 2014/12/02. doi: 10.1016/j.celrep.2014.10.003. PubMed PMID: 25437548.

22.     Chapman SJ, Khor CC, Vannberg FO, Rautanen A, Segal S, Moore CE, et al. NFKBIZ polymorphisms and susceptibility to pneumococcal disease in European and African populations. Genes Immun. 2010;11(4):319-25. Epub 2009/10/03. doi: 10.1038/gene.2009.76. PubMed PMID: 19798075; PubMed Central PMCID: PMCPMC3051152.

610    23.    Baldwin AS, Jr. The NF-kappa B and I kappa B proteins: new discoveries and insights. Annu Rev
611    Immunol. 1996;14:649-83. Epub 1996/01/01. doi: 10.1146/annurev.immunol.14.1.649. PubMed PMID:
612    8717528.
613    24.    Motoyama M, Yamazaki S, Eto-Kimura A, Takeshige K, Muta T. Positive and negative regulation of
614    nuclear factor-kappaB-mediated transcription by IkappaB-zeta, an inducible nuclear protein. J Biol Chem.
615    2005;280(9):7444-51. Epub 2004/12/25. doi: 10.1074/jbc.M412738200. PubMed PMID: 15618216.
616    25.    Yamazaki S, Muta T, Matsuo S, Takeshige K. Stimulus-specific induction of a novel nuclear factor-
617    kappaB regulator, IkappaB-zeta, via Toll/Interleukin-1 receptor is mediated by mRNA stabilization. J Biol
618    Chem. 2005;280(2):1678-87. Epub 2004/11/04. doi: 10.1074/jbc.M409983200. PubMed PMID: 15522867.
619    26.    Round JL, Mazmanian SK. The gut microbiota shapes intestinal immune responses during health and
620    disease. Nat Rev Immunol. 2009;9(5):313-23. Epub 2009/04/04. doi: 10.1038/nri2515. PubMed PMID:
