



This is a repository copy of *Diagnostic host gene signature to accurately distinguish enteric fever from other febrile diseases.*

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/135962/>

Version: Submitted Version

Article:

Blohmke, C., Muller, J. orcid.org/0000-0002-1046-2968, Gibani, M. et al. (15 more authors) (Submitted: 2018) Diagnostic host gene signature to accurately distinguish enteric fever from other febrile diseases. BioRxiv. (Submitted)

<https://doi.org/10.1101/327429>

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1 **Diagnostic host gene signature to accurately distinguish enteric fever from other febrile diseases**

2

3 Blohmke CJ^{1*}, Muller J², Gibani MM¹, Dobinson H¹, Shrestha S¹, Perinparajah S¹, Jin C¹, Hughes H¹,
4 Blackwell L¹, Dongol S³, Karkey A³, Schreiber F⁴, Pickard D⁴, Basnyat B³, Dougan G⁴, Baker S⁵,
5 §Pollard AJ¹, §Darton TC^{1,5,6}

6

7 §Joint senior authorship

8

9 Affiliations:

10 1 Oxford Vaccine Group, Centre for Clinical Vaccinology and Tropical Medicine, Department of
11 Paediatrics, and the Oxford National Institute of Health Research Biomedical Centre, University of
12 Oxford, UK.

13 2 The Jenner Institute, University of Oxford, Old Road Campus Research Building, Oxford, UK.

14 3 Oxford University Clinical Research Unit, Patan Academy of Healthy Sciences, Kathmandu, Nepal.

15 4 Infection Genomics Program, The Wellcome Trust Sanger Institute, Hinxton, UK.

16 5 The Hospital for Tropical Diseases, Wellcome Trust Major Overseas Programme, Oxford
17 University Clinical Research Unit, Ho Chi Minh City, Vietnam.

18 6 Department of Infection, Immunity and Cardiovascular Disease, University of Sheffield, Sheffield,
19 United Kingdom.

20

21 Correspondence:

22 *Christoph J. Blohmke,

23 NIHR Oxford Biomedical Research Centre,

24 The Churchill Hospital, Oxford OX3 7LE, UK.

25 eMail: christoph.blohmke@paediatrics.ox.ac.uk

26

27

28 **ABSTRACT**

29 Misdiagnosis of enteric fever is a major global health problem resulting in patient mismanagement,
30 antimicrobial misuse and inaccurate disease burden estimates. Applying a machine-learning algorithm
31 to host gene expression profiles, we identified a diagnostic signature which could accurately
32 distinguish culture-confirmed enteric fever cases from other febrile illnesses (AUROC>95%).
33 Applying this signature to a culture-negative suspected enteric fever cohort in Nepal identified a
34 further 12.6% as likely true cases. Our analysis highlights the power of data-driven approaches to
35 identify host-response patterns for the diagnosis of febrile illnesses. Expression signatures were
36 validated using qPCR highlighting their utility as PCR-based diagnostic for use in endemic settings.

37 Enteric fever, a disease caused by systemic infection with *S. enterica* serovars Typhi or Paratyphi A,
38 accounts for 13.5 to 26.9 million illness episodes worldwide each year.^{1,2} In resource-limited tropical
39 settings these infections are endemic and the accurate diagnosis of patients presenting with
40 undifferentiated fever is challenging.

41 Diagnostic tests for enteric fever rely on microbiological culture or detection of a serological response
42 to infection, and are often unavailable or insufficiently sensitive and specific.³ Blood culture remains
43 the reference standard against which new diagnostic tests are evaluated, and the sensitivity for this test
44 can reach 80% under optimal conditions⁴ but low blood volumes and uncontrolled antibiotic use often
45 result in decreased sensitive in the field. New diagnostic approaches are urgently needed to enable the
46 accurate detection of enteric fever cases in endemic settings, to guide management of febrile patients,
47 appropriate use of antimicrobials, and to identify populations likely to benefit from vaccine
48 implementation.

49 Most common tests used for acute infectious disease diagnosis employ methods to directly detect the
50 disease-causing pathogen, either by culture, antigen detection or amplification of genetic material by
51 PCR. An alternative approach is to identify a set of human host immune responses, which together
52 may generate a specific pattern associated with individual infections or pathogens. With an increasing
53 quantity of molecular host response data being generated by high-throughput methods – including
54 whole blood gene expression profiling – differences in the activation status of the immune response
55 network during infection may be a tractable diagnostic approach. Recently small sets containing 2-3
56 genes have been described, the expression of which can accurately differentiate between viral or
57 bacterial infection, and active or latent tuberculosis.^{5,6} Merging available well-characterised datasets
58 derived from human clinical samples representative of a variety of fever-causing infections common
59 in tropical settings presents an invaluable resource to identify host immune response patterns specific
60 for enteric fever.

61 As a human restricted infection, the development of enteric fever diagnostics has been hindered by the
62 lack of reliable *in vivo* models. Using data from a series of controlled human infection models
63 (CHIM)^{4,7} or *S. Typhi* or *S. Paratyphi A* infection, whole blood gene transcriptional responses were
64 identified and then further characterised using samples collected from febrile patients in an endemic
65 setting (Kathmandu, Nepal). Integrating these data with publically available human gene transcription
66 datasets, we employed a machine learning algorithm to identify an expression signature that could
67 accurately distinguish blood culture-confirmed EF cases in both the controlled environment (CHIM)
68 and endemic setting from other febrile disease aetiologies and non-infected individuals (healthy
69 controls).⁸⁻¹²

70

71 **Results**

72 *Transcriptional profiles in response to enteric fever are similar in challenge study and endemic*
73 *cohorts*

74 We recently described the molecular response profile of acute enteric fever in individuals
75 participating in the typhoid CHIM, which was characterized by innate immunity, inflammatory and
76 interferon signalling patterns.¹³

77 To compare responses to enteric fever occurring during natural infection in an endemic area, we
78 generated transcriptional profiles in samples collected from culture-confirmed enteric fever patients
79 (*S. Typhi*: ‘03NP-ST’; *S. Paratyphi*: ‘03NP-SPT’), healthy community controls (‘03NP-CTRL’) and
80 febrile, culture-negative suspected enteric fever cases (‘03NP-sEF’) recruited in Nepal (Kathmandu;
81 Study: ‘03NP’) (**Figure 1a**). We detected significant differential expression (DE; FDR<0.05,
82 FC±1.25) of 4,308 and 4,501 genes in enteric fever patients with confirmed *S. Typhi* ($n=19$) and *S.*
83 *Paratyphi* ($n=12$) bacteraemia, respectively, when compared with healthy community controls ($n=47$;
84 **Figure 1b**). Similar numbers of genes were differentially expressed in samples collected at the time of
85 enteric fever diagnosis in healthy adult volunteers challenged with either *S. Typhi* (‘T1-ST’) or *S.*
86 *Paratyphi* (‘P1-SPT’) in a CHIM (**Figure 1b**).^{7,13}

87 As comparison of host responses at the gene level can be difficult to interpret, we performed Gene Set
88 Enrichment Analysis (GSEA)¹⁴ of blood transcriptional modules (BTMs) as a conceptual framework
89 to interpret the host responses in the context of biological pathways and themes.¹⁵ Overall, between 54
90 and 74 BTMs were significantly enriched (BH adjusted $p<0.01$) in blood culture-confirmed enteric
91 fever cases in the CHIM and natural infection and CHIM participants who did not develop enteric
92 fever (measured at day 7 post-challenge - ‘nD7’) (**Supplementary Table 1**). The majority of BTMs
93 enriched in cases from the enteric fever CHIM were also enriched in naturally infected cases from
94 Nepal (56%-69, **Supplementary Table 1**—red squares). Positively enriched modules represented cell
95 cycle (CCY), type I/II interferon and innate antiviral responses (IFN), dendritic cell (DC), innate
96 immunity, inflammation and monocyte (Infl./Mono) signatures. In contrast, T cell (TC) signatures
97 were down-regulated in patients with confirmed enteric fever, as we have previously described
98 (**Figure 1c-e**).¹³ In addition a number of modules including inflammasome receptors (M53),
99 monocyte enrichment (M118.0, M118.1, M81, M4.15, M23, M73, M64, S4) and inflammatory
100 responses (M33) were significantly enriched in the CHIM but not in cases from Nepal. Single sample
101 GSEA (ssGSEA) demonstrated the similar enrichment pattern for a selection of IFN and DC
102 signatures between individuals with confirmed typhoid and paratyphoid fever in the CHIM and
103 naturally infected cases (**Figure 1f**). Overall, we observed marked similarity in the gene transcription
104 responses between acute enteric fever cases from the CHIM and an endemic setting in Nepal.

105

106

107 *Responses of febrile, culture-negative samples in Nepal*

108 In culture-negative, suspected enteric fever patients ('sEF') from Nepal, we detected differential
109 expression of 3,517 genes when compared with healthy community controls (**Supplementary Figure**
110 **1b**). While we observed 2,843 genes as commonly expressed in all three Nepali patient cohorts
111 (03NP-ST, 03NP-SPT and 03NP-sEF), an additional 582, 756 and 183 genes were uniquely expressed
112 by subjects with confirmed *S. Typhi*, *S. Paratyphi* or suspected enteric fever, respectively
113 (**Supplementary Figure 1a&b**). Unsupervised hierarchical clustering of these patients based on their
114 expression of the 500 most variable genes in the Nepal cohort demonstrated clustering into three
115 groups (**Figure 1g**): Group 1 contained mostly healthy control participants; Group 2 contained mostly
116 patients with suspected enteric fever; and Group 3 contained a mixture of patients with suspected
117 enteric fever, and blood culture-confirmed *S. Typhi* or *S. Paratyphi* infection.

118 Using ssGSEA we observed a heterogeneous BTM enrichment pattern with broad variability in
119 normalized enrichment scores across suspected enteric fever patients (depicted by the interdecile
120 range; **Supplementary Figure 1c**). The most consistent positively or negatively enriched modules
121 represented cell cycle, IFN, inflammatory responses, DC and some NK cell signatures (green cluster)
122 and TC and BC related signatures (red cluster), respectively. In contrast, heterogeneous enrichment in
123 which approximately half of participant samples demonstrated up or down regulation was observed in
124 BTMs representing TC activation patterns, protein folding and metabolism (brown cluster), or in
125 innate response and monocyte signatures (purple cluster) (**Supplementary Figure 1c**). These febrile
126 patients were considered clinically to have enteric fever, and were therefore treated as such, however
127 their heterogeneous gene transcription profiles suggest that any one of several different aetiologies
128 may have precipitated hospital presentation. Further evidence to this is that in a recent RCT a higher
129 proportion of culture-negative cases responded to fluoroquinolones rather than a 3rd generation
130 cephalosporin, possibly due to the frequency of murine and scrub typhus in this population, however
131 distinguishing between these infections is currently difficult.

