

# 1 Morse-clustering of a Topological Data Analysis Network 2 Identifies Phenotypes of Asthma Based on Blood Gene 3 Expression Profiles

4 James P R Schofield<sup>1, 2</sup>, Fabio Strazzeri<sup>3, 4</sup>, Jeannette Bigler<sup>5</sup>, Michael Boedigheimer<sup>5</sup>, Ian M  
5 Adcock<sup>6</sup>, Kian Fan Chung<sup>6</sup>, Aruna Bansal<sup>7</sup>, Richard Knowles<sup>8</sup>, Sven-Erik Dahlen<sup>9</sup>, Craig E.  
6 Wheelock<sup>10</sup>, Kai Sun<sup>6</sup>, Ioannis Pandis<sup>11</sup>, John Riley<sup>12</sup>, Charles Auffray<sup>13</sup>, Bertrand De  
7 Meulder<sup>13</sup>, Diane Lefaudeux<sup>13</sup>, Devi Ramanan<sup>14</sup>, Ana R Sousa<sup>12</sup>, Peter J Sterk<sup>15</sup>, Rob. M  
8 Ewing<sup>3</sup>, Ben D Macarthur<sup>4</sup>, Ratko Djukanovic<sup>2</sup>, Ruben Sanchez-Garcia<sup>4</sup> and Paul J Skipp<sup>1</sup>

9 <sup>1</sup>Centre for Proteomic Research, Institute for Life Sciences, University of Southampton,  
10 Southampton, UK

11 <sup>2</sup>NIHR Southampton Respiratory Biomedical Research Unit and Clinical and Experimental  
12 Sciences, Southampton, UK

13 <sup>3</sup>Mathematical Sciences, University of Southampton, Southampton, UK

14 <sup>4</sup>Biological Sciences, University of Southampton, Southampton, UK

15 <sup>5</sup>Amgen, Thousand Oaks, CA, USA

16 <sup>6</sup>National Heart and Lung Institute, Imperial College, London, UK

17 <sup>7</sup>Acclarogen, St John's Innovation Centre, Cambridge, UK

18 <sup>8</sup>Arachos Pharma, Stevenage, UK

19 <sup>9</sup>Karolinska University Hospital and Centre for Allergy Research, Karolinska Institutet,  
20 Stockholm, Sweden

21 <sup>10</sup>Centre for Allergy Research, Karolinska Institutet, Stockholm, Sweden

22 <sup>11</sup>Data Science Institute, Imperial College, London, UK

23 <sup>12</sup>Respiratory Therapeutic Unit, GSK, Stevenage, UK

24 <sup>13</sup>European Institute for Systems Biology and Medicine, CNRS-ENS-UCBL-INSERM, CIRI-  
25 UMR5308, Lyon, France

26 <sup>14</sup>Ayasdi

27 <sup>15</sup>Academic Medical Centre, University of Amsterdam, Amsterdam, The Netherlands

28 To whom correspondence should be addressed: James Schofield at  
29 J.P.Schofield@soton.ac.uk, NIHR Southampton Biomedical Research Centre,  
30 Clinical and Experimental Sciences, Faculty of Medicine, University of  
31 Southampton, UK

32 Sources of funding

33 The U-BIOPRED consortium receives funding from the European Union and from  
34 the European Federation of Pharmaceutical Industries and Associations as an IMI  
35 JU funded project (no. 115010).

36 Conflicts of interest: James Schofield has nothing to disclose. Fabio Strazzeri has  
37 nothing to disclose; Ratko Djukanović has consulted and presented at symposia  
38 organised by TEVA, Novartis, GlaxoSmithKline and AstraZeneca, has shares in  
39 and consults for Synairgen; Charles Auffray reports grants from Innovative  
40 Medicine Initiative; Jeanette Bigler reports that she owns stock in Amgen Inc;  
41 Michael Boedigheimer owns stocks in Amgen Inc; Kian Fan Chung has received  
42 honoraria for participating in Advisory Board meetings of the pharmaceutical  
43 industry regarding treatments for asthma and chronic obstructive pulmonary  
44 disease and has also been remunerated for speaking engagements; Kai Sun has  
45 nothing to disclose; Ioannis Pandis has nothing to disclose. Peter Sterk reports  
46 grants from IMI: Innovative Medicines Initiative, during the conduct of the study;  
47 Aruna T Bansal has nothing to disclose; Ian Adcock has nothing to disclose. Ben  
48 Macarthur has nothing to disclose; Ruben Sanchez-Garcia has nothing to disclose;  
49 Paul Skipp has nothing to disclose.

50

51

52

## 53 Abstract

54 Stratified medicine requires discretisation of disease populations for targeted treatments. We  
55 have developed and applied a discrete Morse theory clustering algorithm to a Topological Data  
56 Analysis (TDA) network model of 498 gene expression profiles of peripheral blood from  
57 asthma and healthy participants. The Morse clustering algorithm defined nine clusters, BC1-9,  
58 representing molecular phenotypes with discrete phenotypes including Type-1, 2 & 17  
59 cytokine inflammatory pathways. The TDA network model and clusters were also  
60 characterised by activity of glucocorticoid receptor signalling associated with different  
61 expression profiles of glucocorticoid receptor (GR), according to microarray probesets targeted  
62 to the start or end of the GR mRNA's 3' UTR; suggesting differential GR mRNA processing  
63 as a possible driver of asthma phenotypes including steroid insensitivity.

64 **Key words:** asthma, topological data analysis, discrete Morse theory, inflammation, cytokines

65

## 66 Introduction

67 Asthma is ranked 16<sup>th</sup> among the leading causes of years lived with disability and affects 339  
68 million people worldwide. Asthma is characterised by an expiratory airflow limitation,  
69 typically reported as forced expiratory volume in one second (FEV<sub>1</sub>). Treated is primarily with  
70  $\beta$ 2-agonists which relax airway smooth muscle, and corticosteroids which reduce underlying  
71 inflammation. Drugs have also been developed to target specific inflammatory pathways such  
72 as the T2 biologics, which reduce asthma exacerbation frequency by around 50%<sup>1,2</sup>. Improved  
73 understanding of asthma disease progression and molecular sub-phenotypes should improve  
74 the use and development of new targeted therapeutics. In this study, we used data from the U-  
75 BIOPRED (Unbiased BIOMarkers for the Prediction of respiratory disease outcomes) project,  
76 the largest multi-centre asthma programme to date, involving 20 academic institutions, 11  
77 pharmaceutical companies and patient groups and charities, with the aim to improve  
78 understanding of the complex molecular mechanisms underpinning asthma and identify useful  
79 biomarkers<sup>3-10</sup>.

80 Asthma is characterized by variability in symptoms and treatment response. Around half of  
81 asthma is thought to arise from T-2 immunity, driven by IL4, IL5 and IL13 cytokine associated  
82 with recruitment of eosinophils into airways<sup>11</sup>. Additionally, high sputum neutrophil counts  
83 are associated with reduced post-bronchodilator FEV<sub>1</sub><sup>12</sup>. Corticosteroids are routinely used to  
84 reduce airway inflammation in asthma by activating glucocorticoid receptor (GR) and  
85 suppressing NF- $\kappa$ B activity which regulates expression of pro-inflammatory cytokines and  
86 cyclo-oxygenase 2 (COX2) as well as inducible nitric oxide synthase (iNOS). However,  
87 patients with severe asthma, particularly T-2-low and T-17-high asthma<sup>13</sup>, respond poorly to  
88 corticosteroids, but it is not known why. The relative expression of GR- $\alpha$  and GR- $\beta$  protein  
89 isoforms, resulting from alternative splicing, influences steroid insensitivity, as GR- $\beta$  does not  
90 bind GC and inhibits GR- $\alpha$  activity by forming a heterodimer<sup>14</sup>. GR protein expression is  
91 further regulated by ARE-mediated degradation of GR mRNA targeting the AU-rich elements  
92 within the 3' UTR<sup>15</sup>.

93 Topological Data Analysis (TDA) is an unsupervised machine learning tool suitable for  
94 analysis of high-dimensional datasets<sup>16,17,18</sup>. Application of TDA via the Mapper algorithm

95 generates a TDA network model, a compressed representation of high-dimensional data with  
96 major features embedded where similar data points are grouped into nodes, and nodes with  
97 common data points are connected by edges. We have previously reported an analysis of  
98 differentially expressed genes (DEGs) from gene expression profiling of 498 gene expression  
99 profiles of peripheral blood from participants in the U-BIOPRED (Unbiased Biomarkers in  
100 Prediction of Respiratory Disease Outcomes) study<sup>10</sup>. Unbiased hierarchical clustering of  
101 DEGs identified two sub-groups, one enriched for patients with severe asthma, use of oral  
102 corticosteroids and blood neutrophilia, and a second cluster composed of mixed-severity  
103 asthmatics and healthy individuals. We generated a Topological Data Analysis (TDA) network  
104 model of the same gene expression data using the Ayasdi TDA software platform and found  
105 these two clusters represented by different regions of the TDA network model. In this study,  
106 we investigated the continuous variation of clinical and molecular biology in the TDA network  
107 model representing the shape of asthma disease pathology; shedding light on possible routes  
108 of disease progression.

109 Stratification of disease allows targeted treatment for improved patient outcome, so we  
110 developed and applied a Morse-clustering algorithm to discretise the continuous TDA network  
111 model of patients into clusters representing different molecular phenotypes of asthma sub-  
112 types. Clusters within TDA networks have typically been delineated by eye<sup>18,19,20</sup>, without  
113 algorithmic reproducibility and few studies have used the standard network clustering  
114 algorithm, community clustering, via the Ayasdi Python SDK. The community clustering  
115 algorithm is limited as it only analyses connectivity between nodes without considering the  
116 density of data points clustered within nodes, an important dimension in TDA network models.  
117 This 3<sup>rd</sup> dimension in the TDA network can be visualised by colouring (Fig. 3A & B) and the  
118 TDA network can, therefore, be considered as a connected 3D map of data points clustered  
119 around peaks that represent conserved sub-types or phenotypes of major features, which in the  
120 study of patient gene expression reflect biological pathway modulations underlying disease  
121 phenotypes. Discrete Morse theory relates the flow (gradients) on a discrete object, such as a  
122 network, with its topology<sup>21</sup>. Here we apply Morse theory to measure the gradients and  
123 connected peaks within a TDA network, thus delineating clusters according to key features of  
124 the dataset. We have developed a Python script to apply Morse-based clustering of TDA  
125 networks in the open source Mapper TDA software or through the Ayasdi Python software  
126 development kit (SDK) which we believe will add value to future analyses. This Morse-  
127 clustering algorithm identified nine clusters, BC1-9, representing discrete molecular  
128 phenotypes characterised by differences in circulating immune cell populations, activation of  
129 T-1, -2 & -17 cytokine inflammatory pathways, and the activity of glucocorticoid receptor  
130 signalling and novel differences in glucocorticoid receptor mRNA isoforms.