621    19343057; PubMed Central PMCID: PMCPMC4095778.
622    27.    Shreiner AB, Kao JY, Young VB. The gut microbiome in health and in disease. Curr Opin
623    Gastroenterol. 2015;31(1):69-75. Epub 2014/11/14. doi: 10.1097/MOG.0000000000000139. PubMed PMID:
624    25394236; PubMed Central PMCID: PMCPMC4290017.
625    28.    Erb-Downward JR, Thompson DL, Han MK, Freeman CM, McCloskey L, Schmidt LA, et al. Analysis of
626    the lung microbiome in the "healthy" smoker and in COPD. PLoS One. 2011;6(2):e16384. Epub 2011/03/03.
627    doi: 10.1371/journal.pone.0016384. PubMed PMID: 21364979; PubMed Central PMCID: PMCPMC3043049.
628    29.    Morris A, Beck JM, Schloss PD, Campbell TB, Crothers K, Curtis JL, et al. Comparison of the
629    respiratory microbiome in healthy nonsmokers and smokers. Am J Respir Crit Care Med. 2013;187(10):1067-
630    75. Epub 2013/03/16. doi: 10.1164/rccm.201210-1913OC. PubMed PMID: 23491408; PubMed Central PMCID:
631    PMCPMC3734620.
632    30.    Hilty M, Burke C, Pedro H, Cardenas P, Bush A, Bossley C, et al. Disordered microbial communities in
633    asthmatic airways. PLoS One. 2010;5(1):e8578. Epub 2010/01/07. doi: 10.1371/journal.pone.0008578.
634    PubMed PMID: 20052417; PubMed Central PMCID: PMCPMC2798952.
635    31.    Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-
636    seq aligner. Bioinformatics. 2013;29(1):15-21. Epub 2012/10/30. doi: 10.1093/bioinformatics/bts635. PubMed