132

133 *Multi-cohort data quality assessment*

134 In order to address the potential over-diagnosis of enteric fever and associated inappropriate
135 antimicrobial use, we next aimed to identify a set of genes whose expression is able to differentiate
136 enteric fever from other common febrile conditions found in tropical settings. We repurposed
137 publically available datasets describing host transcriptional response in two malaria,^{10,16} four
138 tuberculosis,^{8,17} and four dengue cohorts (**Supplementary Table 2**).^{9,18,19} We designed a discovery
139 cohort consisting of control samples from each respective study ($n=220$ community controls or
140 convalescent samples, 'CTRL'), 74 enteric fever ('EF'), 94 blood stage *P. falciparum* ('bsPf'), 67
141 dengue ('DENV') and 54 active pulmonary tuberculosis ('PTB') cases. An independent validation
142 cohort consisted of 109 CTRLs, 50 EF, 19 bsPf, 49 DENV, and 97 PTB samples (**Figure 2**). Finally,

143 a cohort of ‘unknown’ samples was created consisting of febrile culture-negative, febrile suspected
144 EF cases from Nepal (‘sEF’), and samples collected from CHIM study participants who did not elop
145 enteric fever after challenge at day 7 (‘nD7’) and their respective pre-challenge baseline samples
146 (‘D0’) (**Figure 2**). Using Principle Component Analysis (PCA) to assess the variability at the level of
147 gene expression between the cohorts indicated some distinct clustering between cases
148 (**Supplementary Figure 2a**), for each infection whereas no such differences were observed with the
149 comparator CTRL samples (**Supplementary Figure 2b**).

150

151 *Five genes sufficiently distinguish EF from other febrile infections*

152 With these data, we aimed to build a classifier containing a minimum set of genes that could
153 discriminate culture-confirmed enteric fever cases from individuals with other causes of fever (class:
154 ‘Rest’, consisting of CTRLs, DENV, PTB and bsPf) (2-class classification, **Figure 2**) using a Guided
155 Regularized Random Forest (GRRF) algorithm.²⁰ Genes were ranked by frequency of selection in
156 each of 100 iterations, and applying a selection threshold of $\geq 25\%$, we identified a putative diagnostic
157 signature containing *STAT1* (98% of iterations), *SLAMF8* (76%), *PSME2* (39%), *WARS* (37%), and
158 *ALDH1A1* (36%) (**Figure 3a**). With this 5-gene signature we were able to predict which individuals
159 in the validation cohort had enteric fever with a sensitivity and specificity of 97.1% and 88.0%,
160 respectively (AUROC: 96.7%) (**Figure 3b, Supplementary Table 3a**). Of blood culture-confirmed
161 enteric fever cases in the validation cohort, 6/51 were misclassified as ‘Rest’ (i.e. classification
162 probability > 0.5 , **Figure 3c-top**), and 8/274 samples belonging to class ‘Rest’ were classified as
163 enteric fever. These included six tuberculosis and one dengue case, and a pre-challenge baseline
164 sample from a CHIM participant (**Figure 3c-bottom**).

165 To allow comparison between the different disease conditions, we quantified expression of the 5
166 genes identified in each sample using the z-score of the geometric mean of the expression values
167 (expression score). Significant differences in expression scores were observed between the enteric
168 fever samples and all other conditions in both the discovery (top) and the validation (bottom) cohort
169 (**Figure 3d**). Of note, there were no significant differences between the scores calculated for the
170 control samples derived from endemic areas or naïve, healthy controls from the CHIM, indicating the
171 homogeneity of expression to these genes in healthy controls from different study and geographical
172 locations.

173 The design of discovery and validation cohorts is likely to have an impact on the diagnostic signature
174 selected, and we therefore exchanged the validation and discovery cohort and re-ran the analysis.
175 Although in this experiment 4 instead of 5 genes were selected (using a threshold $\geq 25\%$), most genes
176 included were also part of the initial signature (*STAT1*, *SLAMF8*, *WARS*) and the high predictive
177 accuracy was maintained (AUROC: 97.2%) (**Supplementary Figure 3a&b**). These results

178 demonstrate the ability of a small number of genes to accurately predict true EF cases from other
179 febrile illnesses caused by another bacterial pathogen (TB), and of parasitic or viral origin.

180

181 *Multiclass prediction accurately classifies three of five conditions simultaneously*

182 Given the apparent success of small gene expression signatures in classifying two distinct groups, we
183 sought to leverage the overall dataset and the GRRF algorithm to identify a signature that could
184 accurately classify more than two classes simultaneously. We re-analysed the data preserving the
185 original class labels (i.e. CTRL, bsPf, DENV, PTB and EF) and performed the iterative feature
186 selection step using the GRRF algorithm (**Figure 2**–“multiclass classification”). Applying a $\geq 25\%$
187 selection threshold to ranked features identified 7 genes (*RFX7*, *CIQB*, *ANKRD22*, *WARS*, *BATF2*,
188 *STAT1*, and *CIQC*) able to discriminate the classes (**Figure 3e**). Prediction of the validation cohort
189 using this 7-gene signature indicated good sensitivity and specificity for accurately classifying CTRL,
190 bsPf and EF cases, however the identification of DENV and PTB was less accurate (**Figure 3f**,
191 **Supplementary Table 3b**). Analysis of individual gene expression levels in each group indicated that
192 *RFX7* was only upregulated in bsPf samples, while *STAT1*, *WARS* as well as *ANKRD22* and *BATF2*
193 were all strongly upregulated in EF. Expression of these genes in PTB and DENV samples was
194 variable accounting for the lower performance of the signature in these conditions (**Supplementary**
195 **Figure 4a&b**).

196

197 *Prediction of unknown samples*

198 Given the superior performance of the 2-class diagnostic signature, our subsequent analyses focused
199 on using the initial 5-genes identified to ascertain whether enteric fever was the likely true underlying
200 aetiology of suspected febrile, blood culture-negative cases in Nepal (sEF; $n=71$), part of the unknown
201 cohort (**Figure 2**). Included in this cohort were 144 samples originating from the challenge study with
202 known class membership confirming the correct classification of 94.4% of the samples by the GRRF
203 algorithm (**Supplementary Table 4**).

204 Classification of these sEF cases predicted 9/71 (12.6%) febrile, culture-negative patients to be true
205 enteric fever cases and the remaining samples to belong to class ‘Rest’ (**Figure 4a**). Relating the gene
206 expression scores to the predicted class probabilities indicated no clear separation of scores according
207 to the predicted class (**Figure 4b**). Furthermore, comparing the expression score of febrile, culture-
208 negative samples with culture confirmed enteric fever in Nepal showed a marked overlap, indicating
209 that these scores alone are insufficient for 2-class discrimination (**Figure 4c**).

210

211 *Diagnostic validation by qPCR*

212 Finally, to validate the induction of the diagnostic gene signature in blood culture-confirmed enteric
213 fever cases, we performed high-throughput qPCR in samples collected during an independent typhoid
214 CHIM (**Supplementary Table 2**)²¹ and in the Nepali cohort. Transcription of the 5-gene signature
215 was increased at the time of diagnosis in most participants with culture-confirmed enteric fever in
216 both sample sets (**Figure 4d&e**). Two CHIM participants diagnosed with typhoid infection and one
217 patient infected with *S. Paratyphi* in Nepal showed low expression of all genes and a resulting low
218 expression score (**Figure 4e**—black arrows). In contrast, one day 7 sample from a participant not
219 diagnosed with enteric fever demonstrated high expression of the putative diagnostic gene signature
220 (**Figure 4e**—black arrows).

221 As surrogate disease severity markers, temperature showed poor correlation with the expression score
222 in both CHIM and endemic setting culture-confirmed enteric fever cases (**Figure 4f&g**—left). In
223 contrast, C-reactive protein levels (only available for CHIM participants) were significantly
224 associated with the expression score of the 5-gene signature (**Figure 4f&g**—right) thus underlining the
225 relevance of this signature in reflecting the clinical presentation of enteric fever. In the Nepal cohort,
226 gene expression also strongly correlated between the array and qPCR data (**Supplementary Figure**
227 **6**). Overall these results verify the strong expression of the putative diagnostic signatures in samples
228 from patients with acute enteric fever and underline the clinical plausibility through association with
229 disease severity parameters.

230

231 Discussion

232 New approaches to diagnose patients with enteric fever are urgently needed, as currently available
233 methods are antiquated and unreliable. New diagnostic modalities are required, to both improve the
234 immediate management of patients, and to increase the accuracy of disease burden measurements to
235 support targeted vaccine implementation. Here we demonstrate a reproducible host expression
236 signature of 5 genes (*STAT1*, *SLAMF8*, *PSME2*, *WARS*, and *ALDH1A1*) able to discriminate EF cases
237 from other common causes of fever in the tropics with an accuracy of >96%. To our knowledge, this
238 exceeds the performance of all previously described enteric fever diagnostic methods, which often
239 perform less well when assessed using samples collected directly from patients or participants.
240 Moreover, application of high-throughput methods such as functional genomics, to this major health
241 concern,²² underscores the importance and tangible benefits of applying ‘omics-technologies’ to
242 combatting infectious diseases in the most needy populations.²³ While further optimisation work is
243 required, validating the expression of our signature using conventional methods such as qPCR
244 demonstrates feasibility of further development into an affordable diagnostic test for use in endemic
245 settings.²⁴

246 The degree of perturbation of molecular responses occurring during enteric fever can be confounded
247 by the duration of clinical illness (ranging in 12hrs to ≥ 3 days in the CHIM and patients from Nepal
248 , respectively) or the specific pathogen (*S. Typhi* or *S. Paratyphi*). This may hinder identification of a
249 reproducible gene expression signature reliably expressed in various settings. The responses to *S.*
250 *Typhi* and *S. Paratyphi* cases in Nepal were remarkably similar, with the majority of DE genes
251 overlapping between the two groups, which is unsurprising given the close genetic relatedness of both
252 pathogens.²⁵ Enrichment of BTMs resembled responses described previously by us^{13,26} and underlined
253 the concordance between culture-confirmed enteric fever cases from Oxford and Nepal despite the
254 possible differences between challenge and currently circulating strains.

