1    **A new paradigm for leprosy diagnosis based on host gene expression**

2    **Insights from leprosy lesions transcriptomics**

3    Thyago Leal-Calvo[1], Charlotte Avanzi[2,#a], Mayara Abud Mendes[1], Andrej Benjak[2,#b],

4    Philippe Busso[2], Roberta Olmo Pinheiro[1], Euzenir Nunes Sarno[1], Stewart T. Cole[2,3],

5    Milton O. Moraes[1*]

6    **Affiliations**

7    [1] Laboratório de Hanseníase, Instituto Oswaldo Cruz, FIOCRUZ, Rio de Janeiro, Rio

8    de Janeiro, Brazil

9    [2] Global Health Institute, École Polytechnique Fédérale de Lausanne, Lausanne,

10   Switzerland

11   [3] Institut Pasteur, Paris, France

12   [#a]Current address: Department of Microbiology, Immunology and Pathology,

13   Mycobacteria Research Laboratories, Colorado State University, Fort Collins,

14   Colorado, United States of America

15   [#b]Current address: Department for BioMedical Research, Oncogenomics Laboratory,

16   University of Bern, Bern, Switzerland

17   * Corresponding author

18   E-mail: milton.moraes@fiocruz.br (MOM)

19

20

21

## Abstract

Transcriptional profiling is a powerful tool to investigate and detect human diseases. In this study, we used bulk RNA-sequencing (RNA-Seq) to compare the transcriptomes in skin lesions of leprosy patients or controls affected by other dermal conditions such as granuloma annulare, a confounder for paucibacillary leprosy. We identified five genes capable of accurately distinguishing multibacillary and paucibacillary leprosy from other skin conditions. Indoleamine 2,3-dioxygenase 1 (*IDO1*) expression alone was highly discriminatory, followed by *TLR10*, *BLK*, *CD38*, and *SLAMF7*, whereas the *HS3ST2* and *CD40LG* mRNA separated multi- and paucibacillary leprosy. Finally, from the main differentially expressed genes (DEG) and enriched pathways, we conclude that paucibacillary disease is characterized by epithelioid transformation and granuloma formation, with an exacerbated cellular immune response, while multibacillary leprosy features epithelial-mesenchymal transition with phagocytic and lipid biogenesis patterns in the skin. These findings will help catalyze the development of better diagnostic tools and potential host-based therapeutic interventions. Finally, our data may help elucidate host-pathogen interplay driving disease clinical manifestations.

## Author Summary

Despite effective treatment, leprosy is still a significant public health issue in more than 120 countries, with more than 200 000 new cases yearly. The disease is caused mainly by *Mycobacterium leprae*, a slow-growing bacillus still uncultivable in axenic media. This limitation has hampered basic research into host-pathogen

2

44    interaction and the development of new diagnostic assays. Currently, leprosy is

45    diagnosed clinically, with no standalone diagnostic assay accurate enough for all

46    clinical forms. Here, we use RNA-seq transcriptome profiling in leprosy lesions and

47    granuloma annulare to identify mRNA biomarkers with potential diagnostic

48    applications. Also, we explored new pathways that can be useful in further

49    understanding the host-pathogen interaction and how the bacteria bypass host

50    immune defenses. We found that *IDO1*, a gene involved with tryptophan catabolism,

51    is an excellent candidate for distinguishing leprosy lesions from other dermatoses.

52    Additionally, we observed that a previous signature of keratinocyte development and

53    cornification negatively correlates with epithelial-mesenchymal transition genes in the

54    skin, suggesting new ways in which the pathogen may subvert its host to survive and

55    spread throughout the body. Our study identifies new mRNA biomarkers that can

56    improve leprosy diagnostics and describe new insights about host-pathogen

57    interactions in human skin.

## Introduction

58

59 Leprosy is a chronic infectious disease caused mainly by the slow-growing

60 intracellular pathogen *Mycobacterium leprae* that does not grow in axenic media. This

61 bacterium resides preferentially in skin macrophages and Schwann cells in peripheral

62 nerves, inducing dermatosis and/or neuritis. Patients can present several distinct

63 clinical forms according to their immune response, histopathological characterization,

64 and bacterial load. A localized tuberculoid form (TT) is characterized by low bacterial

65 counts and a strong cellular immune response. Conversely, in the opposite

66 lepromatous (LL) pole, a disseminated form, patients exhibit several lesions, a

67 predominantly humoral response, and a high bacterial load in the tissues [1–3].