637    PMID: 23104886; PubMed Central PMCID: PMCPMC3530905.
638    32.    Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, et al. Ribosomal Database Project: data and
639    tools for high throughput rRNA analysis. Nucleic Acids Res. 2014;42(Database issue):D633-42. Epub
640    2013/11/30. doi: 10.1093/nar/gkt1244. PubMed PMID: 24288368; PubMed Central PMCID:
641    PMCPMC3965039.
642    33.    Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments.
643    Genome Biol. 2014;15(3):R46. Epub 2014/03/04. doi: 10.1186/gb-2014-15-3-r46. PubMed PMID: 24580807;
644    PubMed Central PMCID: PMCPMC4053813.
645    34.    Guo M, Wang H, Potter SS, Whitsett JA, Xu Y. SINCERA: A Pipeline for Single-Cell RNA-Seq Profiling
646    Analysis. PLoS Comput Biol. 2015;11(11):e1004575. Epub 2015/11/26. doi: 10.1371/journal.pcbi.1004575.
647    PubMed PMID: 26600239; PubMed Central PMCID: PMCPMC4658017.
648    35.    Newman MEJ, Girvan M. Finding and evaluating community structure in networks. Phys Rev E.
649    2004;69(2). PubMed PMID: WOS:000220255500019.
650    36.    Clauset A, Newman MEJ, Moore C. Finding community structure in very large networks. Phys Rev E.
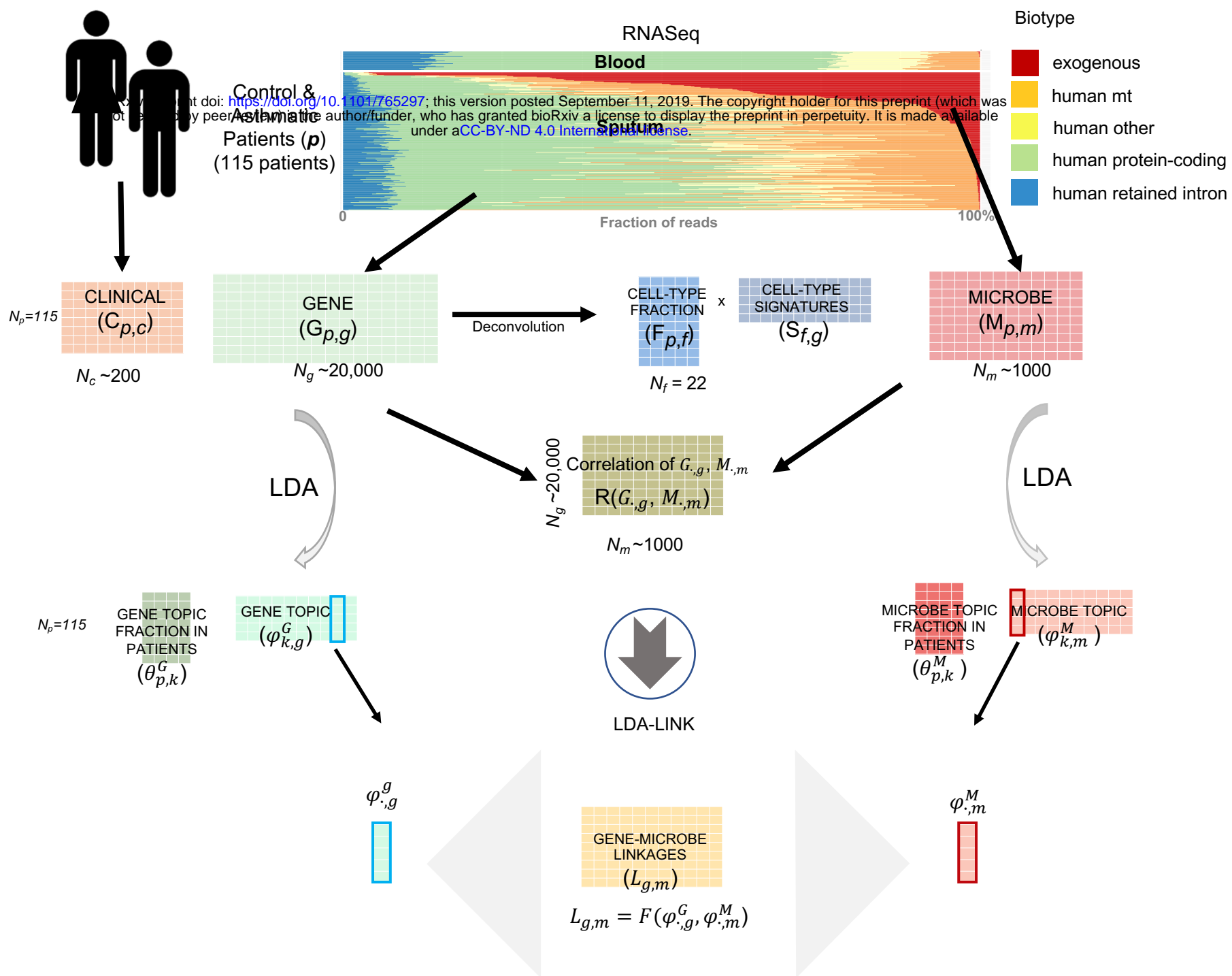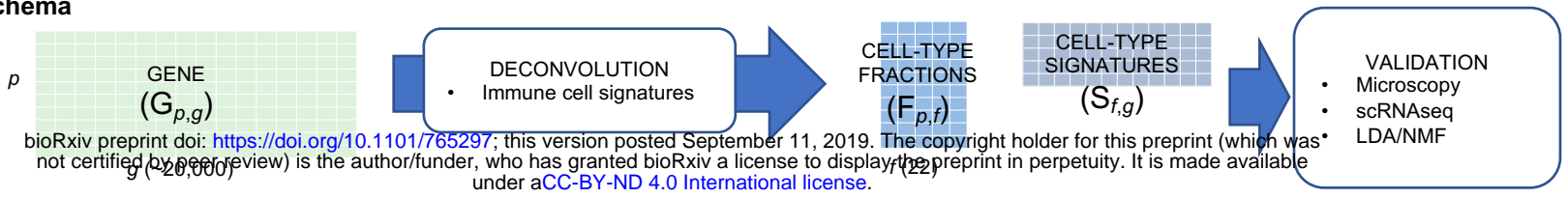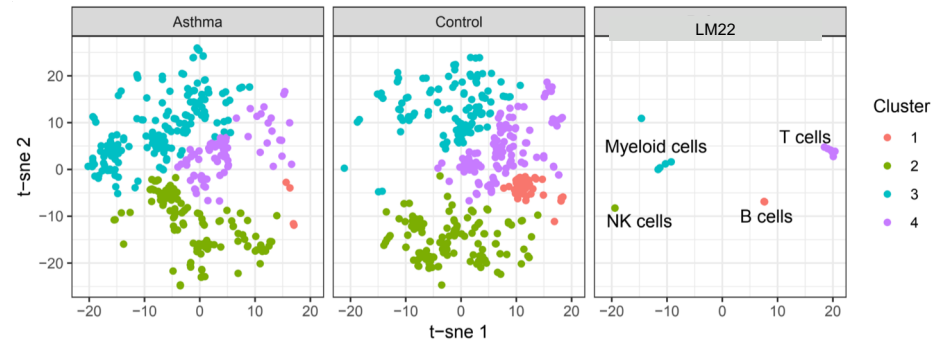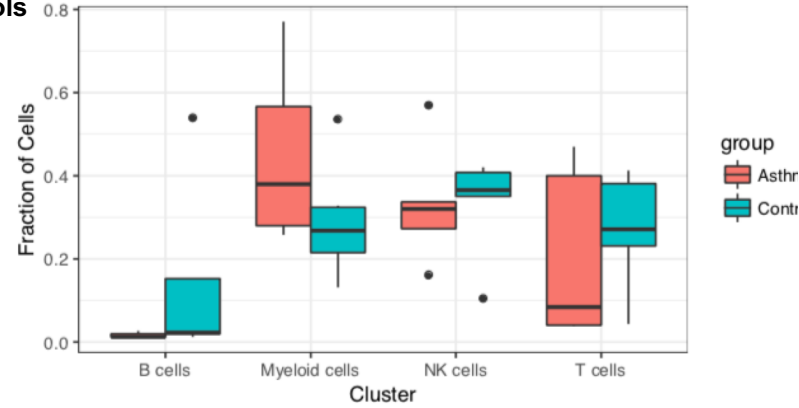651    2004;70(6). PubMed PMID: WOS:000226299200018.
652