255 Despite the multiple redundancies incorporated into human immune pathways driven by successful
256 evolution,²⁷ our data suggest that the pattern of immune response activation is sufficiently specific to
257 allow identification of the causative pathogen. For example, while immune responses during enteric
258 fever and TB are broadly characterized by IFN-signalling, we and others have reported that this
259 response during acute *S. Typhi* infection appears to be skewed towards a type-II pattern likely
260 associated with neutrophils and NK cells rather than the type-I dominated profile found in TB.^{7,8,13,28-}

261 ³¹ Application of computational methods to large datasets including host gene expression has been
262 shown to be an effective approach to capture such differential activation of immune pathways.^{5,6} Two
263 of the genes identified in our 5-gene diagnostic signature are important entities in the IFN- γ signalling
264 cascade (*STAT1*, *WARS*), which has been broadly implicated in the responses to enteric fever, TB,⁸
265 dengue,³² and *P. falciparum*³³ infection. The discriminatory impact of increased expression of these
266 genes identified in our analysis, however, suggests that there are distinct differences during the

267 responses to these very different pathogens sufficient to discriminate underlying disease aetiology^{34,35}
268 possibly based on subtle metabolic differences.^{13,36} While STAT1 and WARS are markers of an IFN- γ
269 response, SLAMF8 is surface-expressed protein³⁷ found in macrophages, DCs and neutrophils and
270 induced by IFN- γ or Gram-negative bacteria.³⁸ SLAMF8 negatively regulates ROS production
271 through inhibition of NADPH oxidase 2 (NOX2) in the bacterial phagosome and reduces ROS-
272 induced inflammatory cell migration.³⁹ While oxidative stress is a common response to infection,
273 *Salmonella* survival is reduced in SLAMF1-deficient mice and can interfere with localization of
274 functional NOX2 in *Salmonella*-containing vacuoles (SCVs), linking SLAM proteins and oxidative
275 stress.⁴⁰ PSME2 is one of two interferon-inducible subunits of the 20S immunoproteasome (IP)
276 regulator 11S and is involved in immune responses and antigen processing.⁴¹ The 20S IP can be
277 induced by oxidative stress and preferentially hydrolyses non-ubiquitinated proteins.^{42,43} Thus, genes
278 involved in these processes may be exploited to distinguish between pathogens inducing oxidative
279 stress from those also triggering ubiquitination.^{44,45} While ALDH1A1 has not specifically been linked
280 with responses to invasive bacterial infections, it is involved in gut-homing of TCs through expression
281 of retinoic acid,^{46,47} a phenotype we have observed following infection with *S. Typhi*.²⁶ C1QB and
282 C1QC are well-known subunits of the complement subcomponent C1q and, together with ANKRD22
283 (involved in cell cycle control⁴⁸), have previously been described as part of a signature able to
284 distinguish active from latent TB.¹⁷ The function of the transcription factor *RFX7* is largely unknown,
285 but has been found to be strongly up-regulated during blood stage malaria and its selection in our 7-
286 gene signature is therefore likely to be driving the classification of malaria cases.

287 Of note, while multiclass classification is difficult to perform and here merely serves as demonstration
288 that data driven approaches may be capable of performing this task, it is interesting to observe
289 increased misclassification rates specifically in the DENV and TB groups. In the validation cohort,
290 the majority of misclassified DENV cases were identified as enteric fever (5/49) or TB (9/49), and
291 misclassified TB samples as enteric fever (13/97) or DENV (23/97), possibly reflecting the
292 overlapping immune response seen due to the intracellular nature of all three pathogens. In the TB
293 group, 15/97 samples were misclassified as controls, compared with one DENV sample being
294 misclassified as such for example, potentially owing to the broad clinical phenotype or lack of
295 inflammatory/immune responses seen in the peripheral blood during tissue specific pulmonary TB
296 infection.

297 Overall, the genes identified in both signatures through our unbiased selection approach are supported
298 by previous studies including those aiming to develop predictive diagnostic signatures.^{8,17,49} In the era
299 of biological 'big data', several studies have explored the utility of gene transcription signatures
300 capable of discriminating viral aetiologies, viral or bacterial infections as well as acute or latent
301 tuberculosis.^{5,6,17,50-53} Only in the tuberculosis studies have such signatures been identified from
302 samples collected in high-incidence, disease endemic settings and been further validated against other

303 disease processes including (but not limited to) pneumonia, sepsis, and streptococcal and
304 staphylococcal infections.^{6,8,50} Herberg *et al.* demonstrated that distinction between viral and bacterial
305 infections could be achieved based on two genes only.⁵ In contrast most efforts undertaken to
306 diagnose active TB employ biomarker signatures ranging in size from 3-86 genes, possibly due to
307 broad and heterologous molecular responses seen in response to differing clinical phenotypes of
308 infection. In our analysis we specifically focused on pathogens with the potential to cause
309 undifferentiated febrile illnesses in tropical settings. While the clinical presentation and epidemiology
310 of the infections chosen may be sufficient to distinguish the aetiologies clinically, enteric fever has a
311 broad differential diagnosis and is frequently over-diagnosed in the absence of confirmatory
312 laboratory results. Notably, despite the high prediction accuracy of the signatures identified in our
313 analysis, this type of data modelling is highly dependent on the quality and availability of suitable
314 input datasets. Although an increasing amount of data is accumulating in the public domain, few well-
315 defined datasets of samples representing a larger repertoire of febrile illnesses are available. For
316 example, rickettsial infection is likely to underlie a large burden of the culture-negative cases in
317 Nepal, however no gene expression datasets exist and the lack of adequate confirmatory diagnostic
318 tests further hinders the inclusion of such data in our analysis.

319 Although the 5-gene signature achieved high accuracy in identifying enteric fever cases, several
320 culture-confirmed cases were misclassified. Metadata from samples collected in the Oxford CHIM
321 indicate that the majority of these misclassified samples had a temperature below 37°C (5/6) and were
322 diagnosed beyond 7 days after challenge (4/6), which, in our CHIM experience, is likely to indicate a
323 less severe disease phenotype. In contrast, six nD7 samples from the Oxford CHIM (part of the
324 unknown cohort) classified as enteric fever showed some sign of response either based on increased
325 cytokines, temperature or a positive stool culture (data not shown). Because our analysis was purely
326 data driven and not motivated by clinical suspicion, we believe that these observations and the
327 significant association of the gene expression scores with CRP provide sufficient evidence that these
328 study participants had infection despite not meeting our study endpoint definitions for enteric fever.

329 In summary, our work demonstrates how a large gene expression dataset derived from challenge study
330 cohorts and settings endemic for febrile infectious diseases can be exploited for diagnostic biomarker
331 discovery. Verification of the putative diagnostic signature using qPCR in independent validation sets
332 indicates that a diagnostic test derived from these gene expression data could be developed for
333 deployment in resource-limited settings. The application of purely data-driven analyses to large and
334 well-defined host-pathogen datasets derived from disease relevant populations may enable us to
335 develop a single, highly accurate diagnostic signature which would allow rapid identification of the
336 main fever-causing aetiologies from readily available biological specimens.

337 **Online Methods**

338 **Typhoid challenge model**

339 Samples included in the discovery cohort were collected during a typhoid dose-escalation study in
340 which 41 healthy adult volunteers ingested a single dose of *S. Typhi* Quail's strain following pre-
341 treatment with 120 mL sodium bicarbonate solution (Study: T1). In this study, one of two doses were
342 administered: $1-5 \times 10^3$ ($n=21$) and $1-5 \times 10^4$ ($n=20$).⁴ Samples used in the validation cohort were
343 collected from a second typhoid challenge model performed as part of a vaccine efficacy study
344 (Study: T2), in which healthy adult volunteers ingested a single dose of *S. Typhi* Quail's strain ($1-5 \times 10^4$, $n=99$) 4 weeks after oral vaccination with Ty21a, M01ZH09 or placebo.⁵⁴ Lastly, samples
345 collected from the control arm of a further vaccine efficacy challenge study, in which participants
346 received meningococcal ACWY-CRM conjugate vaccine (MENVEO[®], GlaxoSmithKline) prior to
347 challenge, were used for the independent qPCR validation experiment.²¹ The clinical and molecular
348 results of these studies have been described previously.^{4,7,21,54} In all typhoid challenge studies
349 participants were treated with a 2-week course of antibiotics at the time of diagnosis (fever $\geq 38^\circ\text{C}$
350 sustained for ≥ 12 hrs and/or positive blood culture), or at day 14 post-challenge if diagnostic criteria
351 were not reached.
352

353

354 **Paratyphoid challenge model**

355 Clinical samples for paratyphoid infection were collected during a dose-escalation study, as
356 previously described (P1).⁷ Briefly, 40 healthy adult volunteers were challenged with a single oral
357 dose of virulent *S. Paratyphi* A (strain NVGH308) bacteria, which as before, was suspended in 30 mL
358 sodium bicarbonate solution [17.5 mg/mL], and after pre-treatment with 120 mL sodium bicarbonate
359 solution. Oral challenge inocula was given at one of two dose levels, low ($n=20$; median
360 [range]= 0.9×10^3 CFU [0.7×10^3 – 1.3×10^3]) or high dose ($n=20$; median [range]= 2.4×10^3 CFU
361 [2.2×10^3 – 2.8×10^3]). Criterion for diagnosis were either microbiological (≥ 1 positive blood culture
362 collected after day 3) and/or clinical (fever $\geq 38^\circ\text{C}$ sustained for ≥ 12 hrs). Participants were ambulatory
363 and followed up as outpatients at least daily after challenge when safety, clinical, and laboratory
364 measurements were performed.⁷

365

366 **Endemic Cohort**

367 To validate the gene transcriptional signatures in a relevant patient cohort, blood samples were
368 collected from three cohorts at Patan Hospital or the Civil Hospital both located in the Lalitpur Sub-
369 Metropolitan City area of Kathmandu Valley in Nepal. Firstly, blood samples were collected as part
370 of a diagnostics study⁵⁵ from febrile patients presenting to hospital and diagnosed with blood culture-
371 confirmed *S. Typhi* ($n=19$) or *S. Paratyphi* A ($n=12$) infection and febrile patients who were blood

372 culture negative ($n=71$). Samples from a cohort of healthy control volunteers ($n=44$) were also
373 collected as part of this study.

374

375 Gene expression arrays sample processing

376 In the human challenge studies (T1, T2, and P1), peripheral venous blood (3mL) was collected in
377 Tempus™ Blood RNA tubes (Applied Biosystems) before challenge (baseline, pre-challenge controls,
378 ‘D0’, $n=166$) and at paratyphoid diagnosis (‘SPT’, $n=18$) or typhoid diagnosis (‘ST’, $n=75$). In those
379 challenged but who did not develop enteric fever within 14 days of challenge, gene expression was
380 measured at the median day of diagnosis of the diagnosed group in the appropriate studies and this
381 day was termed ‘nD7’ ($n=73$). In Nepal, blood was collected when patients presented to hospital
382 ($n=102$) and from healthy controls ($n=44$) (**Figure 1A, Supplementary Table S1**). Total RNA was
383 extracted from all samples using the Tempus™ Spin RNA Isolation kit (Life Technologies). Where
384 applicable, 50ng of RNA was used for hybridization into Illumina HT-12v4 bead-arrays (Illumina
385 Inc.) at the Wellcome Trust Sanger Institute (Hinxton, UK) or The Wellcome Trust Centre for Human
386 Genetics (Oxford, UK) and fluorescent probe intensities captured with the GenomeStudio software
387 (Illumina Inc.). For the paratyphoid CHIM (P1) RNA gene expression was determined using RNA
388 sequencing. Briefly, libraries were prepared using a poly-A selection step to exclude ribosomal RNA
389 species (read length: 75bp paired-end) and samples were subsequently multiplexed in 95 samples/lane
390 over 10 lanes plus one 5-plex pool run on 1 lane and sequenced using a Illumina HiSeq 200 V4.