68 Borderline forms are classified according to their proximity to the poles. For operational

69 and treatment purposes, leprosy is classified by the World Health Organization as

70 paucibacillary (PB) or multibacillary (MB), based on the number of skin lesions,

71 associated with nerve involvement or the bacilli detection in slit-skin smears [4].

72 Early and precise diagnosis is instrumental to leprosy control since delay in

73 diagnosis leads to late multidrug therapy, higher disability risk, and continuing

74 transmission, as highlighted by the 200,000 new cases consistently reported annually

75 in the last 10 years [4,5]. However, bacteriological, immunological, genetics or

76 molecular methods are not sufficient for specific diagnosis when used alone.

77 Diagnosis most commonly relies on clinical evaluation, occasionally complemented

78 with histopathological examination and bacterial counts, but these procedures are

79 mostly performed in national reference centers [4,6].

80   Efforts have been deployed to improve leprosy diagnostics using cutting-edge

81  technologies, such as molecular identification of *M. leprae*, serological tests for

82  specific bacterial antigens, and quantification of host biomarkers in plasma or *in vitro*

83  whole blood assays (WBA) [7–9]. Overall, all methods outperform standard clinical

84  diagnosis and can compensate for the low accuracy in detecting PB patients

85  [4,7,8,10–14]. Yet, until now such investigations involved comparing confirmed leprosy

86  cases against healthy endemic controls, who are not representative of individuals with

87  suspected leprosy. Here, other skin conditions represent a better comparator.

88   Identification of markers for early infection is hindered by our poor

89  understanding of pathogenicity and the mechanism by which patients develop one or

90  the other form of leprosy, and nerve injuries [15]. Gene expression signatures have

91  been used as diagnostic tools for several illnesses, from infectious [10–12,14] and

92  autoimmune diseases [16,17] to cancer [18–20]. Some signatures have already been

93  approved for clinical use [12,21–23]. In leprosy, findings from past studies indicate the

94  great potential of expression profiling for disease diagnosis [24–27]. Nonetheless, they

95  were limited by the number of patients [28], or lacked proper epidemiological controls,

96  such as differential diagnosis groups.

97   Here, we applied a combination of bulk RNA sequencing and quantitative

98  validation by RT-qPCR on RNA extracted from skin biopsies of various leprosy forms

99  and from non-leprosy patients to define a specific leprosy host signature applicable to

100  diagnosis. Then, we explored gene expression patterns to improve our understanding