Figure 1

**A. Schema**

GENE
$(G_{p,g})$
$p$
$g\ (\sim 20,000)$

DECONVOLUTION
• Immune cell signatures

CELL-TYPE FRACTIONS
$(F_{p,f})$
$f\ (22)$

CELL-TYPE SIGNATURES
$(S_{f,g})$

VALIDATION
• Microscopy
• scRNAseq
• LDA/NMF

**B Single cell RNAseq of sputum to identify cell types**

Asthma | Control | LM22

Myeloid cells, T cells, NK cells, B cells

t-sne 2 / t-sne 1

Cluster: 1, 2, 3, 4

**C. Differences between scRNAseq cell types in Asthmatics and controls**

Fraction of Cells

B cells, Myeloid cells, NK cells, T cells

Cluster

group: Asthma, Control

**D**

Microscopy

Eosinophils: r = 0.53, p = 1.8e−07
Lymphocytes: r = 0.25, p = 0.021
Macrophages: r = 0.27, p = 0.012
Neutrophils: r = 0.41, p = 8e−05

Cell-type fraction

**E. cell type hematopoiesis vs topic (components)**

$R\left(\ \begin{matrix}\text{CELL-TYPE}\\\text{FRACTION}\\(F_{\cdot,f})\end{matrix}\ ,\ \begin{matrix}\text{GENE TOPIC}\\\text{COMPONENTS}\\(\theta^G_{\cdot,t})\end{matrix}\ \right)$

Neutrophils, Eosinophils, Mast.cells.activated, Mast.cells.resting, Dendritic.cells.activated, Dendritic.cells.resting, Macrophages.M2, Macrophages.M1, Macrophages.M0, Monocytes, NK.cells.activated, NK.cells.resting, T.cells.gamma.delta, T.cells.regulatory..Tregs., T.cells.follicular.helper, T.cells.CD4.memory.activated, T.cells.CD4.memory.resting, T.cells.CD4.naive, T.cells.CD8, Plasma.cells, B.cells.memory, B.cells.naive

topic1, topic2, topic3, topic4, topic5, topic6, topic7, topic8, topic9, topic10

**F**

$R\left(\ \begin{matrix}\text{CELL-TYPE}\\\text{FRACTION}\\(F_{\cdot,f})\end{matrix}\ ,\ \begin{matrix}\text{CLINICAL}\\(C_{\cdot,c})\end{matrix}\ \right)$

Neutrophils, Eosinophils, Mast.cells.activated, Mast.cells.resting, Dendritic.cells.activated, Dendritic.cells.resting, Macrophages.M2, Macrophages.M1, Macrophages.M0, Monocytes, NK.cells.activated, NK.cells.resting, T.cells.gamma.delta, T.cells.regulatory..Tregs., T.cells.follicular.helper, T.cells.CD4.memory.activated, T.cells.CD4.memory.resting, T.cells.CD4.naive, T.cells.CD8, Plasma.cells, B.cells.memory, B.cells.naive