391 Data pre-processing

392 Paired-end reads were adapter removed and trimmed from 75 to 65bp using trimmomatic v0.35⁵⁶ and
393 only reads exceeding a mean base quality 5 within all sliding windows of 5bp were mapped to the
394 Gencode v25/hg38 transcriptome using STAR aligner v2.5.2b keeping only multi mapped reads
395 mapping to at most 20 locations. featureCounts from the subread set of tools v1.5.1 was used to
396 quantify reads in Gencode v25 basic gene locations with parameters -C -B -M -s 2 -p -S fr. Between-
397 sample normalization was performed using TMM (Trimmed Mean of M-values) normalization as
398 implemented in the edgeR⁵⁷ package and we used principle component analysis (PCA) as quality
399 control step and excluded 2 samples, which were clear outliers due to also failing QC during the
400 library preparation. Counts were converted into \log_2 counts per million (cpm) values with 0.5 prior
401 counts to avoid taking the logarithm of zero and were then taken forward to the multi-cohort quality
402 control. Illumina HT-12v4 bead array data were pre-processed by background subtraction, quantile
403 normalization and \log_2 -transformation using the limma package in R.⁵⁸ Probes were collapsed to
404 HUGO gene identifiers keeping only the highest expressed probe.

405 Data download

406 Previously published whole blood transcriptional array data was downloaded from the Gene
407 Expression Omnibus (GEO) data repository. In this study we specifically focused on studies

408 investigating blood stage *Plasmodium falciparum* (bsPf; two cohorts of blood-stage, HIV-negative
409 malaria cohorts; children and adults),^{10,16} acute uncomplicated dengue (DENV; four adult South-East
410 Asian cohorts of uncomplicated dengue fever patients),^{9,18,19} and active pulmonary tuberculosis (PTb;
411 four cohorts of active, pulmonary TB HIV-negative adults from Africa and the UK),^{8,17} all infections
412 which present with undifferentiated fever and are relevant to areas where enteric fever is endemic
413 (**Supplementary Table S2**). Raw data were downloaded from GEO using the getGEO-function⁵⁹ and
414 quantile normalization with detection p-values and control probes where available. Probes were
415 collapsed to HUGO gene identifiers keeping only the highest expressed probe.

416 Data processing and cohort Quality Control

417 Probe sequences on microarrays may not correspond to the most recent release of the human reference
418 genome that was used for the RNAseq alignment. In order to mitigate this potential discrepancy we
419 re-annotated the probes to the Gencode v25/hg38. The new annotations were used as gene names for
420 each probe. To avoid uninformative genes and gender bias only probes common to all datasets, not
421 located on sex chromosomes and with an expression above the lowest tertile of the average expression
422 (12,821 probes) were used and a ‘superset’ was created by merging the expression data from all
423 studies into one large data matrix. In order to avoid platform or study related artefacts between the
424 data we applied surrogate variable analysis (sva)⁶⁰ to remove batch effects based on study ID while
425 preserving the disease condition (i.e. control or individual infection).

426 Diagnostic signature identification

427 For classification analyses, we separated the superset into a discovery cohort and a validation cohort.
428 To ensure heterogeneity and optimal feature identification we restricted the discovery cohort to
429 samples solely generated on Illumina platforms and ensured inclusion of EF samples from Oxford and
430 Nepal. In order to establish a validation cohort we casted a wider net and permitted studies generated
431 on other platforms including Affymetrix due to the limited amount of suitable datasets available in the
432 public domain. In addition, to predict unknown samples by applying the signatures identified in this
433 study, we separated the febrile, culture-negative suspected enteric fever cases, samples at day 7 after
434 challenge of those who stayed well and their respective pre-challenge control samples from the
435 superset into a cohort of samples of unknown aetiology (Unknown Cohort) (**Figure 2**).

436 Only the discovery cohort was used for feature selection using Guided Regularized Random Forest
437 (GRRF)²⁰ as implemented in the R package RRF v1.7⁶¹ with gamma = 0.5 and parameter mtry tuning
438 was performed using the tuneRRF command. Feature selection was repeated on 100 iterations of
439 bootstrapped subsets of about 70% of the data in the discovery cohort. To assess model performance,
440 predictions on the held out 30% of the discovery cohort were performed and balanced accuracies⁶²
441 were recorded to account for class imbalances. Genes were then ranked by the frequency of positive
442 gene selection by GRRF (based on mean Gini) during the 100 iterations and only genes included in at

443 least 25% of the selection rounds were included in the diagnostic signature and used for prediction of
444 the independent validation cohort as well as the samples belonging to the unknown cohort (**Figure 2**).

445 High-throughput qPCR validation

446 We performed TaqMan gene expression assays to validate gene expression levels in samples from
447 Nepal and a subset of individuals from the Oxford challenge studies. A panel of 24 probes were
448 measured in triplicates on a 192.24 Fluidigm chip using the Biomark at the Weatherall Institute for
449 Molecular Medicine (WIMM) single cell facility. Four samples and one probe failed in the quality
450 control and were removed from the analysis. Raw Ct values were normalized to the housekeeping
451 gene cyclophilin A (PPIA) ($^{\Delta}Ct$ values) and subsequently to control samples (healthy controls) to
452 achieve $\Delta\Delta Ct$ values.

453 Statistical analysis

454 All data were processed in R version 3.2.4. Comparison of groups in Figure 3d were performed using
455 Student's t-Test and correlations between clinical parameters and expression scores were performed
456 using Pearson correlation and correlation between array and qPCR expression as performed using
457 Spearman correlations (alternative: two-sided).

458 **Data deposition:**

459 The datasets generated in these studies were deposited at GEO: GSE113867.

460

461 **Acknowledgements**

462 We gratefully acknowledge the assistance of the participants who have taken part in the study both in
463 Oxford and Nepal. We are also grateful for the support from Laura B. Martin and GSK Vaccines for
464 Global Health for setting up the paratyphoid challenge model and providing the *S. Paratyphi* O:2
465 antigen and the *S. Paratyphi* challenge strain (NVGH308). We are grateful for the support from
466 Myron M. Levine for providing the *S. Typhi* Quailles strain used in the typhoid CHIM.

467 This work was supported by the Bill and Melinda Gates Foundation (OPP1089317 and OPP1084259);
468 funding for the challenge studies was provided by a Wellcome Trust Strategic Translational Award
469 (grant number 092661 to AJP); the European Vaccine Initiative (ref: PIM); European Commission
470 FP7 grant “Advanced Immunization Technologies” (ADITEC); and the NIHR Oxford Biomedical
471 Research Centre (Clinical Research Fellowship to TCD; oxfordbrc.nihr.ac.uk). For the support of the
472 fluidigm experiment we acknowledge the Weatherall Institute of Molecular Medicine (WIMM) single
473 cell facility, which was supported by the MRC funded Oxford Consortium for Single-cell Biology
474 (MR/M00919X/1), and the Oxford-Wellcome Trust Institutional Strategic Support Fund.

475 TCD, CJ, CJB, CSW, AJP are supported by the NIHR Oxford Biomedical Research Centre (Oxford
476 University Hospitals NHS Trust, Oxford), TCD is an NIHR funded Academic Clinical Lecturer; SB is
477 a Sir Henry Dale Fellow, jointly funded by the Wellcome Trust and the Royal Society
478 (100087/Z/12/Z); AJP is a Jenner Investigator, James Martin Senior Fellow and an NIHR Senior
479 Investigator. The views expressed in this article are those of the author(s) and not necessarily those of
480 the NHS, the NIHR, or the Department of Health.

481

482 **Author Contributions:** The study in Nepal was designed by TCD, AJP and CJB and run by BB, AK,
483 SD, SB and TCD. Challenge studies at Oxford were designed and performed by AJP, TCD, MMG,
484 HD, CJ, CJB. Laboratory work and data generation was performed by CJB, SS, SP, HH, LB, FS, DP
485 and GD. The computational analysis was conceptualized and executed by JM and CJB. The
486 manuscript was conceptualized and written by CJB and TCD. All authors critically reviewed the
487 manuscript.

488 **Competing Interest Statement:** All authors declare not competing interests.

- 490 1 Mogasale, V. *et al.* Burden of typhoid fever in low-income and middle-income countries: a
 491 systematic, literature-based update with risk-factor adjustment. *Lancet Glob Health* **2**, e570-
 492 580, doi:10.1016/S2214-109X(14)70301-8 (2014).
- 493 2 Buckle, G. C., Walker, C. L. & Black, R. E. Typhoid fever and paratyphoid fever: Systematic
 494 review to estimate global morbidity and mortality for 2010. *J Glob Health* **2**, 010401,
 495 doi:10.7189/jogh.02.010401 (2012).
- 496 3 Parry, C. M., Wijedoru, L., Arjyal, A. & Baker, S. The utility of diagnostic tests for enteric fever
 497 in endemic locations. *Expert Rev Anti Infect Ther* **9**, 711-725, doi:10.1586/eri.11.47 (2011).
- 498 4 Waddington, C. S. *et al.* An outpatient, ambulant-design, controlled human infection model
 499 using escalating doses of Salmonella Typhi challenge delivered in sodium bicarbonate
 500 solution. *Clin Infect Dis* **58**, 1230-1240, doi:10.1093/cid/ciu078 (2014).
- 501 5 Herberg, J. A. *et al.* Diagnostic Test Accuracy of a 2-Transcript Host RNA Signature for
 502 Discriminating Bacterial vs Viral Infection in Febrile Children. *JAMA* **316**, 835-845,
 503 doi:10.1001/jama.2016.11236 (2016).
- 504 6 Sweeney, T. E., Braviak, L., Tato, C. M. & Khatry, P. Genome-wide expression for diagnosis of
 505 pulmonary tuberculosis: a multicohort analysis. *Lancet Respir Med* **4**, 213-224,
 506 doi:10.1016/S2213-2600(16)00048-5 (2016).
- 507 7 Dobinson, H. C. *et al.* Evaluation of the clinical and microbiological response to Salmonella
 508 Paratyphi A infection in the first paratyphoid human challenge model. *Clin Infect Dis*,
 509 doi:10.1093/cid/cix042 (2017).
- 510 8 Berry, M. P. *et al.* An interferon-inducible neutrophil-driven blood transcriptional signature in
 511 human tuberculosis. *Nature* **466**, 973-977, doi:10.1038/nature09247 (2010).
- 512 9 Hoang, L. T. *et al.* The early whole-blood transcriptional signature of dengue virus and features
 513 associated with progression to dengue shock syndrome in Vietnamese children and young
 514 adults. *J Virol* **84**, 12982-12994, doi:10.1128/JVI.01224-10 (2010).
- 515 10 Idaghdour, Y. *et al.* Evidence for additive and interaction effects of host genotype and infection
 516 in malaria. *Proc Natl Acad Sci U S A* **109**, 16786-16793, doi:10.1073/pnas.1204945109 (2012).
- 517 11 Naim, A. N., Tolfvenstam, T., Fink, K. & Hibberd, M. L. Genome-wide gene expression analysis
 518 of human whole-blood samples during acute dengue disease and early convalescence.
 519 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE28991> (2016).
- 520 12 Obermeyer, Z. & Emanuel, E. J. Predicting the Future - Big Data, Machine Learning, and Clinical
 521 Medicine. *N Engl J Med* **375**, 1216-1219, doi:10.1056/NEJMp1606181 (2016).
- 522 13 Blohmke, C. J. *et al.* Interferon-driven alterations of the host's amino acid metabolism in the
 523 pathogenesis of typhoid fever. *J Exp Med* **213**, 1061-1077, doi:10.1084/jem.20151025 (2016).
- 524 14 Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for
 525 interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-15550,
 526 doi:10.1073/pnas.0506580102 (2005).
- 527 15 Li, S. *et al.* Molecular signatures of antibody responses derived from a systems biology study
 528 of five human vaccines. *Nat Immunol* **15**, 195-204, doi:10.1038/ni.2789 (2014).
- 529 16 Subramaniam, K. S. *et al.* The T-Cell Inhibitory Molecule Butyrophilin-Like 2 Is Up-regulated in
 530 Mild Plasmodium falciparum Infection and Is Protective During Experimental Cerebral
 531 Malaria. *J Infect Dis* **212**, 1322-1331, doi:10.1093/infdis/jiv217 (2015).
- 532 17 Kaforou, M. *et al.* Detection of tuberculosis in HIV-infected and -uninfected African adults
 533 using whole blood RNA expression signatures: a case-control study. *PLoS Med* **10**, e1001538,
 534 doi:10.1371/journal.pmed.1001538 (2013).
- 535 18 Kwissa, M. *et al.* Dengue virus infection induces expansion of a CD14(+)CD16(+) monocyte
 536 population that stimulates plasmablast differentiation. *Cell Host Microbe* **16**, 115-127,
 537 doi:10.1016/j.chom.2014.06.001 (2014).