131

132

133

134

## 135 **Results**

136 The TDA network model of peripheral blood gene expression from 498 participants in the U-  
137 BIOPRED asthma study consisted of a hub with an increased prevalence of healthy participants  
138 and connected flares with increased prevalence of severe asthma and decreased FEV<sub>1</sub>,  
139 reflecting multiple interconnected possible routes of disease progression (Fig. 1). Regions of  
140 the TDA network with highest eosinophil counts (Fig. 1G) had high prevalence of severe  
141 asthma (Fig. 1E) and were associated with high COX2, NF- $\kappa$ B, IL5, IL13 (Fig. 1J, N, O, P),  
142 and low IFN- $\gamma$  and GR mRNA (Fig. 1T, Q, R). There was a distinct pattern across the TDA  
143 network model of GR mRNA expression according to probesets targeting the start of the 3'  
144 UTR (probesets 201865\_x\_at and 211671\_s\_at, illustrated as  $\Delta$ x NR3C1 mRNA in Fig. 1R)  
145 and a different pattern according to probesets targeting towards the end of the 3' UTR  
146 (probesets 201866\_s\_at and 216321\_s\_at, illustrated as FL NR3C1 mRNA in Fig. 1Q). The  
147 binding locations of the Affymetrix NR3C1 probes and corresponding NCBI RefSeq sequences  
148 are shown mapped onto the Human genome in figure 2. We hypothesized that the  $\Delta$ x NR3C1  
149 mRNA has a truncated 3' UTR compared to the FL NR3C1; meaning  $\Delta$ x NR3C1 has fewer  
150 AU-rich elements (AREs), and is missing a miR 486 target sequence, compared to the FL  
151 NR3C1 mRNA. The TDA network was polarised by FL NR3C1 (Fig. 1Q) and associated GR-  
152 responsive genes, COX2, ANXA1 and IFN $\gamma$  (Fig. 1J, L, T). Probesets targeting the start of the  
153 3' UTR of GR mRNA indicated a different pattern of expression across the TDA model (Fig.  
154 1R) and corresponded to OCS dose (Fig. 1I) and GR-responsive gene expression, ZPF36,  
155 GILZ, FKBP5 (Fig. 1K, M, S).

156 To define groups of people with similar gene expression signatures from the TDA network  
157 model, we developed and applied a Morse-clustering algorithm. The Morse-clustering  
158 algorithm identified 9 clusters which we termed BC1 to 9. The reporter operating characteristic  
159 (ROC) area under the curve (AUC) for the 9 clusters ranged from 0.76 to 0.97, representing  
160 very good to excellent prediction of cluster classification in the test set based on a logistic  
161 regression model identifying predictors of the cluster in the training set (Fig. 4). BC1-9 were  
162 found to have activation of cytokine-mediated inflammatory pathways consistent with their  
163 distribution on the TDA network model with trends identified in pathway and upstream  
164 regulator activation across the clusters (Table 1 & 2). BC1 was predominantly severe  
165 asthmatics, with reduced lung function, represented by low FEV<sub>1</sub>. BC1 also had a T-17  
166 signature of gene expression<sup>22</sup>, with increased expression of IL17A, IL21 and IL22 ( $q = 1.31E^{-5}$ ,  
167  $7.99E^{-4}$ ,  $1.71E^{-3}$ ). BC1 had decreased expression of  $\beta$ -2 adrenergic receptor (ADRB2) mRNA  
168 the protein product of which is involved in smooth muscle relaxation and bronchodilatation.  
169 Cystatin D (CST5) was predicted as the most activated upstream regulator of gene expression  
170 in BC1 but was also highly activated in BC9 and 8 (Table 2).

## 171 Discussion

172 The TDA network model identified familiar phenotypes of asthma and gave insight into  
173 potential routes of disease progression. For example, the furthest eosinophilic region from the  
174 'healthy hub' was associated with high T-17 markers, TGF $\beta$ , IL17A, IL21, IL22 (Fig. 1D, V,  
175 W, X) and increased neutrophilia (Fig. 1H). The T-17 region was connected to the 'healthy  
176 hub' via the solely T-2 high region, suggesting disease progression from healthy to T-17 high  
177 via an only T-2-high phenotype. Differential expression of FL NR3C1 and  $\Delta$ x NR3C1 and  
178 corresponding expression patterns of GR-responsive genes suggests different functional  
179 responses to steroids across the TDA network model, associated with differential expression  
180 of GR mRNA isoforms.

181 The Morse-clustering algorithm identified 9 clusters, however, clusters BC4, 6 and 8 were  
182 small ( $n=35, 37, 33$ , respectively), with correspondingly low representation in the training and  
183 test sets which resulted in ROC curves whose shapes were not smooth and may have  
184 represented overfitting. The identified clusters represented groups of patients with significant  
185 differences in the activation of pathways related to inflammation, including pathways  
186 associated with glucocorticoid receptor (GR) signalling, Type (T)-2, T-1 and T-17  
187 inflammatory responses. Transglutaminase (TGM2), a marker of T-2 inflammation<sup>23</sup>, was  
188 predicted in this study as the most activated upstream regulator of gene expression in BC2, 3,  
189 7 and 8 (Table 2). It is known to catalyse the serotonin transamidation of glutamines  
190 (serotonylation), which regulates cell signalling and actin polymerization. BC2 and 3 were  
191 characterised by high TGM2-mediated gene expression, including Toll-like receptors (TLR)  
192 and iNOS signalling. TGM2 is also implicated in recruitment of eosinophils into asthmatic  
193 airways<sup>11</sup>, which was reflected in the highest sputum eosinophil count in BC2, but high sputum  
194 eosinophils counts were not seen in BC3 (Table 3). Melatonin, the end product of the serotonin  
195 pathway is a free radical scavenger, acting to suppress inflammation<sup>24</sup>. Pathways associated  
196 with tryptophan metabolism were enriched in cluster BC1; serotonin degradation was the most  
197 activated pathway identified by IPA (Table 1). Serotonin levels are known to be implicated in  
198 asthma pathology, and serum serotonin levels tend to be increased in patients with active  
199 asthma<sup>25</sup>. The increased activation of melatonin degradation in BC1 may contribute to the  
200 severe asthma phenotype.

201 T-cell acute lymphocytic leukemia protein 1 (TAL1) was identified as the top upstream  
202 regulator of gene expression in BC9, together with miR-486, which has previously been  
203 identified as a potential marker of childhood asthma in plasma<sup>26</sup> and a promoter of NF- $\kappa$ B  
204 activity<sup>27</sup>. Our analysis predicted CD24 as the most activated upstream regulator of gene  
205 expression in BC6, 4, and 5. CD24 can reflect activity of one of its key transcription factors,  
206 c-myc, whose expression is inhibited by CST5. BC5 had high expression of IFN- $\gamma$  mRNA (Fig.  
207 1T), indicative of a T-1 response; however, IFN- $\gamma$ -mediated gene expression was not  
208 upregulated in this group (Table 3).

209 The shape of the TDA network and patterns of gene expression representative of differentially  
210 activated pathways reflected both corticosteroids use and expression of GR mRNA.  
211 Clusters BC1-3, mostly representing those of the Severe Asthma enriched cluster previously  
212 reported<sup>10</sup> (Fig. 1C), had the highest percentages of patients on OCS (Table 3). These clusters  
213 were also characterised by enrichment for patients on high doses of OCS, but other clusters  
214 were also enriched for patients with high OCS dose; particularly cluster BC5 (Fig. 1I). We  
215 observed common patterns of gene expression under the control of glucocorticoid response  
216 elements (GRE) that were differentially expressed between clusters, although the patterns were  
217 not necessarily consistent between GRE genes. This suggests different types of steroid response  
218 between the clusters. We did not find GR-signalling as a top upstream regulator of gene  
219 expression using IPA, because there are two signatures of GR-signalling which are alternately  
220 up and down regulated in the TDA structure. The expression of GRE genes, glucocorticoid-  
221 induced leucine zipper (GILZ), FK506-binding protein 5 (FKBP5) and Tristetraprolin (ZFP36)  
222 (Fig. 1M, S and K) were similarly distributed across Morse-clusters high in neutrophilic  
223 clusters of the top of the TDA network, BC1, 2, 3 & 4 and higher in the predominantly healthy  
224 cluster, BC7. However, the expression of Annexin A1, a classical indicator of steroid response,  
225 was very differently distributed between clusters (Fig. 1L) and was significantly higher in BC5  
226 when compared to the other patients ( $q = 2.3E^{-10}$ ). Serotonin degradation, which is  
227 interdependent on GR signalling, was identified as the top canonical pathway enriched in BC1  
228 (Table 1). In clusters BC1-3, there was increased expression of the RNA-binding protein,

229 tristetraprolin (TTP), a negative regulator of mRNA half-life, binding to AREs in the 3' UTR  
230 of target genes (Fig. 1K). Since the expression of TTP is regulated by a GRE site, GR-signalling  
231 causes increased ARE-mediated mRNA decay.

232 BC1 had low expression of short ( $\Delta x$  NR3C1) and long (FL NR3C1) GR mRNA and low  
233 expression of steroid-inducible anti-inflammatory mRNAs ANXA1 (Fig. 1L), SOCS1 and high  
234 expression of pro-inflammatory COX genes (Fig. 1J). We detected mixed levels of GILZ and  
235 FKBP5 (Fig. 1M & S). There was moderate expression of DUSP1 mRNA, another marker of  
236 GR activity. In the clusters on the left side of the TDA network there was high expression of  
237 NUPR1 which increases expression of p38MAPK, a key regulator of asthma pathogenesis<sup>28</sup>.  
238 Additionally, NUPR1 is known to activate phosphatidylinositol 3-kinases (PI3K)<sup>29</sup> which  
239 activate phosphoinositide pathways; inositol-related metabolism was highly upregulated in  
240 BC5 and 6, where the expression of phosphoinositol (PI) phosphatases was increased relative  
241 to health. Conversely, the expression of PI phosphatases was decreased when compared to  
242 health in BC8 and 9. Clusters BC5 and 6 showed increased expression of the enzyme which  
243 catalyses the dephosphorylation of 1D-myo-inositol (3)-monophosphate to myo-inositol,  
244 inositol-1 (or 4)-monophosphatase, when compared to health, whereas BC1, 7, 8 and 9 had  
245 decreased expression relative to health. It has previously been reported that myo-inositol is  
246 increased in animal asthma models following steroid treatment<sup>30</sup>, suggesting differential  
247 steroid responses between these clusters. In contrast to BC1, BC5 and 6 had gene expression  
248 profiles characteristic of low GR responses, as indicated by activation of CD24-mediated gene  
249 expression and inactivation of CST5-mediated gene expression. CST5 is activated by vitamin  
250 D receptor (VDR) expression<sup>31</sup>, whose expression is regulated by steroid-induced  
251 GR signalling<sup>32</sup> (Fig. 5). The enriched expression of inositol pathways in BC5 and 6 provided  
252 further support of a low GR response. Contraction of airway smooth muscle is initiated by  
253 increased cytosolic calcium ions ( $Ca^{2+}$ ), so this may, in part, explain the reduced FEV<sub>1</sub> seen in  
254 these clusters.