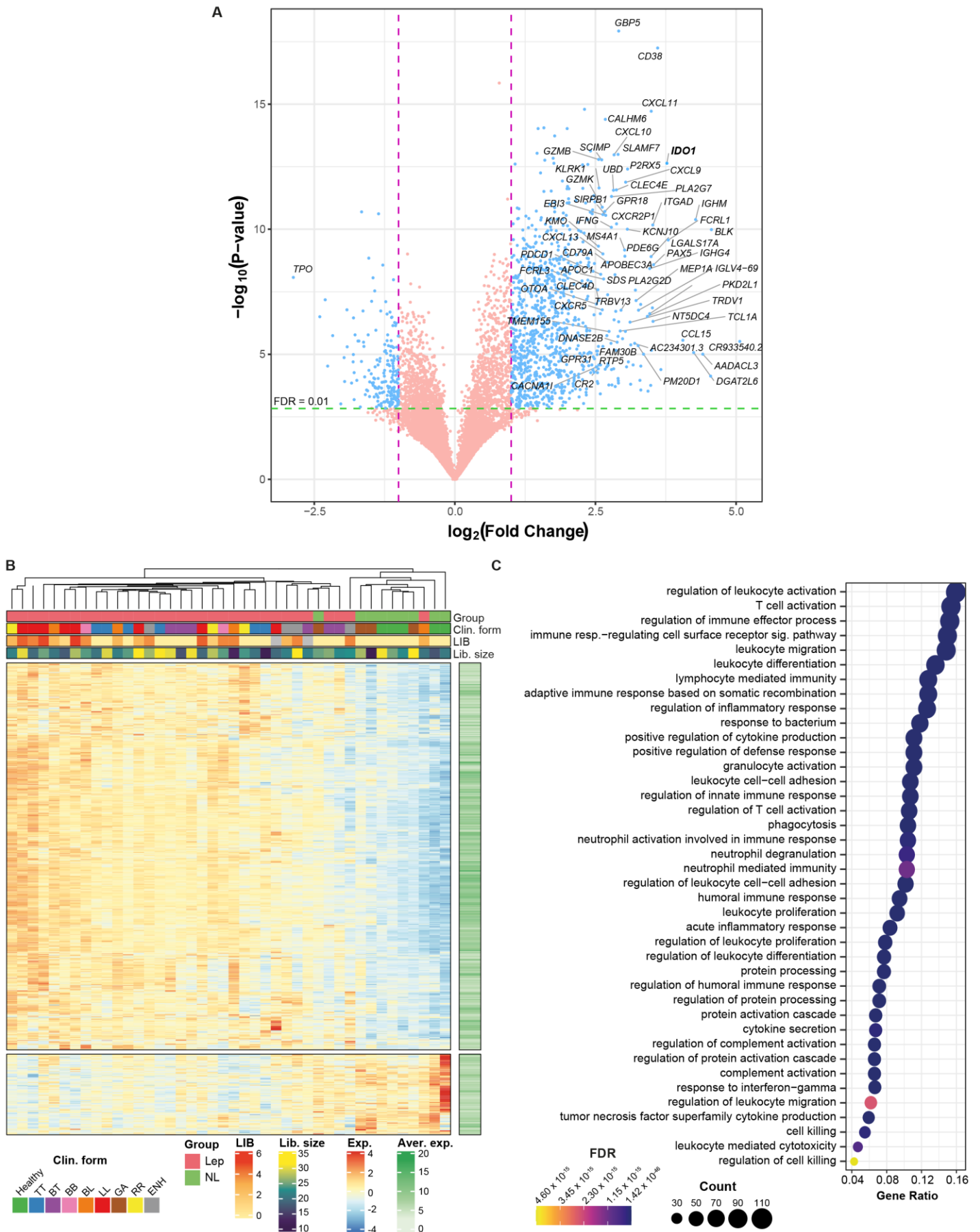101  of the immunopathogenic mechanisms towards leprosy polarization.