538 19 Tolfvenstam, T. *et al.* Characterization of early host responses in adults with dengue disease. *BMC Infect Dis* **11**, 209, doi:10.1186/1471-2334-11-209 (2011).

539

540 20 Deng, H. & Runger, G. Gene selection with guided regularized random forest *Pattern*

541 *Recognition* **46**, 3483-3489 (2013).

542 21 Jin, C. *et al.* Efficacy and immunogenicity of a Vi-tetanus toxoid conjugate vaccine in the

543 prevention of typhoid fever using a controlled human infection model of Salmonella Typhi: a

544 randomised controlled, phase 2b trial. *Lancet*, doi:10.1016/S0140-6736(17)32149-9 (2017).

545 22 Escadafal, C. *et al.* New Biomarkers and Diagnostic Tools for the Management of Fever in Low-

546 and Middle-Income Countries: An Overview of the Challenges. *Diagnostics (Basel)* **7**,

547 doi:10.3390/diagnostics7030044 (2017).

548 23 Baker, S. Genomic medicine has failed the poor. *Nature* **478**, 287, doi:10.1038/478287a

549 (2011).

550 24 Jiang, L. *et al.* Solar thermal polymerase chain reaction for smartphone-assisted molecular

551 diagnostics. *Sci Rep* **4**, 4137, doi:10.1038/srep04137 (2014).

552 25 McClelland, M. *et al.* Comparison of genome degradation in Paratyphi A and Typhi, human-

553 restricted serovars of Salmonella enterica that cause typhoid. *Nat Genet* **36**, 1268-1274,

554 doi:10.1038/ng1470 (2004).

555 26 Salerno-Goncalves, R. *et al.* Challenge of Humans with Wild-type Salmonella enterica Serovar

556 Typhi Elicits Changes in the Activation and Homing Characteristics of Mucosal-Associated

557 Invariant T Cells. *Front Immunol* **8**, 398, doi:10.3389/fimmu.2017.00398 (2017).

558 27 Nish, S. & Medzhitov, R. Host defense pathways: role of redundancy and compensation in

559 infectious disease phenotypes. *Immunity* **34**, 629-636, doi:10.1016/j.immuni.2011.05.009

560 (2011).

561 28 Spees, A. M. *et al.* Neutrophils are a source of gamma interferon during acute Salmonella

562 enterica serovar Typhimurium colitis. *Infect Immun* **82**, 1692-1697, doi:10.1128/IAI.01508-13

563 (2014).

564 29 Thompson, L. J. *et al.* Transcriptional response in the peripheral blood of patients infected

565 with Salmonella enterica serovar Typhi. *Proc Natl Acad Sci U S A* **106**, 22433-22438,

566 doi:10.1073/pnas.0912386106 (2009).

567 30 Blohmke, C. J. *et al.* Induction of Cell Cycle and NK Cell Responses by Live-Attenuated Oral

568 Vaccines against Typhoid Fever. *Front Immunol* **8**, 1276, doi:10.3389/fimmu.2017.01276

569 (2017).

570 31 Manca, C. *et al.* Hypervirulent M. tuberculosis W/Beijing strains upregulate type I IFNs and

571 increase expression of negative regulators of the Jak-Stat pathway. *J Interferon Cytokine Res*

572 **25**, 694-701, doi:10.1089/jir.2005.25.694 (2005).

573 32 De La Cruz Hernandez, S. I. *et al.* A strong interferon response correlates with a milder dengue

574 clinical condition. *J Clin Virol* **60**, 196-199, doi:10.1016/j.jcv.2014.04.002 (2014).

575 33 Miller, J. L., Sack, B. K., Baldwin, M., Vaughan, A. M. & Kappe, S. H. Interferon-mediated innate

576 immune responses against malaria parasite liver stages. *Cell Rep* **7**, 436-447,

577 doi:10.1016/j.celrep.2014.03.018 (2014).

578 34 Munoz-Jordan, J. L., Sanchez-Burgos, G. G., Laurent-Rolle, M. & Garcia-Sastre, A. Inhibition of

579 interferon signaling by dengue virus. *Proc Natl Acad Sci U S A* **100**, 14333-14338,

580 doi:10.1073/pnas.2335168100 (2003).

581 35 Obiero, J. M. *et al.* Impact of malaria preexposure on antiparasite cellular and humoral

582 immune responses after controlled human malaria infection. *Infect Immun* **83**, 2185-2196,

583 doi:10.1128/IAI.03069-14 (2015).

584 36 Zhang, Y. J. *et al.* Tryptophan biosynthesis protects mycobacteria from CD4 T-cell-mediated

585 killing. *Cell* **155**, 1296-1308, doi:10.1016/j.cell.2013.10.045 (2013).

586 37 van Driel, B. J., Liao, G., Engel, P. & Terhorst, C. Responses to Microbial Challenges by SLAMF

587 Receptors. *Front Immunol* **7**, 4, doi:10.3389/fimmu.2016.00004 (2016).

588 38 Wang, G. *et al.* Cutting edge: Slamf8 is a negative regulator of Nox2 activity in macrophages. *J*
589 *Immunol* **188**, 5829-5832, doi:10.4049/jimmunol.1102620 (2012).

590 39 Wang, G. *et al.* Migration of myeloid cells during inflammation is differentially regulated by
591 the cell surface receptors Slamf1 and Slamf8. *PLoS One* **10**, e0121968,
592 doi:10.1371/journal.pone.0121968 (2015).

593 40 Fang, F. C. Antimicrobial actions of reactive oxygen species. *MBio* **2**, doi:10.1128/mBio.00141-
594 11 (2011).

595 41 Nandi, D., Jiang, H. & Monaco, J. J. Identification of MECL-1 (LMP-10) as the third IFN-gamma-
596 inducible proteasome subunit. *J Immunol* **156**, 2361-2364 (1996).

597 42 Seifert, U. *et al.* Immunoproteasomes preserve protein homeostasis upon interferon-induced
598 oxidative stress. *Cell* **142**, 613-624, doi:10.1016/j.cell.2010.07.036 (2010).

599 43 Dubiel, W., Pratt, G., Ferrell, K. & Rechsteiner, M. Purification of an 11 S regulator of the
600 multicatalytic protease. *J Biol Chem* **267**, 22369-22377 (1992).

601 44 Cirillo, S. L. *et al.* Protection of Mycobacterium tuberculosis from reactive oxygen species
602 conferred by the mel2 locus impacts persistence and dissemination. *Infect Immun* **77**, 2557-
603 2567, doi:10.1128/IAI.01481-08 (2009).

604 45 Spooner, R. & Yilmaz, O. The role of reactive-oxygen-species in microbial persistence and
605 inflammation. *Int J Mol Sci* **12**, 334-352, doi:10.3390/ijms12010334 (2011).

606 46 Iwata, M. *et al.* Retinoic acid imprints gut-homing specificity on T cells. *Immunity* **21**, 527-538,
607 doi:10.1016/j.immuni.2004.08.011 (2004).

608 47 Molotkov, A. & Duester, G. Genetic evidence that retinaldehyde dehydrogenase Raldh1
609 (Aldh1a1) functions downstream of alcohol dehydrogenase Adh1 in metabolism of retinol to
610 retinoic acid. *J Biol Chem* **278**, 36085-36090, doi:10.1074/jbc.M303709200 (2003).

611 48 Yin, J. *et al.* ANKRD22 promotes progression of non-small cell lung cancer through
612 transcriptional up-regulation of E2F1. *Sci Rep* **7**, 4430, doi:10.1038/s41598-017-04818-y
613 (2017).

614 49 Zak, D. E. *et al.* A blood RNA signature for tuberculosis disease risk: a prospective cohort study.
615 *Lancet* **387**, 2312-2322, doi:10.1016/S0140-6736(15)01316-1 (2016).

616 50 Anderson, S. T. *et al.* Diagnosis of childhood tuberculosis and host RNA expression in Africa. *N*
617 *Engl J Med* **370**, 1712-1723, doi:10.1056/NEJMoa1303657 (2014).

618 51 Andres-Terre, M. *et al.* Integrated, Multi-cohort Analysis Identifies Conserved Transcriptional
619 Signatures across Multiple Respiratory Viruses. *Immunity* **43**, 1199-1211,
620 doi:10.1016/j.immuni.2015.11.003 (2015).

621 52 Mahajan, P. *et al.* Association of RNA Biosignatures With Bacterial Infections in Febrile Infants
622 Aged 60 Days or Younger. *JAMA* **316**, 846-857, doi:10.1001/jama.2016.9207 (2016).

623 53 Zaas, A. K. *et al.* Gene expression signatures diagnose influenza and other symptomatic
624 respiratory viral infections in humans. *Cell Host Microbe* **6**, 207-217,
625 doi:10.1016/j.chom.2009.07.006 (2009).

626 54 Darton, T. C. *et al.* Using a Human Challenge Model of Infection to Measure Vaccine Efficacy:
627 A Randomised, Controlled Trial Comparing the Typhoid Vaccines M01ZH09 with Placebo and
628 Ty21a. *PLoS Negl Trop Dis* **10**, e0004926, doi:10.1371/journal.pntd.0004926 (2016).

629 55 Darton, T. C. *et al.* Identification of Novel Serodiagnostic Signatures of Typhoid Fever Using a
630 Salmonella Proteome Array. *Front Microbiol* **8**, 1794, doi:10.3389/fmicb.2017.01794 (2017).