255 We propose that Morse clustering can be applied to TDA networks of patient 'omics data to  
256 identify sub-phenotypes of disease, thereby offering new insights into disease mechanisms and  
257 stratification of patients for more targeted drug development based on molecular mechanisms.

258

## 259 **Materials and Methods**

### 260 **Study population**

261 U-BIOPRED is a multi-centre prospective cohort study, involving 16 clinical centres in 11  
262 European countries. Blood samples were analysed from 498 study participants; 246 non-  
263 smoking severe asthmatics, 88 smoking severe asthmatics, 77 non-smoking mild/moderate  
264 asthmatics and 87 non-smoking non-asthmatic individuals. It is registered on  
265 ClinicalTrials.gov (identifier: NCT01982162).

### 266 **Ethics Statement**

267 The study was conducted in accordance with the principles expressed in the Declaration of  
268 Helsinki. It was approved by the Institutional Review Boards of all the participating

269 institutions; Academic Medical Centre (AMC), Amsterdam; University Hospital Southampton  
270 NHS Trust; South Manchester Healthcare Trust; Protisvalor Méditerranée SAS; Karolinska  
271 University Hospital; Nottingham University Hospital; NIHR-Wellcome Trust Clinical  
272 Research Facility; and adhered to the standards set by the International Conference on  
273 Harmonization and Good Clinical Practice. All participants provided written informed consent.

## 274 **Microarray Analysis**

275 RNA was isolated using the PAXgene Blood RNA kit (Qiagen, Valencia, CA) with on-column  
276 DNase treatment (Qiagen). RNA integrity was assessed using a Bioanalyzer 2100 (Agilent  
277 Technologies, Santa Clara, CA). Samples with  $RIN \geq 6$  were processed for microarray as  
278 described (19) and hybridized onto Affymetrix HT HG-U133PM+ arrays (Affymetrix, Santa  
279 Clara, CA) using a GeneTitanR according to Affymetrix technical protocols. The microarray  
280 data are deposited in GEO under GSE69683.

## 281 **Training and Test Data Analysis Sets**

282 The 498 samples available for analysis were randomized into training ( $n = 328$ ) and validation  
283 sets ( $n = 170$ ).

## 284 **Topological Data Analysis**

### 285 **Generating TDA graphs in Ayasdi Platform**

286 The transcriptomics data were clustered by topological data analysis (TDA) as previously  
287 reported<sup>10</sup>, using Ayasdi Platform with a norm correlation metric and two Neighbourhood  
288 lenses. Correlation was measured using normalised values for the expression of each probeset  
289 (Metric: norm correlation). The space for clustering was generated using 100 bins in each  
290 dimension according to t-SNE -calculated vectors and 60% overlap between neighbouring bins  
291 (Fig 3A): two neighbourhood lenses, resolution = 100; gain,  $\times 6$ ).

### 292 **Clustering of high patient density regions of TDA graphs**

293 Using the Ayasdi TDA Platform, the magnitude of nodes was represented by a colour heatmap  
294 where the colour spectrum from blue to red represent the range from the lowest to highest  
295 levels. Discrete Morse theory was applied to cluster TDA nodes according to patient density.  
296 Data from each node's neighbours were also used in calculating the annotation function, giving  
297 context to where a node lies within the broader topology, effectively 'smoothing' the data,  
298 decreasing noise and allowing identification of the most prominent peaks. To each node we  
299 assigned the annotation  $f: V \rightarrow \mathbb{R}^2$  where for each node  $C_i$  we have

$$300 \quad f(C_i) = \left( s(C_i), \left( s(C_i) + \sum s(C_j) \right) * Corr(C_i) \right),$$



301 and  $Corr(C_i)$  is the average correlation among all the patient in cluster-node  $C_i$ . Differently  
302 from other clustering algorithms, as k-nearest neighbours, we do not assume that cluster-nodes  
303 with similar value with respect to  $f$  are similar, neither we expect that  $f$  is a kernel-based  
304 function which fits the data. Our approach instead assumes that  $f$  gives the cluster-nodes a  
305 hierarchical structure and the nodes' connectivity is supplied by the Mapper network. In this  
306 way, with Morse, each cluster of nodes in the network has a structure of rooted tree and each  
307 leaf connects a cluster-node to a higher one (with respect to  $f$ ) with the root the highest cluster-  
308 node.

### 309 **Robustness of TDA network clusters evaluated by ROC analysis**

310 We applied logistic regression to test the tightness of the clusters according to key features  
311 identified by logistic regression. A logistic regression model was trained on a pre-defined  
312 training set of ( $n = 328$ ) and the classification accuracy tested on a test data set ( $n = 170$ ).  
313 Accuracy of the logistic regression reflects reproducibility in the clustering, ie. robust  
314 classification assigned by clustering results in accurate classification of test data by an  
315 independently trained logistic regression model.

316 Affymetrix probes for NR3C1 were aligned with NCBI RefSeq genes using the Ensembl  
317 Genome browser 94.

### 318 **Pathway analysis identified trends and discrete molecular features of clusters**

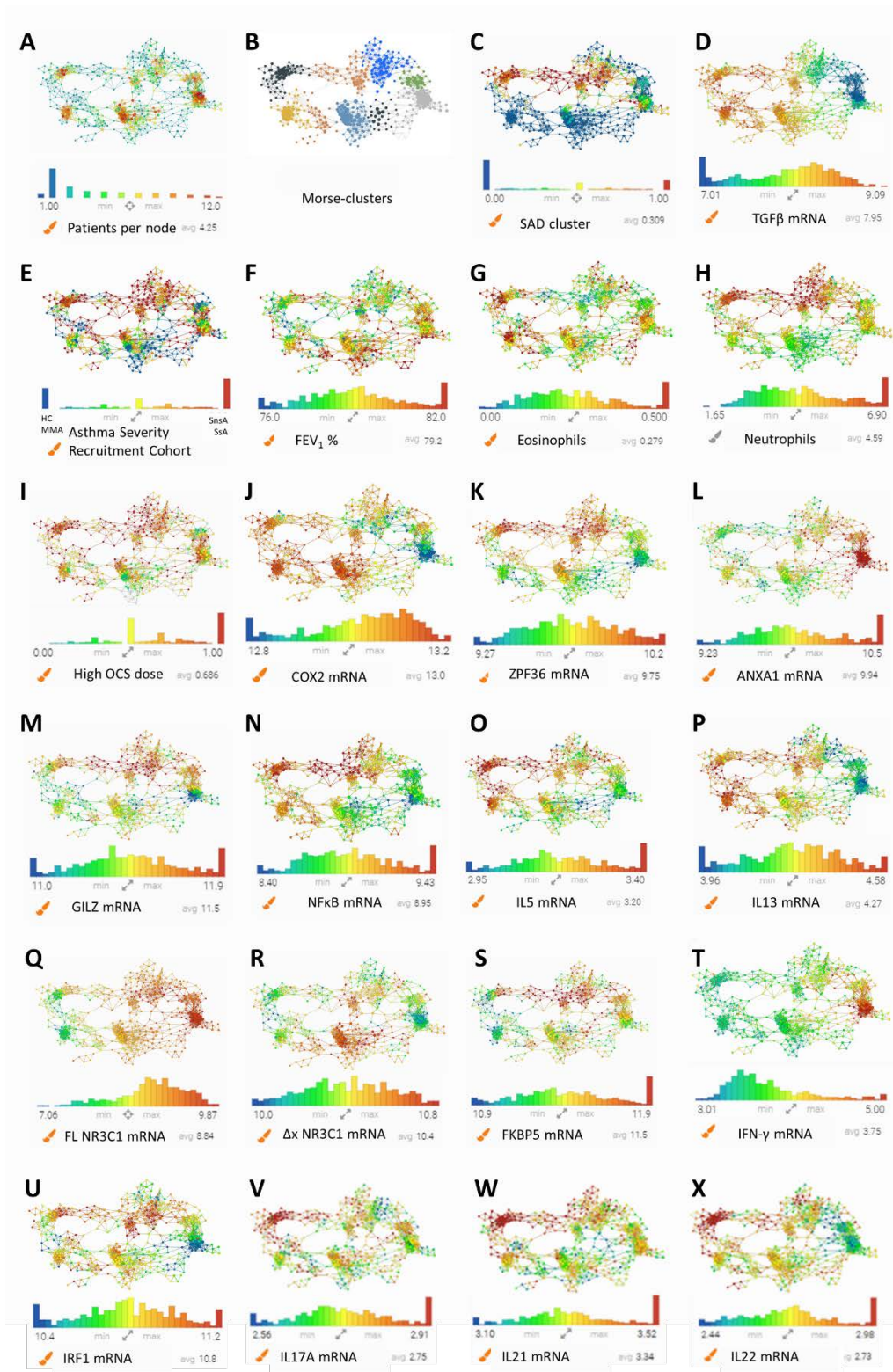
319 The shape of data represented by a TDA network is defined by the lenses (t-SNE in this study),  
320 which are implicitly used as coordinates for plotting the network. These coordinates focus on  
321 differentially activated pathways because genes of a common pathway are more likely to be  
322 co-expressed, and patients are clustered by similarity in key features in a TDA network.  
323 Ingenuity pathway analysis (IPA) was used to identify pathways with enriched gene expression  
324 within each of the clusters (Table 1), many of which were activated in clusters neighbouring  
325 each other in the TDA network, reflecting a trend in the activation of key pathways across the  
326 TDA network.

327

328

329 **Figure 1** Selected gene expression distribution across the TDA network

330 **Figure 1.** Selected gene expression distribution across the TDA network. Colours in legends denote the  
331 concentrations of the gene expression, ranging from blue (low) to red (high).