ACT, Age, Age.DX, Age.SX.Onset, BDR, BMI, FENO, FEV1.FVC.postBD, FEV1.FVC.preBD, HIL, HPY, ICS, Intubations, Number.of.OCS, OCS.Total, platelets, Total.Pack.Years, white.count
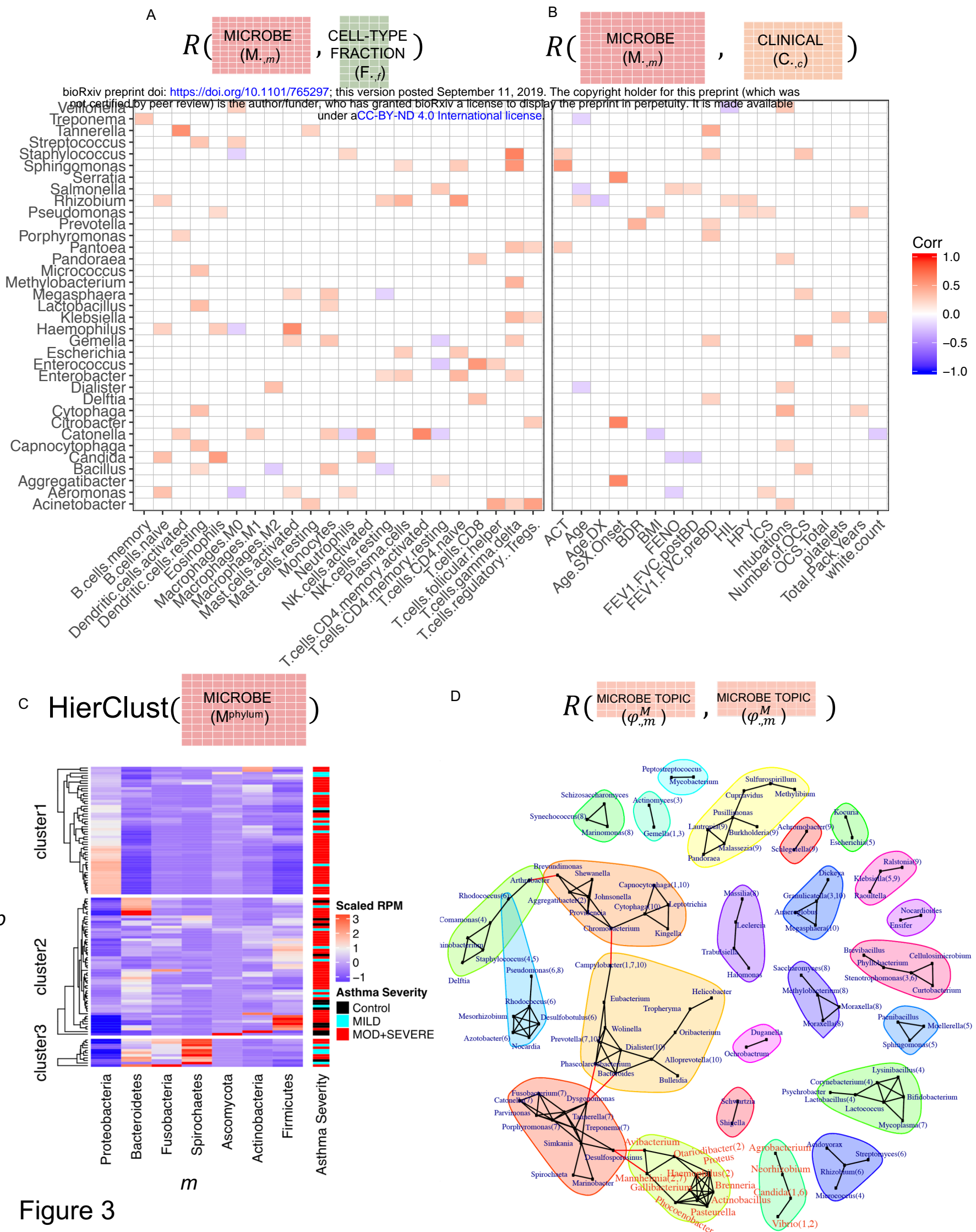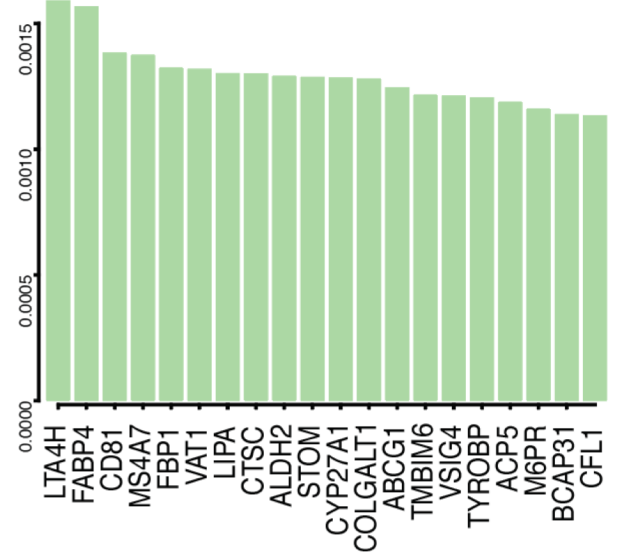
Corr: 1.0, 0.5, 0.0, −0.5, −1.0