631 56 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
632 data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170 (2014).

633 57 Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential
634 expression analysis of digital gene expression data. *Bioinformatics* **26**, 139-140,
635 doi:10.1093/bioinformatics/btp616 (2010).

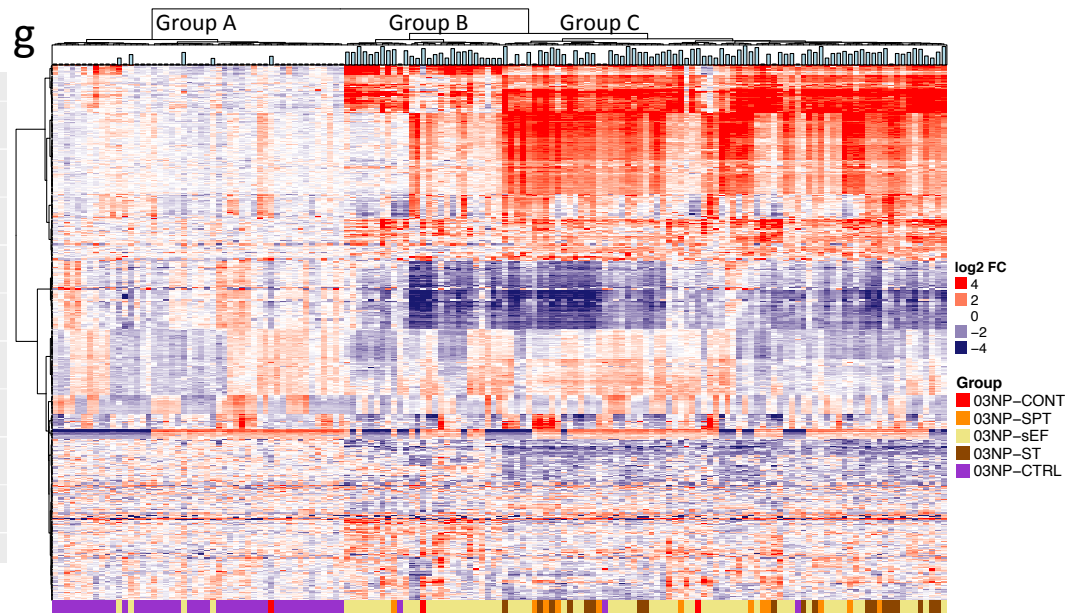
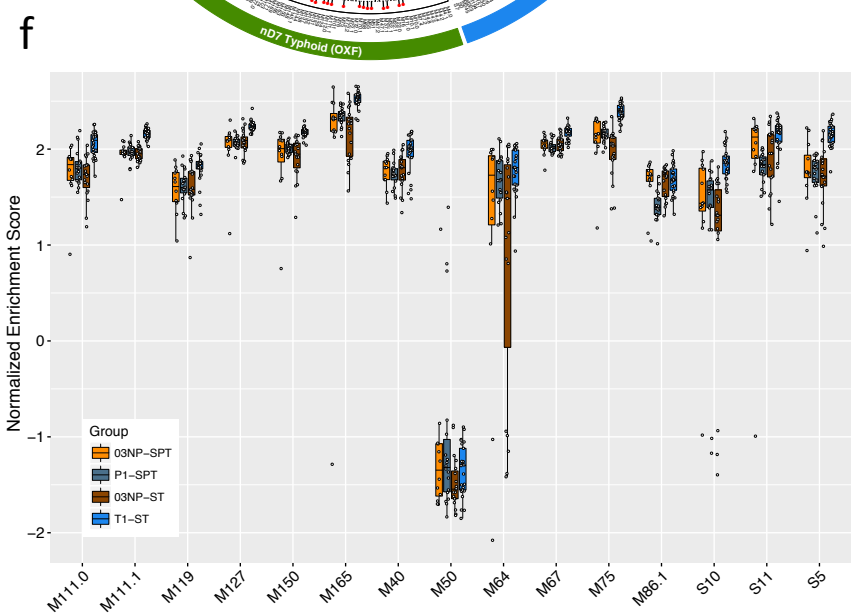
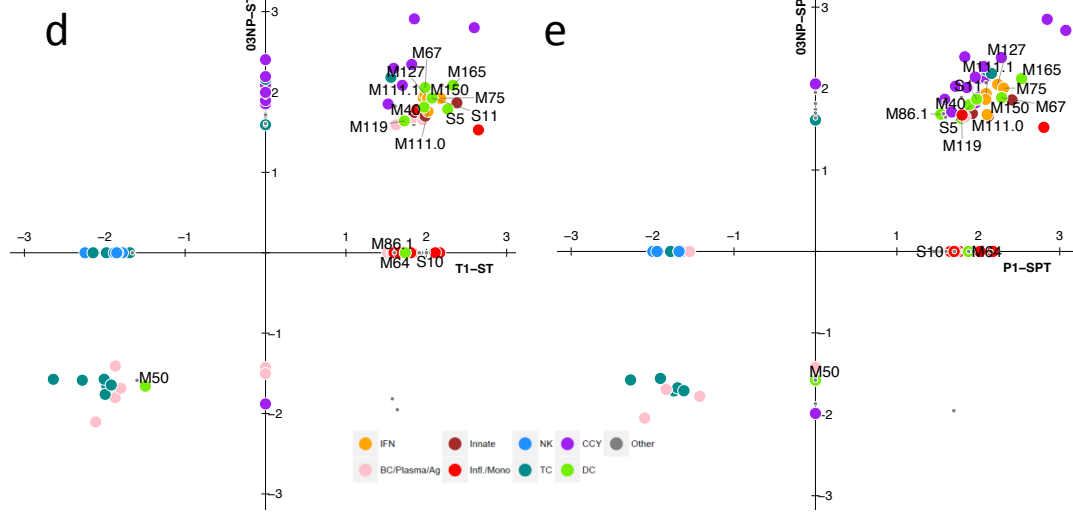
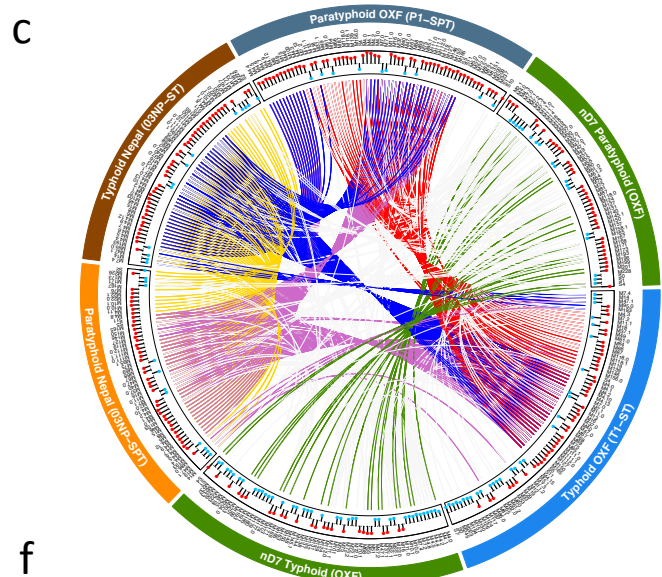
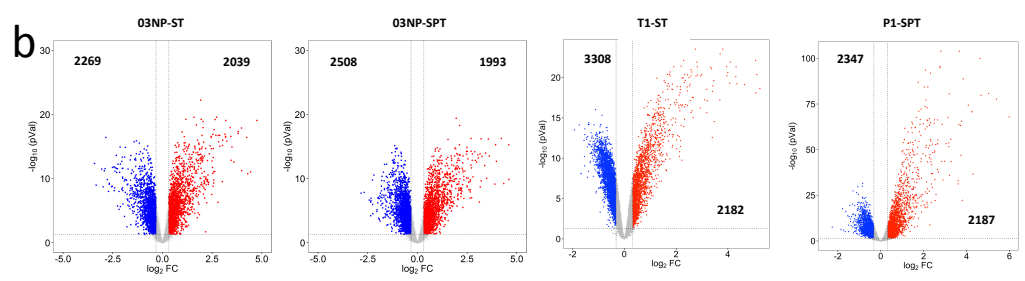
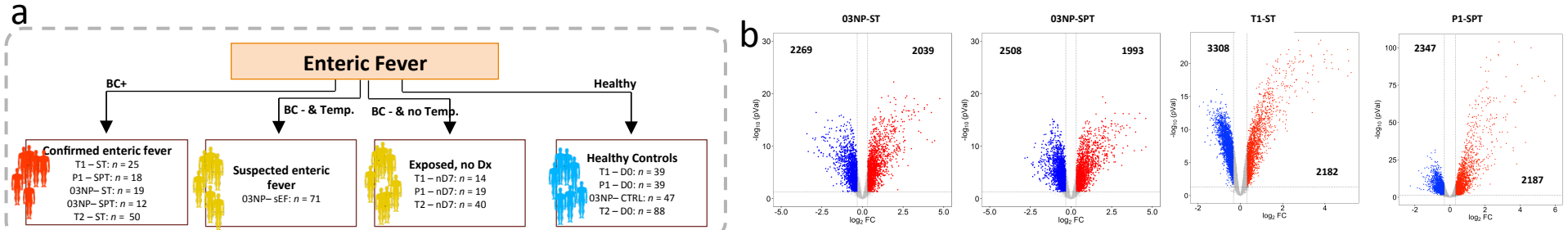
636 58 Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and
637 microarray studies. *Nucleic Acids Res* **43**, e47, doi:10.1093/nar/gkv007 (2015).

638 59 Davis, S. & Meltzer, P. S. GEOquery: a bridge between the Gene Expression Omnibus (GEO)
639 and BioConductor. *Bioinformatics* **23**, 1846-1847, doi:10.1093/bioinformatics/btm254 (2007).
640 60 Leek, J. T. *et al.* sva: Surrogate Variable Analysis. *R Package Version 3.24.4* (2017).
641 61 Deng, H. Guided Random Forest in the RRF Package. *arXiv:1306.0237* (2013).
642 62 Brodersen, K. H., Ong, C. S., Stephan, K. E. & Buhmann, J. M. The balanced accuracy and its
643 posterior distribution. . *Pattern Recognition (ICPR), 2010 20th International Conference on,*
644 *3121-3124*, doi:10.1109/ICPR.2010.764 (2010).