332

333 **Figure 2:** The chromosome binding locations of the Affymetrix NR3C1 probes



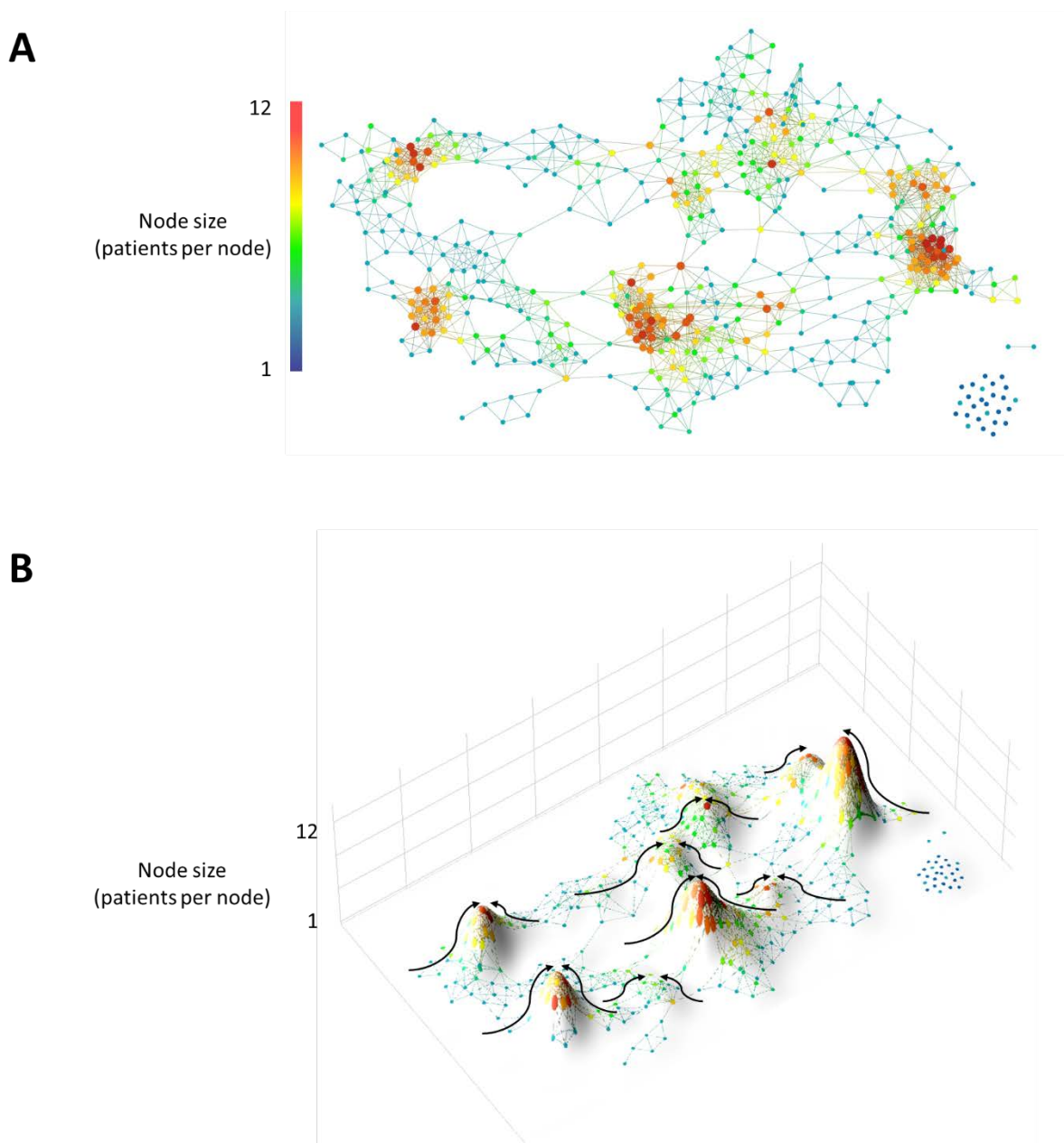
334 **Figure 2.** The binding locations of the Affymetrix NR3C1 probes and corresponding NCBI RefSeq sequences aligned to the  
335 Human genome. NR3C1 probesets 201865\_x\_at and 211671\_s\_at target isoforms with truncated 3' UTR:  $\Delta x$  NR3C1.  
336 Probesets 201866\_s\_at and 216321\_s\_at target NR3C1 mRNAs towards the end of the 3' UTR annotated in the RefSeq genes.  
337 Image generated using the Ensembl Genome Browser: <https://genome.ucsc.edu>  
338

339

340

341

342 **Figure 3:** Morse-clustering of the TDA network of UBIOPRED gene expression profiling of  
343 peripheral blood



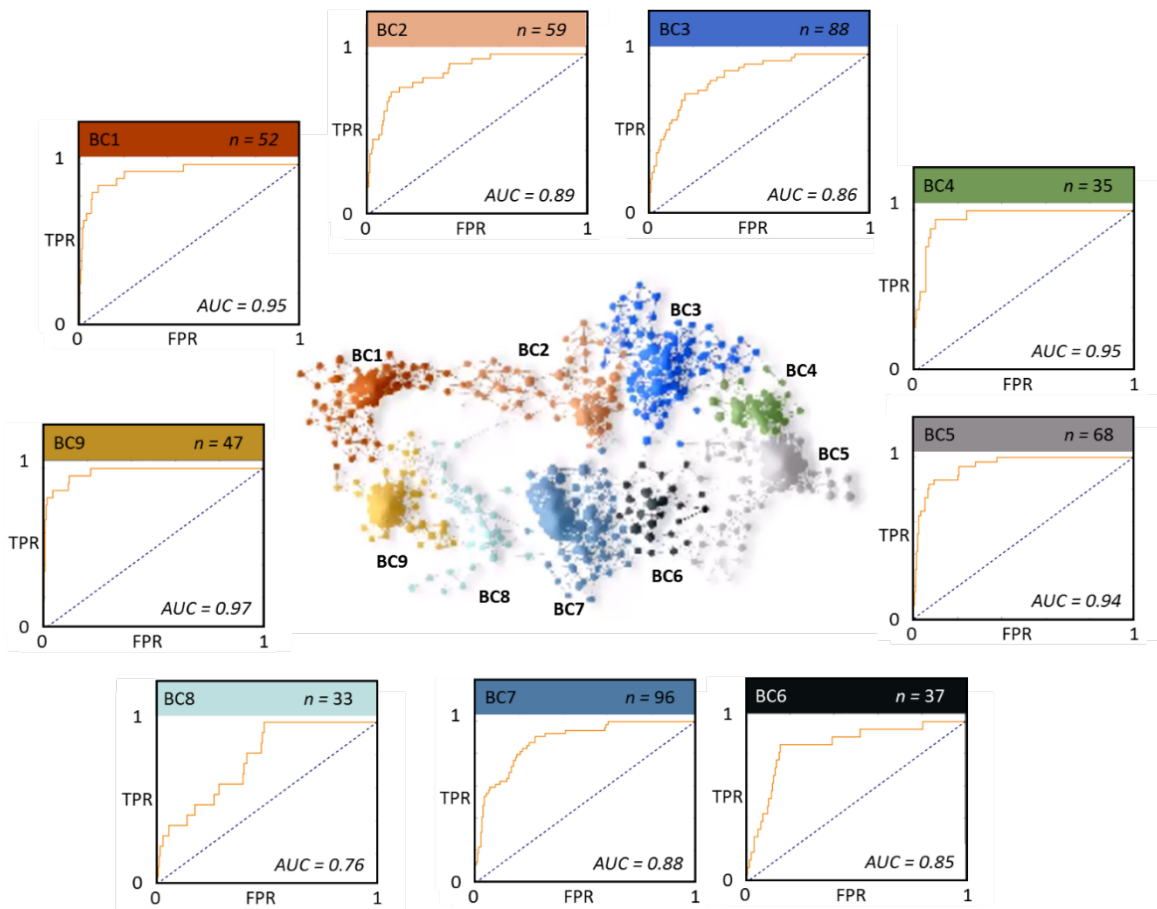
344

345 **Figure 3.** TDA network landscape of correlated gene expression (54,613 probesets, n = 498). Metric:  
346 norm correlation. Lenses: neighbourhood lens 1 (resolution, 100 bins; gain,  $\times 6$ ), neighbourhood lens 2  
347 (resolution, 100 bins; gain,  $\times 6$ ) (A). The vector (node value) is a 3<sup>rd</sup> dimension in TDA networks, in a  
348 standard heatmap colouring of a TDA network, the colour represents the 3<sup>rd</sup> dimension (B). Arrows  
349 indicate the gradients of the 3-dimensional topology measured by Morse-based clustering identifying  
350 the 'peaks' as clusters of subjects with similar profiles of analysed variables.

351

352

353 **Figure 4:** Clusters identified by Morse-clustering of the TDA network



354

355 **Figure 4. Centre:** TDA network coloured by clusters (BC1-9) identified using the Morse-based  
356 algorithm. **Outside:** Colour-coded ROC curves of cluster prediction success representative of cluster  
357 robustness.