Figure 2

Figure 3

A Microbe-gene linkage

Figure 4

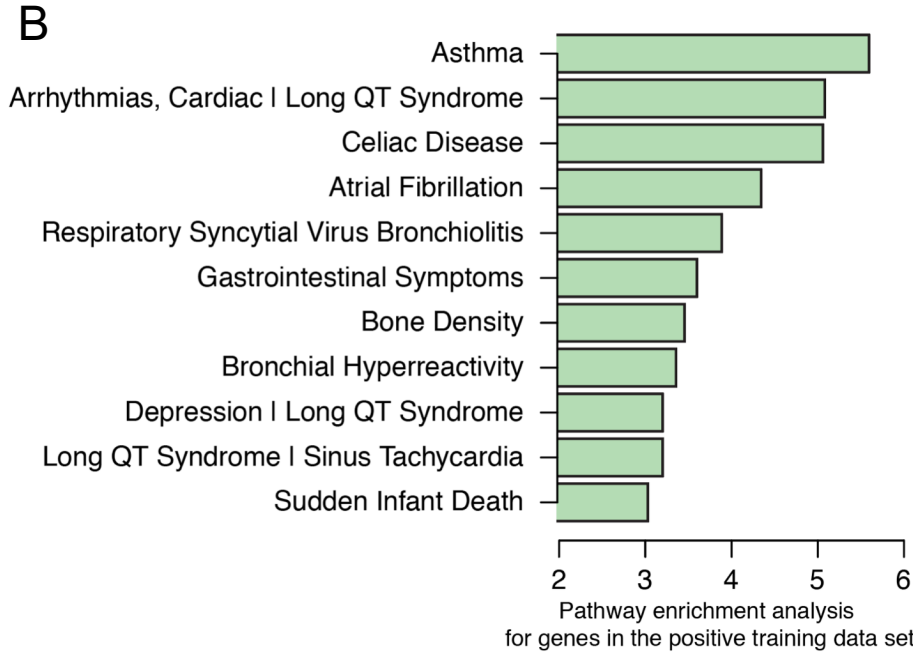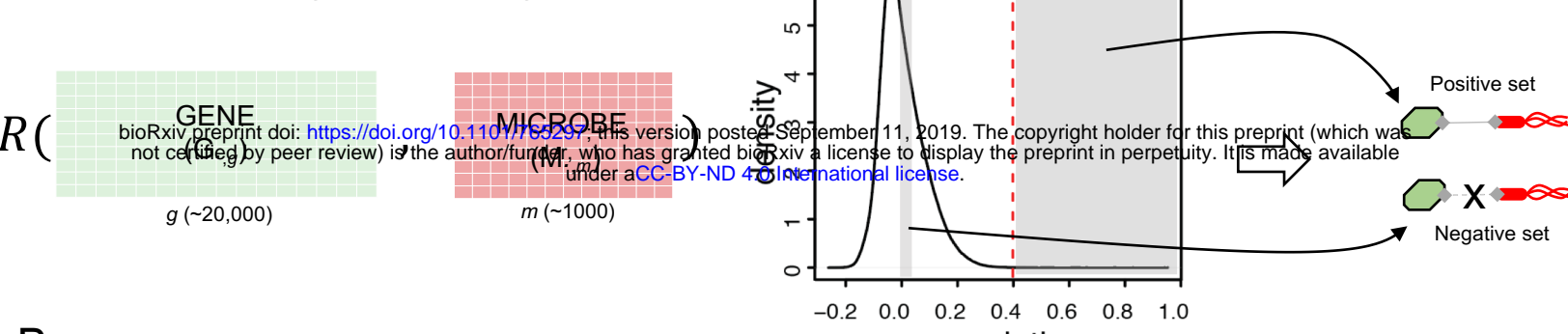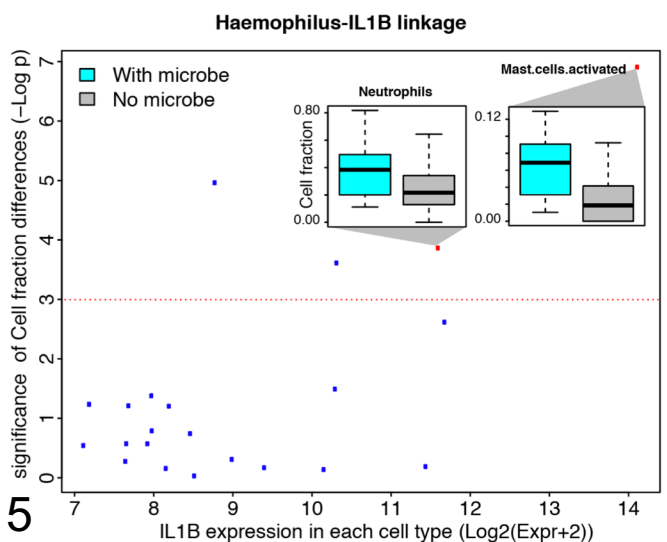# A  Subset of the gene-microbe linkages defined by the LDA-link model

Figure 5