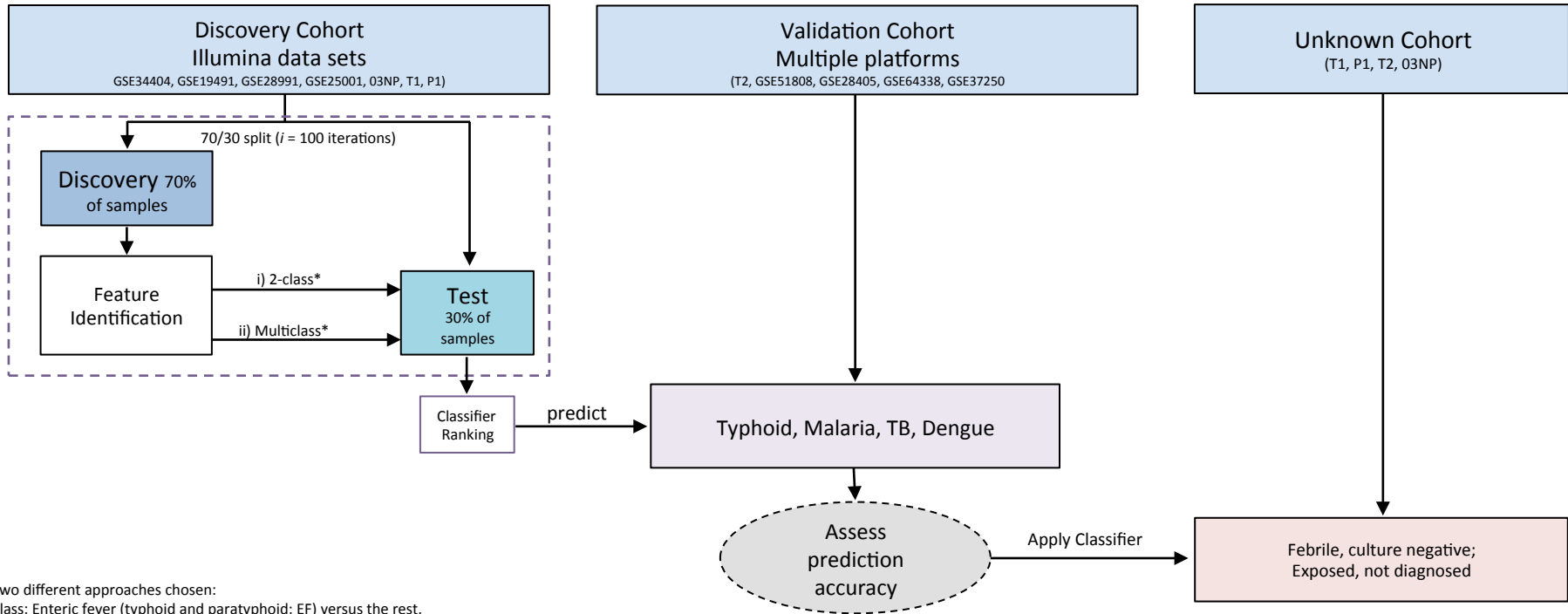
645

646 **Figure Legends and Figures:**

647 **Figure 1: Overview of Oxford and Nepal comparison.** (a) Overview of enteric fever cohorts used
648 in this study (T1: Typhoid CHIM study 1; T2: Typhoid CHIM study 2; P1: Paratyphoid CHIM; 03NP:
649 Nepali cohort. ST: *S. Typhi*; SPT: *S. Paratyphi*; sEF: suspected Enteric Fever; D0: day of challenge
650 which represents the control samples in the Oxford CHIM; CTRL: endemic community controls;
651 nD7: day 7 after challenge in participants who stayed well in the CHIM; BC+: Blood-culture positive;
652 BC-: Blood-culture negative; Dx: Diagnosis). (b) Volcano plots of up (red) and down (blue) regulated
653 genes in *S. Typhi* and *S. Paratyphi* positive individuals (Nepal and Oxford). Black numbers indicate
654 the up- and down-regulated genes. (c) Circular plot depicting the overlap of BTMs between enteric
655 fever and nD7 samples from Oxford and Nepal. Tracks (from outer to inner): cohort and samples;
656 BTM labels; direction of enrichment (blue: down; red: up). Cords represent overlap of enrichment in
657 given cohorts (red: overlap between P1-SPT and T1-ST; green: overlap between T1-nD7 and P1-nD7;
658 blue: overlap of 03NP-ST with P1-SPT and T1-ST; purple: overlap of 03NP-SPT with P1-SPT and
659 T1-ST; yellow: overlap between 03NP-SPT and 03NP-ST). (d-e) Scatter plots of BTMs enriched
660 ($p > 0.05$) in blood-culture positive samples in Nepal (y-axis) versus Oxford (x-axis) for typhoid fever
661 (d) and paratyphoid fever (e). (f) Single-sample GSEA Normalised Enrichment Scores (NES) of IFN
662 and DC BTMs of individuals with blood-culture confirmed enteric fever in Nepal and Oxford. (g)
663 Heatmap of the 500 most variably expressed genes in samples of the Nepali cohort. Bar graph on top
664 of the heatmap shows temperature of each individual at the time of sampling.



665 **Figure 2: Flow diagram of machine learning analysis.** The discover cohort consisted of only
666 Illumina datasets and was used for feature selection using the GRRF algorithm. For the validation
667 cohort Affymetrix datasets were also included. A cohort of unknown samples consisted of pre-
668 challenge baseline samples of participants who stayed well following challenge, their respective nD7
669 samples (7 days after challenge), and febrile, culture-negative suspected enteric fever (sEF) cases
670 from Nepal. Refer to Supplementary Table S2 for study identifiers. 03NP: Nepali cohort. T1: Oxford
671 typhoid CHIM study 1. T2: Oxford typhoid CHIM study 2; P1: Oxford paratyphoid CHIM.

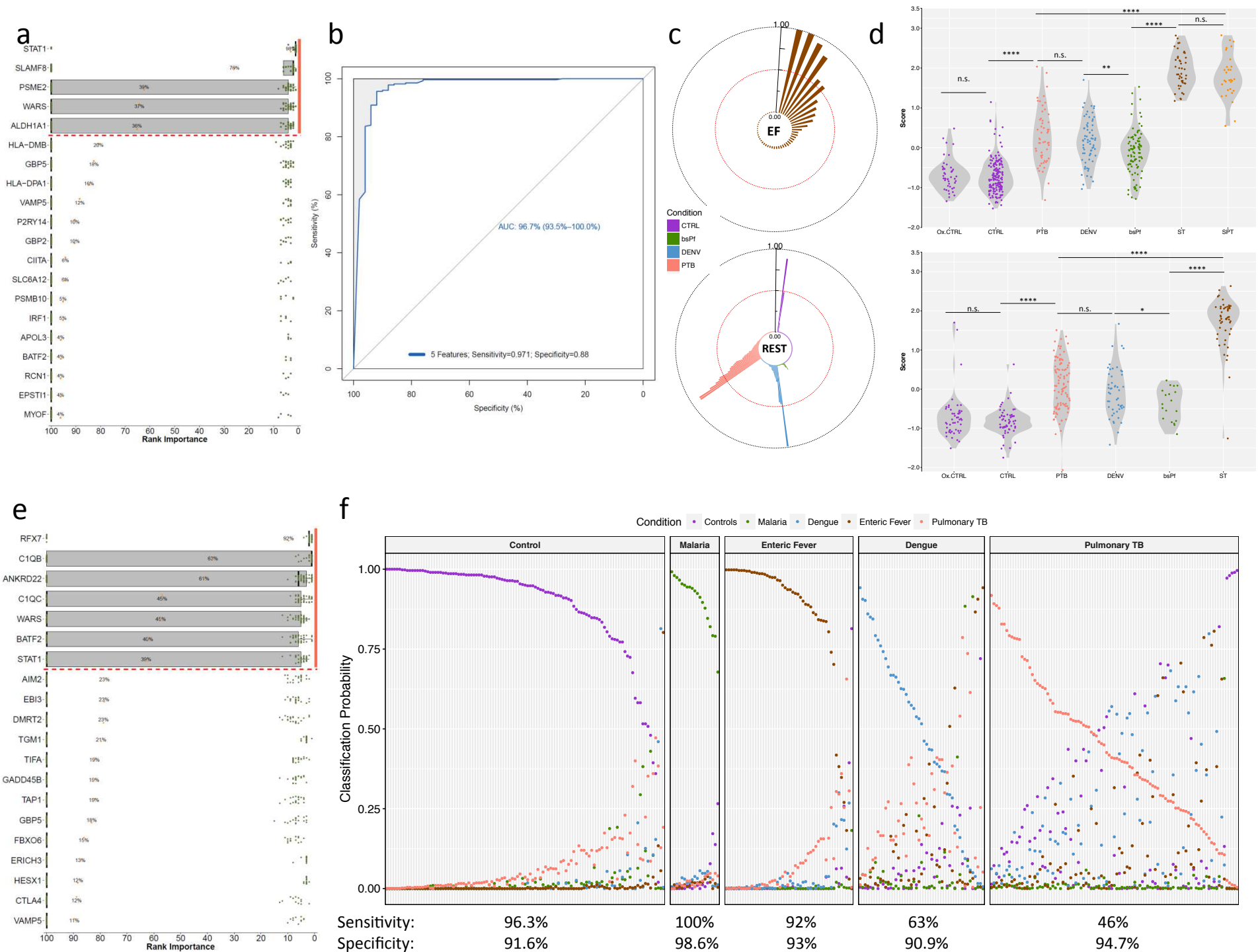


* Two different approaches chosen:

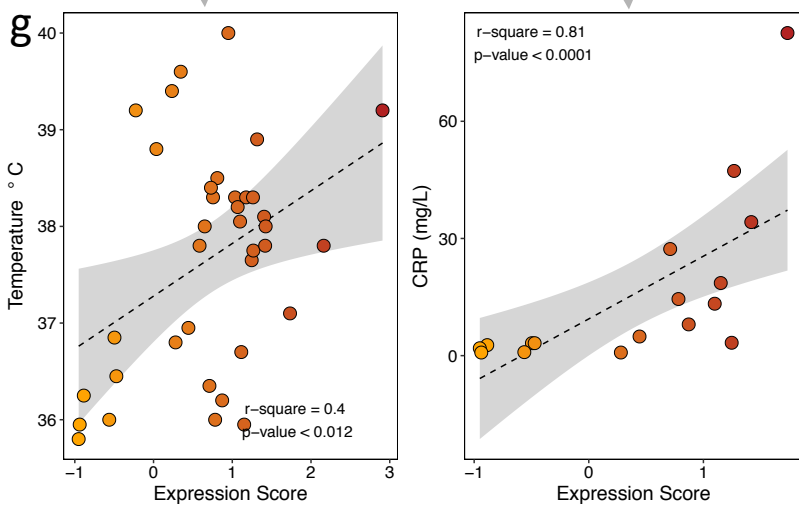
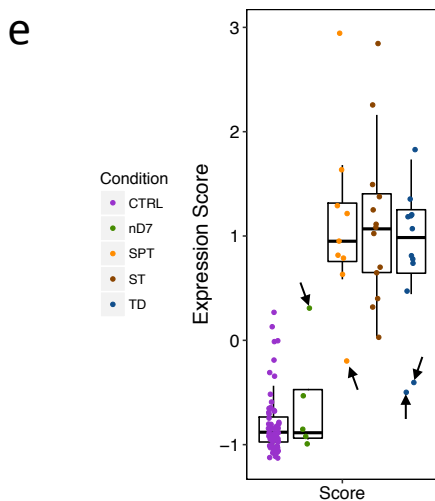
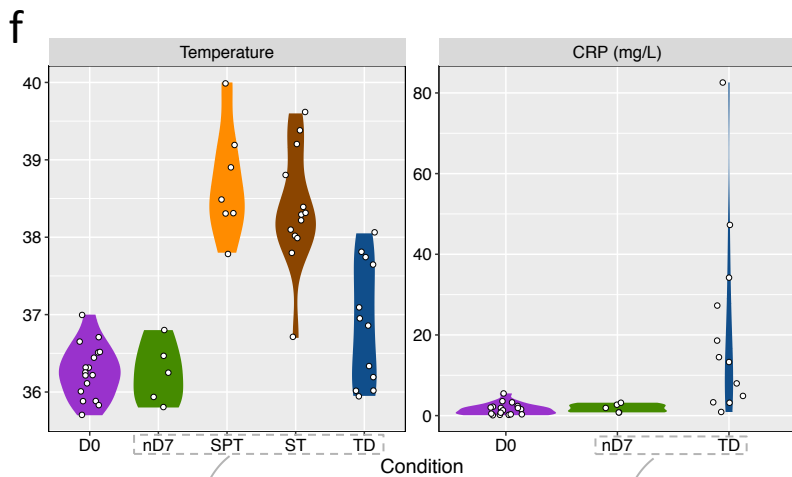
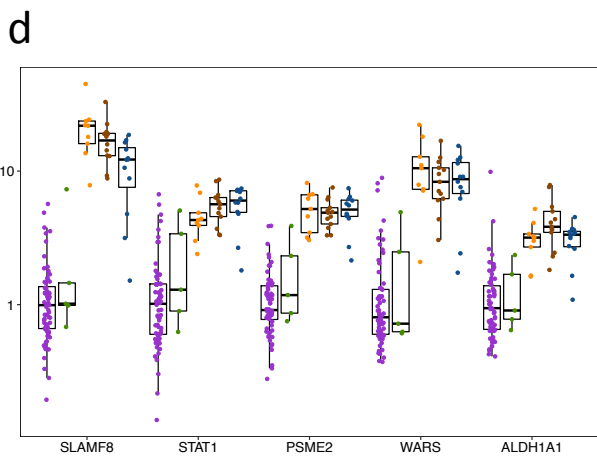
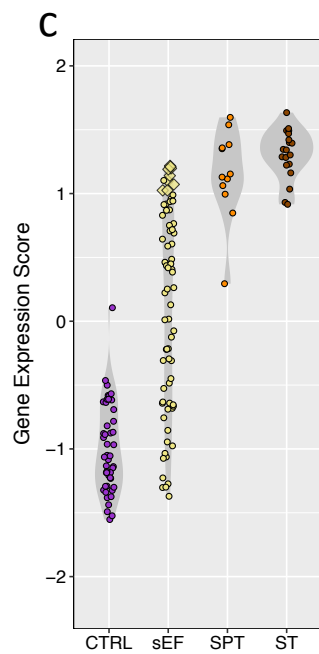
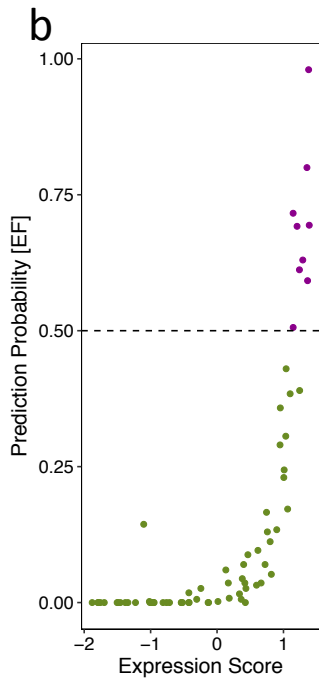
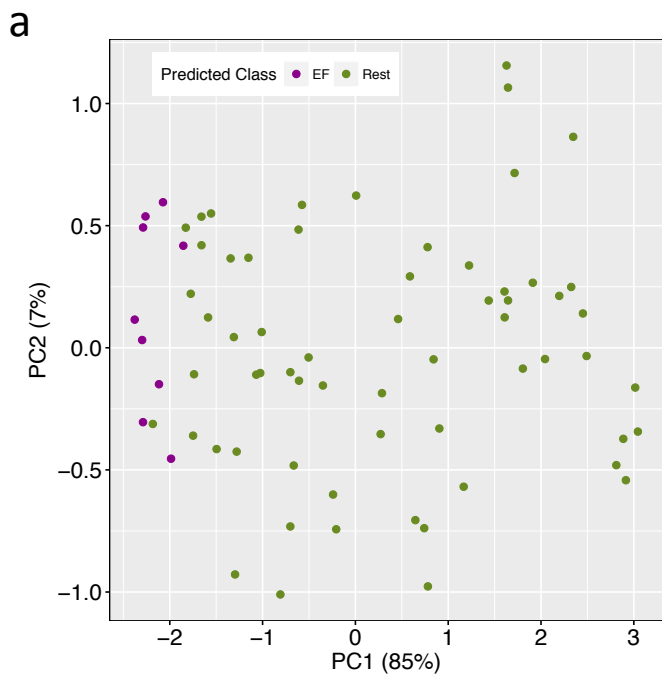
2-class: Enteric fever (typhoid and paratyphoid; EF) versus the rest.

Multiclass: Simultaneously differentiation between all classes (EF, bsPf, DENV, PTB, CTRL).

672 **Figure 3: Identification of diagnostic signatures.** (a) Ranking of genes by their selection frequency
673 into the diagnostic signature out of 100 iterations during the 2-class classification. A cut-off of 25%
674 was selected to detect 5-gene putative diagnostic signature (orange bar). (b) Performance of the 5-
675 gene classifier when predicting the class membership of the validation cohort. (c) Top: Probability of
676 an EF sample to be classified as non-EF (>0.5). Bottom: Probability of sample belonging to 'rest' to
677 be classified as EF (>0.5). (d) Combined expression score for samples based on the 5-gene signature
678 for samples in the discovery cohort (top) and validation cohort (bottom). Ox CTRL: Oxford controls
679 (D0); CTRL: Nepali control samples. PTB: pulmonary TB; DENV: Dengue samples; bsPf: blood-
680 stage *P. falciparum*; SPT: *S. Paratyphi*; ST: *S. Typhi*. (e) Ranking of genes by their selection
681 frequency into the diagnostic signature out of 100 iterations during the multiclass classification. A
682 cut-off of 25% was selected to detect a 7-gene putative diagnostic signature (orange bar). (f)
683 Classification probabilities for each sample of the validation cohort based on the 7-gene signature.
684 Significance levels in panel d were determined using the Student's *t*-Test: * $p < 0.05$; ** $p < 0.01$;
685 *** $p < 0.0001$.



686 **Figure 4: Prediction of Nepali unknown samples using the 2-class and qPCR validation.** (a) PCA
687 of sEF samples based on the 5-gene signature coloured by predicted class membership (EF: purple;
688 green: rest). (b) Dot plot of prediction probability of being class EF versus the expression score
689 calculated on the bases of the 5-gene signature. (c) qPCR gene expression scores of the 5-gene
690 signature ($\Delta\Delta\text{CT}$ over PPIA) for CTRLs, sEF, SPT and ST samples from Nepal. Yellow diamonds in
691 the sEF category represent the 9 patients classified as EF based on the RF algorithm. (d) qPCR
692 expression values ($\Delta\Delta\text{Ct}$ over PPIA) of the 5-gene signature in control samples (Oxford and Nepal),
693 samples at day 7 after challenge of participants who stayed well following challenge with *S. Typhi*
694 (nD7), *S. Paratyphi* (SPT) or *S. Typhi* (ST) in Nepal, or typhoid diagnosis after challenge (TD).
695 Colour legend in panel (e). (e) Combined qPCR expression score of the 5-gene signature. Black
696 arrows indicate outlier samples. (f) Temperature and CRP for samples of which data was available
697 (CRP was only measured in the Oxford CHIM). (g) Spearman's rank correlation of the 5-gene
698 combined expression score and temperature (left; only nD7 and TD samples from the Oxford CHIM
699 and SPT and ST cases from Nepal were included) and CRP (right; CRP was only available for Oxford
700 CHIM samples and we excluded D0 baseline measures) at presentation to hospital (Nepal), diagnosis
701 (Oxford CHIM) or day 7 after challenge in those who stayed well (Oxford CHIM).



702 **Supplementary Figure Legends**

703 **Supplementary Figure 1: Differentially expressed genes and BTMs of sEF cases from Nepal. (a)**

704 Volcano plot of differentially expressed genes in the Nepali sEF cohort. **(b)** Venn diagram
705 representing the overlap of DE genes in the Nepali ST, SPT and sEF cases. **(c)** ssGSEA heatmap of
706 BTMs significantly expressed ($p < 0.05$) in at least 60% of sEF samples. Bar plot panel represents the
707 interdecile range (IDR) for each BTM across all sEF cases. NES: Normalized Enrichment Score.

708 **Supplementary Figure 2: Superset quality control. (a)** PCA plot based on the 500 most variable

709 genes (IQR) of enteric fever cases (EF), malaria cases (bsPf), dengue cases (DENV) and TB cases
710 (PTB) after batch correction. **(b)** PCA of all control samples for each disease cohort after batch
711 correction.

712 **Supplementary Figure 3: Signature identification using a re-designed discovery and validation**

713 **cohort. (a)** Ranking of genes by their selection frequency into the diagnostic signature out of 100
714 iterations during the 2-class classification. A cut-off of 25% was chosen to detect a putative diagnostic
715 signature consisting of 4 genes (orange bar). **(b)** Prediction of the validation cohort using the 4 genes
716 identified in (a).

717 **Supplementary Figure 4:** Expression of the 7 target genes identified during the multiclass

718 classification analysis in each sample of the discovery **(a)** and validation cohort **(b)**.

719 **Supplementary Figure 5: Prediction of Oxford CHIM samples part of the unknown cohort. (a)**

720 PCA of Oxford pre-challenge baseline samples and nD7 samples based on the expression values for
721 the 5-gene diagnostic signature (2-class classification) coloured by predicted class membership
722 (green: REST, purple: EF). **(b)** Dot plot of prediction probabilities against a combined expression
723 score for each sample coloured by predicted class membership.

724 **Supplementary Figure 6:** Spearman correlation of expression values of the 5-gene diagnostic

725 signature derived from microarrays or qPCR.

726

727 **Supplementary Tables:**

728 **Supplementary Table 1:** Overlap of BTMs between different study groups (in percent).

729 **Supplementary Table 2:** Datasets included in this study.

730 **Supplementary Table 3. (a)** Contingency table of class membership following the 2-class

731 classification. **(b)** Contingency table of class membership following the multiclass classification.

732 **Supplementary Table 4:** Prediction accuracy and overview of misclassified samples following
733 prediction using the 5-gene 2-class signature of the Oxford samples included in the unknown cohort.

734 **Supplementary Table 5:** Class memberships of Oxford CHIM and Nepali samples included in the
735 unknown cohort following the prediction using the 7-gene multiclass signature.

736

737