358

359

360

361

362

363

364

365

366

367 **Table 1.** Molecular pathways enriched in the 9 clusters

Canonical Pathway	Sub-phenotype								
	BC1	BC9	BC8	BC2	BC7	BC3	BC6	BC4	BC5
Serotonin Degradation	3.1								
Superpathway of Melatonin Degradation	2.5								
Melatonin Degradation I	2.5								
Glutamate Receptor Signaling	2.4								
Neuropathic Pain Signaling In Dorsal Horn Neurons	2.4		-0.9	-1.8	0.0		-0.6		
Oxidative Phosphorylation		3.5	3.5		4.0	-4.4		-4.1	
Glycolysis I			3.0		2.8	-1.9		-2.5	
Role of p14/p19ARF in Tumor Suppression			1.4	0.0	0.3	-0.3		0.5	-0.9
Cyclins and Cell Cycle Regulation	2.1								
TNFR1 Signaling	1.9	1.7	0.9	2.1	0.3	-1.6	-2.2	-0.3	
tRNA Charging	1.4	2.7		3.1	-2.7		-1.6		
Gluconeogenesis I						-1.1		-1.7	
iNOS Signaling	0.8	2.3	3.3	3.5	3.1	-2.5	-2.2		
Toll-like Receptor Signaling			3.2		3.5				
Type I Diabetes Mellitus Signaling	1.0	2.1	3.0	3.3	2.4	-2.6	-2.9	-0.4	
TREM1 Signaling			2.9		3.7				
Neuroinflammation Signaling Pathway			2.7		2.3	-2.2			
IL-1 Signaling	-1.0	-0.2	2.5	1.5	2.8		-0.8	1.1	
Inflammasome pathway			2.4		2.6				
D-myo-inositol (1,4,5,6)-Tetrakisphosphate Biosynthesis	-3.0	-0.3		0.0	0.9	2.8	0.3	4.4	
D-myo-inositol (3,4,5,6)-tetrakisphosphate Biosynthesis	-3.0	-0.3		0.0	0.9	2.8	0.3	4.4	
3-phosphoinositide Biosynthesis	-3.8	-0.7		0.5	0.2	2.7	-0.3	4.5	
3-phosphoinositide Degradation	-3.0	-0.1		0.6	0.7	2.4	0.1	4.0	
Superpathway of Inositol Phosphate Compounds	-3.5	-0.5		1.2	0.0	2.1	-0.9	4.2	
Cell Cycle: G1/S Checkpoint Regulation	-1.7				0.6		2.0	1.4	
Antioxidant Action of Vitamin C	0.0		-0.7		-0.9		2.0		
HIPPO signaling	0.7	-0.5		-1.5		1.2		0.0	
Cardiac $\beta$ -adrenergic Signaling	-1.1		-2.2	-1.0					
ERK5 Signaling	-3.3	-1.3		0.2	1.1	1.8	1.6	2.0	
D-myo-inositol-5-phosphate Metabolism	-2.5	-0.2		0.8	0.7	2.1	0.0	4.3	

368 **Table 1.** IPA identified significantly enriched ( $p < 0.05$ ) canonical pathways of gene expression in  
 369 clusters (the top 5 pathways for clusters BC1-9 are shown). Values are z-scores, reflecting both the  
 370 enrichment of specific transcription factor-regulated genes in the pathways and the degree of  
 371 activation/inhibition. The z-scores are coloured blue (greatest downregulated transcription factor-  
 372 regulated gene expression) to red (greatest upregulated transcription factor-regulated gene expression).

373

374 **Table 2.** Activated upstream regulators enriched in the clusters

Upstream regulator	Sub-phenotype								
	BC1	BC9	BC8	BC2	BC7	BC3	BC6	BC4	BC5
CST5	3.45	2.56	2.01	3.24	1.69	2.02	-2.6	-1.5	-3.4
TP63						1.79	0.17		
HSF1			1.31			2.13			
TGM2				5.91		3.85	-4.4		
ERG		-1.6		-0.3			-1.4		-0.9
TAL1		3.31	2.42						
miR-486-5p (and other miRNAs w/seed CCUGUAC)		2.91		0.37	1.33	-1.2	-3.3	-2	-2.6
mir-486		2.89		0.24		-1.2	-3.3	-2.1	-2.6
NUPR1	0.76	2.86		2.98		2.54			
RAE1	1.34	2.83		0.45			-1.9		
SPP1		2.37				-2.2			
TFEB			2.98						
IL15		1.15	2.67	1.22		-0.8	-1.3	-1.5	
miR-30a-3p (and other miRNAs w/seed UUUUCAGU)		2.82	2.63	1.63			-1.3	-1.6	-2.2
EIF2AK2				3.05		1.44			
CEBPA				2.77		2.8			
PCGEM1					2.28	-1.2		-1.4	
LINC01139				1	2.24	0.45			
PLA2R1		1.25	1.04					-1.6	
LDL				1.39		1.93			
PPRC1						3.46			
PDGF BB						3.31			
TNF						3.11			
IL5				1.26					
CD24		-5.3	-5.2		-3.9	1	4.41	4.67	5.11
MYC	-2.9	-2.6		-4.5		-2	3.06	0.74	
HELLS		-1					2.45		2.24
MAPK1				-2					
SAFB				-2.1		-1.9	2.35		
SLC29A1		-1.2				1.63		2.65	1.41
WT1			-1.6	-1.1		1.61	-0.2		
FSH		-2.1	-2.3	-0.4		0.43	1.96	2.62	2.72
TCR				-0.7		-0.8		-1.8	2.49
THOC5		-2.2					1.63		2.45

375

376 **Table 2.** Upstream regulators of gene expression ( $p < 0.05$ ) in clusters predicted by IPA (the top 5  
 377 upstream regulators for clusters BC1-9 are shown). Values shown are z-scores, reflecting both the  
 378 enrichment of specific transcription factor-regulated genes in the pathways and the degree of  
 379 activation/inhibition. The z-scores are coloured blue of varying intensity (greatest downregulated  
 380 transcription factor-regulated gene expression) to varying red (greatest upregulated transcription factor-  
 381 regulated gene expression).

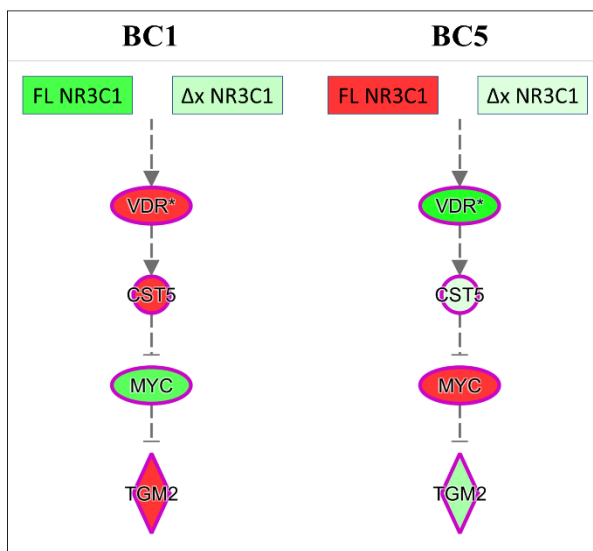
382

**Table 3.** Clinical characteristics of the clusters

Cluster	BC1	BC2	BC3	BC4	BC5	BC6	BC7	BC8	BC9
Number of participants	52 (10.44%)	59 (11.84%)	88 (17.67%)	35 (7.02%)	68 (13.65%)	37 (7.42%)	96 (19.27%)	33 (6.62%)	47 (9.43%)
FEV <sub>1</sub> (%)	72.21 ± 24.64	66.04 ± 20.76	67.89 ± 25.16	76.06 ± 23.28	79.63 ± 23.62	87 ± 21.03	83.33 ± 23.83	78.57 ± 22.97	71.69 ± 24
FVC (%)	88.97 ± 20.9	88.83 ± 19.43	85.52 ± 23.42	95.12 ± 24.11	98.13 ± 21.01	99.77 ± 17.12	98.67 ± 19.96	95.12 ± 22.5	86.68 ± 21.35
Severe Asthma (non-smoker) (%)	69.2	50.8	38.6	42.8	33.8	43.2	18.7	51.5	51
Severe Asthma (smoker) (%)	9.6	23.7	21.5	17.1	19.1	10.8	15.6	18.1	17
Mild-moderate Asthma (%)	9.6	11.8	9	22.8	13.2	8.1	25	18.1	10.6
Healthy (%)	11.5	1.6	9	17.1	22	37.8	23.9	12.1	21.2
Severe Asthma cluster (%)	75	81	39	22	25	5	5	3	2
Age	51.44 ± 14.73	53.03 ± 14.44	51.07 ± 14.45	46.88 ± 16.45	44.07 ± 13.97	44.51 ± 14.87	45.22 ± 14.95	47.57 ± 15.47	50.8 ± 15.58
Smoking (Pack Years)	3.3 ± 11.44	6.38 ± 16.00	5.05 ± 11.69	3.64 ± 7.11	4.59 ± 10.87	2.66 ± 7.69	3.69 ± 10.56	5.07 ± 10.72	5.87 ± 14.52
Mean ACQ5	1.69 ± 1.49	1.95 ± 1.23	1.83 ± 1.34	1.46 ± 1.51	1.44 ± 1.39	1.03 ± 1.41	1.18 ± 1.23	1.58 ± 1.36	1.65 ± 1.48
Mean ACQ7	2 ± 1.65	2.31 ± 1.36	2.17 ± 1.5	1.66 ± 1.65	1.67 ± 1.52	1.15 ± 1.52	1.4 ± 1.37	1.82 ± 1.46	1.98 ± 1.61
Mean AQLQ	3.68 ± 2.24	4.64 ± 1.57	4.08 ± 2	3.6 ± 2.52	3.98 ± 2.35	3.16 ± 2.81	3.78 ± 2.59	3.74 ± 2.48	3.36 ± 2.24
Admitted to ICU (%)	0.25 ± 0.4	0.2 ± 0.54	0.17 ± 0.37	0.17 ± 0.17	0.23 ± 0.19	0.05 ± 0.13	0.13 ± 0.13	0.18 ± 0.18	0.17 ± 0.19
Oral steroids (%)	40.38 ± 46.57	54.24 ± 38.46	37.50 ± 40.45	17.14 ± 41.23	19.12 ± 39.79	13.51 ± 45.32	13.54 ± 44.21	18.18 ± 46.09	19.15 ± 44.31
Blood periostin (ng/ml)	46.57 ± 24.62	38.46 ± 23.24	40.45 ± 27.57	41.23 ± 27.09	39.79 ± 22.02	45.32 ± 24.13	44.21 ± 21.45	46.09 ± 19.88	44.31 ± 23.59
Atopy (% positive)	0.65 ± 29.81	0.67 ± 31.71	0.67 ± 32.66	0.68 ± 36.58	0.72 ± 31.78	0.56 ± 30.74	0.67 ± 33.75	0.66 ± 28.72	0.8 ± 26.34
Exhaled NO (ppb)	29.81 ± 22.04	31.71 ± 30.11	32.66 ± 26.52	36.58 ± 32.73	31.78 ± 30.61	30.74 ± 32.05	33.75 ± 31.02	28.72 ± 26.51	26.34 ± 14.71
Blood eosinophils (x10 <sup>3</sup> /μL)	0.31 ± 0.3	0.18 ± 0.17	0.25 ± 0.28	0.21 ± 0.14	0.25 ± 0.25	0.23 ± 0.21	0.23 ± 0.2	0.29 ± 0.24	0.35 ± 0.33
Blood neutrophils (x10 <sup>3</sup> /μL)	5.63 ± 2.3	6.78 ± 2.94	5.41 ± 2.35	4.35 ± 1.52	4.18 ± 1.86	3.32 ± 1.37	3.42 ± 1.09	3.99 ± 1.2	4.06 ± 1.75
Blood lymphocytes (x10 <sup>3</sup> /μL)	2.06 ± 0.7	1.57 ± 0.7	1.83 ± 0.76	2 ± 0.47	1.91 ± 0.82	2.03 ± 0.73	1.87 ± 0.46	2.22 ± 0.66	2.14 ± 0.75
Sputum Eosinophils (%)	1.67 ± 5.16	6.37 ± 14.89	2.33 ± 9.42	1.77 ± 8.27	3.84 ± 12.49	5.79 ± 16.41	5.28 ± 12.42	4.47 ± 10.25	3.32 ± 12.41
Sputum Neutrophils (%)	30.18 ± 36.16	29.48 ± 34.25	5.7 ± 17.38	3.45 ± 12.12	17.37 ± 25.54	21.7 ± 28.31	28.65 ± 28.88	28.48 ± 29.83	24.83 ± 31.74
Sputum Macrophages (%)	13.65 ± 20.15	12.66 ± 17.96	3.16 ± 10.9	2.79 ± 9.48	17.89 ± 27.24	25.88 ± 33.44	30.62 ± 30.57	29.99 ± 30.84	26.38 ± 32.85
Sputum Lymphocytes (%)	0.62 ± 1.26	0.61 ± 1.06	0.15 ± 0.65	0.53 ± 2.29	0.57 ± 0.99	0.64 ± 0.84	1.04 ± 1.34	0.68 ± 0.95	0.74 ± 1.21

**Table 3.** Clinical features associated with the TDA-defined asthma phenotypes. Values are shown as means and are colour coded on a heat scale for each variable; highest variable value is in red, lowest value in blue. FEV<sub>1</sub>: forced expiratory volume in one second (measured by spirometry). FVC: forced vital capacity. (%) Severe Asthma cluster (%) is the percentage of study participants previously identified in the severe asthma enriched cluster identified by hierarchical clustering<sup>10</sup>. ACQ5 or 7: asthma quality questionnaire consisting of 5 or 7 questions. AQLQ: asthma quality of life questionnaire. Sputum cells are shown as percentages of total inflammatory cells.





**Figure 5.** The regulatory gene pathway of NR3C1 transcript variants, and VDR, CST5, MYC & TGM2; identified as top upstream regulators by IPA (Table 2). Colours indicate gene expression relative to healthy participants, where green represents lower gene expression and red represents higher gene expression, white indicates no change (negative, positive and zero-fold change). Left column shows gene expression in cluster BC1, right column shows gene expression in BC5. Image generated using IPA.

## References

- 1 Farne HA, Wilson A, Powell C, Bax L, Milan SJ. Anti-IL5 therapies for asthma. *Cochrane Libr* 2017.
- 2 Walker S, Monteil M, Phelan K, Lasserson TJ, Walters EH. Anti-IgE for chronic asthma in adults and children. *Cochrane Database Syst Rev* 2006; **2**.
- 3 Wheelock CE, Goss VM, Balgoma D, *et al*. Application of 'omics technologies to biomarker discovery in inflammatory lung diseases. *Eur Respir J* 2013; **42**: 802–25.
- 4 Shaw DE, Sousa AR, Fowler SJ, *et al*. Clinical and inflammatory characteristics of the European U-BIOPRED adult severe asthma cohort. *Eur Respir J* 2015; **46**: 1308–21.
- 5 Fleming L, Murray C, Bansal AT, *et al*. The burden of severe asthma in childhood and adolescence: results from the paediatric U-BIOPRED cohorts. *Eur Respir J* 2015; **46**: 1322–33.
- 6 Wilson SJ, Ward JA, Sousa AR, *et al*. Severe asthma exists despite suppressed tissue inflammation: findings of the U-BIOPRED study. *Eur Respir J* 2016; **48**: 1307–19.
- 7 Loza MJ, Adcock I, Auffray C, *et al*. Longitudinally stable, clinically defined clusters of patients with asthma independently identified in the ADEPT and U-BIOPRED asthma studies. *Ann Am Thorac Soc* 2016; **13**: S102–3.
- 8 Kuo C-HS, Pavlidis S, Loza M, *et al*. A transcriptome-driven analysis of epithelial brushings and bronchial biopsies to define asthma phenotypes in U-BIOPRED. *Am J Respir Crit Care Med* 2017; **195**: 443–55.
- 9 Lefaudeux D, De Meulder B, Loza MJ, *et al*. U-BIOPRED clinical adult asthma clusters linked to a subset of sputum omics. *J Allergy Clin Immunol* 2017; **139**: 1797–807.
- 10 Bigler J, Boedigheimer M, Schofield JPR, *et al*. A severe asthma disease signature from gene expression profiling of peripheral blood from U-BIOPRED cohorts. *Am J Respir Crit Care Med* 2017; **195**: 1311–20.
- 11 Soveg F, Abdala-Valencia H, Campbell J, Morales-Nebreda L, Mutlu GM, Cook-Mills JM. Regulation of allergic lung inflammation by endothelial cell transglutaminase 2. *Am J Physiol Cell Mol Physiol* 2015; **309**: L573–83.
- 12 Shaw DE, Berry MA, Hargadon B, *et al*. Association between neutrophilic airway inflammation and airflow limitation in adults with asthma. *Chest* 2007; **132**: 1871–5.
- 13 Woodruff PG, Modrek B, Choy DF, *et al*. T-helper type 2–driven inflammation defines major subphenotypes of asthma. *Am J Respir Crit Care Med* 2009; **180**: 388–95.
- 14 Bamberger CM, Bamberger A-M, de Castro M, Chrousos GP. Glucocorticoid receptor beta, a potential endogenous inhibitor of glucocorticoid action in humans. *J Clin Invest* 1995; **95**: 2435–41.

- 15 Schaaf MJM, Cidlowski JA. AUUUA motifs in the 3' UTR of human glucocorticoid receptor  $\alpha$  and  $\beta$  mRNA destabilize mRNA and decrease receptor protein expression. *Steroids* 2002; **67**: 627–36.
- 16 Lum PY, Singh G, Lehman A, *et al.* Extracting insights from the shape of complex data using topology. *Sci Rep* 2013; **3**: 1236.
- 17 Nielson JL, Paquette J, Liu AW, *et al.* Topological data analysis for discovery in preclinical spinal cord injury and traumatic brain injury. *Nat Commun* 2015; **6**: 8581.
- 18 Nicolau M, Levine AJ, Carlsson G. Topology based data analysis identifies a subgroup of breast cancers with a unique mutational profile and excellent survival. *Proc Natl Acad Sci U S A* 2011; **108**: 7265–70.
- 19 Landi C, Bargagli E, Carleo A, *et al.* A system biology study of BALF from patients affected by idiopathic pulmonary fibrosis (IPF) and healthy controls. *Proteomics Clin Appl* 2014; **8**: 932–50.
- 20 Hinks TSC, Zhou X, Staples KJ, *et al.* Innate and adaptive T cells in asthmatic patients: Relationship to severity and disease mechanisms. *J Allergy Clin Immunol* 2015; **136**: 323–33.
- 21 Forman R. A discrete Morse theory for cell complexes. In: in “Geometry, Topology 6 Physics for Raoul Bott. Citeseer, 1995.
- 22 Manni ML, Robinson KM, Alcorn JF. A tale of two cytokines: IL-17 and IL-22 in asthma and infection. *Expert Rev Respir Med* 2014; **8**: 25–42.
- 23 Yamaguchi M, Zacharia J, Laidlaw TM, Balestrieri B. PLA2G5 regulates transglutaminase activity of human IL-4-activated M2 macrophages through PGE2 generation. *J Leukoc Biol* 2016; **100**: 131–41.
- 24 Reiter RJ, Calvo JR, Karbownik M, Qi W, Tan DX. Melatonin and its relation to the immune system and inflammation. *Ann N Y Acad Sci* 2000; **917**: 376–86.
- 25 Kang BN, Ha SG, Bahaie NS, *et al.* Regulation of serotonin-induced trafficking and migration of eosinophils. *PLoS One* 2013; **8**: e54840.
- 26 Wang Y, Yang L, Li P, *et al.* Circulating microRNA signatures associated with childhood asthma. *Clin Lab* 2015; **61**: 467–74.
- 27 Song L, Lin C, Gong H, *et al.* miR-486 sustains NF- $\kappa$ B activity by disrupting multiple NF- $\kappa$ B-negative feedback loops. *Cell Res* 2013; **23**: 274.
- 28 Chung KF. p38 mitogen-activated protein kinase pathways in asthma and COPD. *Chest* 2011; **139**: 1470–9.
- 29 Vincent AJ, Ren S, Harris LG, *et al.* Cytoplasmic translocation of p21 mediates NUPR1-induced chemoresistance: NUPR1 and p21 in chemoresistance. *FEBS Lett* 2012; **586**: 3429–34.

- 30 Saude EJ, Obiefuna IP, Somorjai RL, *et al.* Metabolomic biomarkers in a model of asthma exacerbation: urine nuclear magnetic resonance. *Am J Respir Crit Care Med* 2009; **179**: 25–34.
- 31 Valle N, García JM, Peña C, *et al.* Cystatin D is a candidate tumor suppressor gene induced by vitamin D in human colon cancer cells. *J Clin Invest* 2009; **119**: 2343–58.
- 32 Hidalgo AA, Deeb KK, Pike JW, Johnson CS, Trump DL. Dexamethasone enhances 1 $\alpha$ , 25-dihydroxyvitamin D3 effects by increasing vitamin D receptor transcription. *J Biol Chem* 2011; : jbc-M111.

## The U-BIOPRED Study Group

A. Bautmans<sup>16</sup>, A. Chaiboonchoe<sup>13</sup>, A. Mazein<sup>13</sup>, A. Sogbesan<sup>17</sup>, A. Meiser<sup>4</sup>, A. Menzies-Gow<sup>17</sup>, A. Berglind<sup>18</sup>, A.-S. Lantz<sup>7</sup>, A.J. James<sup>8</sup>, A. Petré<sup>8</sup>, A.F. Behndig<sup>19</sup>, A. Dijkhuis<sup>13</sup>, A. Postle<sup>20</sup>, A. Rowe<sup>21</sup>, A. Vink<sup>22</sup>, A. Pacino<sup>23</sup>, A. Aliprantis<sup>24</sup>, A. Wagener<sup>13</sup>, A. Braun<sup>25</sup>, A. D'Amico<sup>26</sup>, A. Woodcock<sup>27</sup>, B. Smids<sup>13</sup>, B. Lambrecht<sup>28</sup>, B. Nicholas<sup>20</sup>, B. Nordlund<sup>18</sup>, B. Thornton<sup>29</sup>, A. Roberts<sup>30</sup>, B. Flood<sup>30</sup>, C. Mathon<sup>31</sup>, C. Smith<sup>32</sup>, C. Holweg<sup>33</sup>, C. Compton<sup>10</sup>, C. von Garnier<sup>34</sup>, C. Rossios<sup>4</sup>, C. Barber<sup>15</sup>, C.S. Murray<sup>27</sup>, C. Wiegman<sup>4</sup>, C. Schoelch<sup>35</sup>, C. Faulenbach<sup>36</sup>, C. Coleman<sup>30</sup>, C. Gomez<sup>8</sup>, D. Erzen<sup>35</sup>, D. Balgoma<sup>8</sup>, D. Gibeon<sup>4</sup>, D. Myles<sup>10</sup>, D. Supple<sup>30</sup>, D. Campagna<sup>37</sup>, D. Burg<sup>1</sup>, D.E. Shaw<sup>38</sup>, D. Staykova<sup>20</sup>, E. Bel<sup>13</sup>, E. Henriksson<sup>39</sup>, E. Yeyasingham<sup>40</sup>, E. Ray<sup>32</sup>, E.J. Kennington<sup>30</sup>, F. Singer<sup>41</sup>, F. Wald<sup>35</sup>, F. Baribaud<sup>42</sup>, G. Galffy<sup>43</sup>, G. Pennazza<sup>26</sup>, G. Santini<sup>44</sup>, G. Roberts<sup>45</sup>, G. Bochenek<sup>46</sup>, G. Hedlin<sup>18</sup>, H. Bisgaard<sup>47</sup>, H. Ahmed<sup>11</sup>, H. Gallart<sup>8</sup>, H. Knobel<sup>22</sup>, I. Horvath<sup>43</sup>, I. De Lepeleire<sup>16</sup>, I. Delin<sup>8</sup>, J. Musial<sup>46</sup>, J. Martin<sup>32</sup>, J. Versnel<sup>30</sup>, J. Hohlfeld<sup>25</sup>, J. Edwards<sup>30</sup>, J. Smith<sup>30</sup>, J.P. Carvalho da Purificação Rocha<sup>17</sup>, J. Kolmert<sup>8</sup>, J.G. Matthews<sup>33</sup>, J. Haughney<sup>48</sup>, J.-O. Thörngren<sup>49</sup>, J. Konradsen<sup>18</sup>, J. Thorsen<sup>47</sup>, J. Ward<sup>50</sup>, J. Brandsma<sup>20</sup>, J. Beleta<sup>51</sup>, J. De Alba<sup>51</sup>, J. Östling<sup>52</sup>, J. Vestbo<sup>27</sup>, J. Gent<sup>53</sup>, J. Corfield<sup>54</sup>, J. Kamphuis<sup>55</sup>, K. Tariq<sup>56</sup>, K. Strandberg<sup>39</sup>, A. Knox<sup>57</sup>, K.M. Smith<sup>57</sup>, K. Riemann<sup>35</sup>, K. Nething<sup>35</sup>, K. van Drunen<sup>13</sup>, K. Dyson<sup>58</sup>, K. Gove<sup>15</sup>, K. Russell<sup>4</sup>, K. Alving<sup>59</sup>, K. Bønnelykke<sup>47</sup>, K. Fichtner<sup>35</sup>, K. Zwinderman<sup>13</sup>, K. Wetzel<sup>35</sup>, L. Ravanetti<sup>13</sup>, L. Larsson<sup>60</sup>, L. Pahus<sup>61</sup>, L. Metcalf<sup>30</sup>, L. Carayannopoulos<sup>29</sup>, L. Tamasi<sup>43</sup>, L. Krueger<sup>62</sup>, L. Marouzet<sup>32</sup>, L. Hewitt<sup>32</sup>, L.J. Fleming<sup>4</sup>, M. Kupczyk<sup>8</sup>, M. Ericsson<sup>63</sup>, M. Rahman-Amin<sup>30</sup>, M. Santoninco<sup>26</sup>, M. Sjödin<sup>8</sup>, A. Berton<sup>52</sup>, M. Gerhardsson de Verdier<sup>52</sup>, M. Mikus<sup>64</sup>, M. van de Pol<sup>13</sup>, M. van Geest<sup>52</sup>, M. Gahlemann<sup>65</sup>, Basel, Switzerland<sup>66</sup>, M. Robberechts<sup>16</sup>, M. Szentkereszty<sup>43</sup>, M. Caruso<sup>37</sup>, M.J. Loza<sup>67</sup>, M. Klüglich<sup>35</sup>, M. Kots<sup>68</sup>, M. Rutgers<sup>55</sup>, M. Miralpeix<sup>51</sup>, N. Mores<sup>44</sup>, N. Vissing<sup>47</sup>, N. Rao<sup>69</sup>, N. Fitch<sup>70</sup>, N. Gozzard<sup>71</sup>, N. Lazarinis<sup>39</sup>, N. Adriaens<sup>13</sup>, N. Krug<sup>25</sup>, P.J. Carvalho<sup>4</sup>, P. Söderman<sup>72</sup>, P. Montuschi<sup>44</sup>, P. Chanez<sup>73</sup>, P. Dennison<sup>56</sup>, P. Brinkman<sup>13</sup>, P. Bakke<sup>74</sup>, P. Howarth<sup>75</sup>, P. Nilsson<sup>64</sup>, P. Monk<sup>76</sup>, P. Badorrek<sup>36</sup>, P.-P. Hekking<sup>13</sup>, P. de Boer<sup>55</sup>, P. Powell<sup>77</sup>, R. Sigmund<sup>35</sup>, R. Lutter<sup>13</sup>, R. Hu<sup>3</sup>, R. Middelveld<sup>8</sup>, R. Chaleckis<sup>31</sup>, R. Emma<sup>37</sup>, S. Lone-Latif<sup>13</sup>, S. Meah<sup>4</sup>, S. Valente<sup>44</sup>, S. Walker<sup>30</sup>, S. Pink<sup>32</sup>, S. Masefield<sup>77</sup>, S. Kuo<sup>4</sup>, S. Wagers<sup>70</sup>, S. Naz<sup>8</sup>, S. Williams<sup>48</sup>, S. Hu<sup>4</sup>, S. Hashimoto<sup>13</sup>, S. Reinke<sup>8</sup>, S. Pavlidis<sup>4</sup>, S.J. Fowler<sup>27</sup>, S.J. Wilson<sup>50</sup>, S. Palkonen<sup>79</sup>, S.-E. Dahlén<sup>8</sup>, T. Dekker<sup>13</sup>, T. Geiser<sup>80</sup>, T. Sandström<sup>19</sup>, T. Higgenbottam<sup>81</sup>, U. Nihlen<sup>52</sup>, U. Frey<sup>82</sup>, U. Hoda<sup>83</sup>, V. Hudson<sup>30</sup>, V. Erpenbeck<sup>84</sup>, W. Yu<sup>3</sup>, W. Zetterquist<sup>18</sup>, W. van Aalderen<sup>13</sup>, W. Seibold<sup>35</sup>, X. Yang<sup>4</sup>, X. Hu<sup>3</sup>, Y.-k. Guo<sup>9</sup>, Z. Weiszhart<sup>85</sup>.

<sup>16</sup>MSD, Brussels, Belgium

<sup>17</sup>Royal Brompton and Harefield NHS Foundation Trust, London, UK

<sup>18</sup>Dept of Women's and Children's Health and Centre for Allergy Research, Karolinska Institutet, Stockholm, Sweden

<sup>19</sup>Dept of Public Health and Clinical Medicine, Umeå University, Umeå, Sweden

<sup>20</sup>University of Southampton, Southampton, UK

<sup>21</sup>Janssen R&D, High Wycombe, UK

<sup>22</sup>Philips Research Laboratories, Eindhoven, The Netherlands

<sup>23</sup>Lega Italiano Anti Fumo, Catania, Italy

<sup>24</sup>Merck Research Laboratories, Boston, MA, USA

<sup>25</sup>Fraunhofer Institute for Toxicology and Experimental Medicine, Hannover, Germany

<sup>26</sup>University of Rome "Tor Vergata", Rome, Italy

<sup>27</sup>Centre for Respiratory Medicine and Allergy, Institute of Inflammation and Repair, University of Manchester and University Hospital of South Manchester, Manchester Academic Health Sciences Centre, Manchester, UK

<sup>28</sup>University of Gent, Gent, Belgium

<sup>29</sup>MSD, Kenilworth, NJ, USA

<sup>30</sup>Asthma UK, London, UK

<sup>31</sup>Centre of Allergy Research, Karolinska Institutet, Stockholm, Sweden

<sup>32</sup>NIHR Southampton Respiratory Biomedical Research Unit, Southampton, UK

<sup>33</sup>Respiratory and Allergy Diseases, Genentech, San Francisco, CA, USA

<sup>34</sup>University Hospital Bern, Bern, Switzerland

<sup>35</sup>Boehringer Ingelheim Pharma, Biberach, Germany

<sup>36</sup>Fraunhofer ITEM, Hannover, Germany

<sup>37</sup>Dept of Clinical and Experimental Medicine, University of Catania, Catania, Italy

<sup>38</sup>Respiratory Research Unit, University of Nottingham, Nottingham, UK

<sup>39</sup>Karolinska University Hospital and Karolinska Institutet, Stockholm, Sweden

<sup>40</sup>UK Clinical Operations, GSK, Uxbridge, UK

<sup>41</sup>University Children's Hospital, Zurich, Switzerland

<sup>42</sup>Janssen R&D, Spring House, PA USA

<sup>43</sup>Semmelweis University, Budapest, Hungary

<sup>44</sup>Università Cattolica del Sacro Cuore, Milan, Italy

<sup>45</sup>NIHR Southampton Respiratory Biomedical Research Unit, Clinical and Experimental Sciences and Human Development and Health, Southampton, UK

<sup>46</sup>II Dept of Internal Medicine, Jagiellonian University Medical College, Krakow, Poland

<sup>47</sup>COPSAC, Copenhagen Prospective Studies on Asthma in Childhood, Herlev and Gentofte Hospital, University of Copenhagen, Copenhagen, Denmark

<sup>48</sup>International Primary Care Respiratory Group, Aberdeen, UK

<sup>49</sup>Karolinska University Hospital, Sweden

<sup>50</sup>Histochemistry Research Unit, Faculty of Medicine, University of Southampton, Southampton, UK

<sup>51</sup>Almirall, Barcelona, Spain

<sup>52</sup>AstraZeneca, Mölndal, Sweden

<sup>53</sup>Royal Brompton and Harefield NHS Foundation Trust, UK

<sup>54</sup>Areteva R&D, Nottingham, UK

<sup>55</sup>Longfonds, Amersfoort, The Netherlands

<sup>56</sup>NIHR Southampton Respiratory Biomedical Research Unit, Clinical and Experimental Sciences, NIHR-Wellcome Trust Clinical Research Facility, Faculty of Medicine, University of Southampton, Southampton, UK

<sup>57</sup>University of Nottingham, Nottingham, UK

<sup>58</sup>CromSource, Stirling UK

<sup>59</sup>Dept of Women's and Children's Health, Uppsala University, Sweden

<sup>60</sup>AstraZeneca, Mölndal, Sweden

<sup>61</sup>Assistance publique des Hôpitaux de Marseille, Clinique des bronches, allergies et sommeil Espace Éthique Méditerranéen, Aix-Marseille Université, Marseille, France

<sup>62</sup>University Children's Hospital Bern, Bern, Switzerland

<sup>63</sup>Karolinska University Hospital, Stockholm, Sweden

<sup>64</sup>Science for Life Laboratory and The Royal Institute of Technology, Stockholm, Sweden

<sup>65</sup>Boehringer Ingelheim

<sup>66</sup>Janssen R&D, Spring House, PA, USA

<sup>67</sup>Chiesi Pharmaceuticals, Parma, Italy

<sup>68</sup>Janssen R&D, San Diego, CA, USA

<sup>69</sup>BioSci Consulting, Maasmechelen, Belgium

<sup>70</sup>UCB, Slough, UK

<sup>71</sup>Dept of Women's and Children's Health, Karolinska Institutet, Stockholm, Sweden

<sup>72</sup>Assistance publique des Hôpitaux de Marseille, Clinique des bronches, allergies et sommeil, Aix-Marseille Université, Marseille, France

<sup>73</sup>Dept of Clinical Science, University of Bergen, Bergen, Norway

<sup>74</sup>NIHR Southampton Respiratory Biomedical Research Unit, Clinical and Experimental Sciences, Southampton, UK

<sup>75</sup>Synairgen Research, Southampton, UK

<sup>76</sup>European Lung Foundation, Sheffield, UK

<sup>77</sup>European Federation of Allergy and Airways Diseases Patient's Associations, Brussels, Belgium

<sup>78</sup>Dept of Respiratory Medicine, University Hospital Bern, Switzerland

<sup>79</sup>Allergy Therapeutics, Worthing, UK

<sup>80</sup>University Children's Hospital, Basel, Switzerland

<sup>81</sup>Imperial College, London, UK

<sup>82</sup>Translational Medicine, Respiratory Profiling, Novartis Institutes for Biomedical Research, Basel, Switzerland

<sup>83</sup>Semmelweis University, Budapest, Hungary

<sup>84</sup>University of Southampton, Southampton

## Author Contribution statement

JPRS and PJSkipp wrote the main manuscript text. JPRS prepared all figures. JPRS, FS, PJSkipp & R S-G developed the methodology for Morse clustering in a TDA network. KS and IP processed, integrated and curated gene expression and patient clinical and demographic data. JPRS, JB, MB, IA, KFC, AB, RK, S-ED, CW, JR, CA, BDM, DL, DR, AS, PJSterk, RE, BM, RD, R S-G and PJS planned the investigation and contributed to revising the manuscript.

## Acknowledgments

This paper is presented on behalf of the U-BIOPRED Study Group with input from the U-BIOPRED Patient Input Platform, Ethics Board and Safety Management Board. We thank all the members of each recruiting centre for their dedicated effort, devotion, promptness and care in the recruitment and assessment of the participants in this study. U-BIOPRED is supported through an Innovative Medicines Initiative Joint Undertaking under grant agreement no. 115010, resources of which are composed of financial contribution from the European Union's Seventh Framework Programme (FP7/2007–2013) and EFPIA companies' in-kind contribution ([www.imi.europa.eu](http://www.imi.europa.eu)). We would also like to acknowledge help from the IMI funded eTRIKS project (EU Grant Code No.115446).

The U-BIOPRED consortium wishes to acknowledge the help and expertise of the following individuals and groups without whom, the study would not have been possible. Investigators and contributors: Nora Adriaens, Antonios Aliprantis, Kjell Alving, Per Bakke, David Balgoma, Clair Barber, Frédéric Baribaud, Stewart Bates, An Bautmans, Jorge Beleta, Grazyna Bochenek, Joost Brandsma, Armin Braun, Dominic Burg, Leon Carayannopoulos, João Pedro Carvalho da Purificação Rocha, Romanas Chaleckis, Arnaldo D'Amico, Jorge De Alba, Tamara Dekker, Annemiek Dijkhuis, Aleksandra Draper, Rosalia Emma, Magnus Ericsson, Breda Flood, Hector Gallart, Kerry Gove, Neil Gozzard, Lorraine Hewitt, Jens Hohlfeld, Cecile Holweg, Richard Hu, Sile Hu, Juliette Kamphuis, Erika J. Kennington, Dyson Kerry, Hugo Knobel, Johan Kolmert, Maxim Kots, Scott Kuo, Maciej Kupczyk, Bart Lambrecht, Saeeda Lone-Latif, Lisa Marouzet, Jane Martin, Sarah Masefield, Caroline Mathon, Sally Meah, Andrea Meiser, Leanne Metcalf, Montse Miralpeix, Shama Naz, Ben Nicholas, Peter Nilsson, Jörgen Östling, Antonio Pacino, Susanna Palkonen, Stelios Pavlidis, Giorgio Pennazza, Anne Petré, Sandy Pink, Anthony Postle, Malayka Rahman-Amin, Navin Rao, Lara Ravanetti, Emma Ray, Stacey Reinke, Leanne Reynolds, John Riley, Martine Robberechts, Amanda Roberts, Kirsty Russell, Michael Rutgers, Marco Santoninco, Corinna Schoelch, James P.R. Schofield, Marcus Sjödin, Paul J. Skipp, Barbara Smids, Caroline Smith, Jessica Smith, Doroteya Staykova, Kai Sun, John-Olof Thörngren, Bob Thornton, Jonathan Thorsen, Marianne van de Pol, Marleen van Geest, Anton Vink, Frans Wald, Samantha Walker, Jonathan Ward, Zsoka Weiszhart, Kristiane Wetzel, Craig E. Wheelock, Coen Wiegman, Siân Williams, Susan J. Wilson, Ashley Woodcock, Xian Yang, Elizabeth Yeyasingham.

Partner organisations: Novartis Pharma AG; University of Southampton, Southampton, UK; Academic Medical Centre, University of Amsterdam, Amsterdam, The Netherlands; Imperial College London, London, UK; University of Catania, Catania, Italy; University of Rome 'Tor Vergata', Rome, Italy; Hvidovre Hospital, Hvidovre, Denmark; Jagiellonian Univ. Medi.College, Krakow, Poland; University Hospital, Inselspital, Bern, Switzerland; Semmelweis University, Budapest, Hungary; University of Manchester, Manchester, UK; Université d'Aix-Marseille, Marseille, France; Fraunhofer Institute, Hannover, Germany; University Hospital, Umea,



Sweden; Ghent University, Ghent, Belgium; Ctr. Nat. Recherche Scientifique, Villejuif, France; Università Cattolica del Sacro Cuore, Rome, Italy; University Hospital, Copenhagen, Denmark; Karolinska Institutet, Stockholm, Sweden; Nottingham University Hospital, Nottingham, UK; University of Bergen, Bergen, Norway; Netherlands Asthma Foundation, Leusden, NL; European Lung Foundation, Sheffield, UK; Asthma UK, London, UK; European Fed. of Allergy and Airways Diseases Patients' Associations, Brussels, Belgium; Lega Italiano Anti Fumo, Catania, Italy; International Primary Care Respiratory Group, Aberdeen, Scotland; Philips Research Laboratories, Eindhoven, NL; Synairgen Research Ltd, Southampton, UK; Aerocrine AB, Stockholm, Sweden; BioSci Consulting, Maasmechelen, Belgium; Almira; AstraZeneca; Boehringer Ingelheim; Chiesi; GlaxoSmithKline; Roche; UCB; Janssen Biologics BV; Amgen NV; Merck Sharp & Dohme Corp.

Third Parties to the project, contributing to the clinical trial: Academic Medical Centre (AMC), Amsterdam (In the U-BIOPRED consortium the legal entity is AMC Medical Research BV (AMR); AMR is a subsidiary of both AMC and the University of Amsterdam; AMC contribute across the U-BIOPRED project); University Hospital Southampton NHS Trust (third party of the University of Southampton and contributor to the U-BIOPRED clinical trial); South Manchester Healthcare Trust (third party to the University of Manchester, South Manchester Healthcare Trust, contributor to the U-BIOPRED clinical trial and to the U-BIOPRED Biobank); Protisvalor Méditerranée SAS (third party to University of the Mediterranean; contributor to the U-BIOPRED clinical trial); Karolinska University Hospital (third party Karolinska Institutet (KI), contributor to the U-BIOPRED clinical trial); Nottingham University Hospital (third party to University of Nottingham, contributor to the U-BIOPRED clinical trial); NIHR-Wellcome Trust Clinical Research Facility.

Members of the ethics board: Jan-Bas Prins, biomedical research, LUMC, the Netherlands; Martina Gahlemann, clinical care, BI, Germany; Luigi Visintin, legal affairs, LIAF, Italy; Hazel Evans, paediatric care, Southampton, UK; Martine Puhl, patient representation (co-chair), NAF, the Netherlands; Lina Buzermaniene, patient representation, EFA, Lithuania; Val Hudson, patient representation, Asthma UK; Laura Bond, patient representation, Asthma UK; Pim de Boer, patient representation and pathobiology, IND; Guy Widdershoven, research ethics, VUMC, the Netherlands; Ralf Sigmund, research methodology and biostatistics, BI, Germany.

The patient input platform: Amanda Roberts, UK; David Supple (chair), UK; Dominique Hamerlijnck, The Netherlands; Jenny Negus, UK; Juliëtte Kamphuis, The Netherlands; Lehanne Sergison, UK; Luigi Visintin, Italy; Pim de Boer (co-chair), The Netherlands; Susanne Onstein, The Netherlands.

Members of the safety monitoring board: William MacNee, clinical care; Renato Bernardini, clinical pharmacology; Louis Bont, paediatric care and infectious diseases; Per-Ake Wecksell, patient representation; Pim de Boer, patient representation and pathobiology (chair); Martina Gahlemann, patient safety advice and clinical care (co-chair); Ralf Sigmund, bio-informatician.

This work was partially funded by the Engineering and Physical Sciences Research Council, UK (EP/N014189: Joining the Dots, from Data to Insight).