# Results

## Discrimination of leprosy *vs.* non-leprosy lesions based on mRNA expression

RNA sequencing was used for pinpointing host candidate genes capable of differentiating leprosy lesions from one of the commonest differential diagnoses of leprosy, granuloma annulare (GA), and from healthy skin. RNA from skin lesions of all leprosy clinical forms (n=33), plus GA (n=4) and healthy skin (n=5) were sequenced (S1 Table). Differentially expressed genes (DEG) in leprosy *vs.* non-leprosy (GA + healthy skin) samples resulted in 1160 DEG with a $|log_2FC| \geq 1$ and FDR $\leq 0.01$, with 961 upregulated in leprosy forms compared to non-leprosy (Fig 1A-B and S2 Table). Exploratory hierarchical clustering of the DEG with $|log_2FC| \geq 1$ and FDR $< 0.01$ grouped all patients' samples into roughly two clusters, except for two: one BL leprosy and one GA that clustered apart from samples with the same diagnosis (Fig 1C). Gene Ontology enrichment analysis of up-regulated genes in leprosy compared to non-leprosy showed enrichment for biological processes associated with leukocyte activation, T-cell activation, immune response, response to the bacterium, neutrophil degranulation, cell killing, cytokine secretion, purinergic receptor signaling pathway, and regulation of defense response to viruses by the host (Fig 1D and S3 Table).
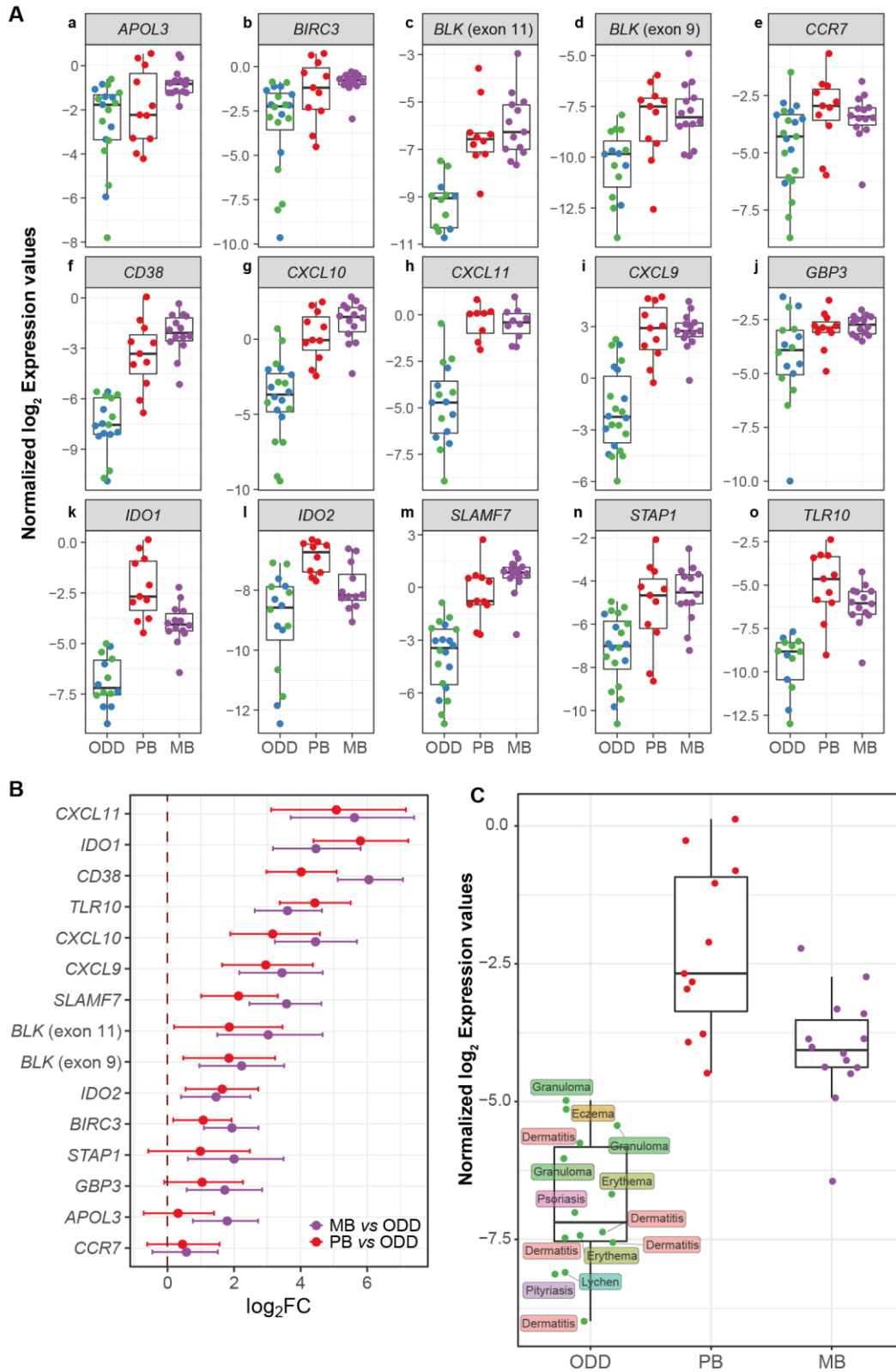
123 **Fig 1. Differentially expressed genes from RNA-seq in leprosy *vs.* GA and**

124    **leprosy *vs.* non-leprosy.** (A) Volcano plot depicting DEG from leprosy *vs.* non-

125    leprosy, where violet dashed line marks $|log_2FC| = 1$. For clarity, gene symbols are

126    shown only for the largest $log_2FC$. (B) Heatmap with hierarchical clustering of samples

127    based on expression of the DEG from leprosy *vs.* non-leprosy comparison. Color scale

128    ranges from lower expression (blue) to higher expression (red). Library size is given

129    in millions. LIB, logarithmic index of bacilli. (C) Biological processes from GO enriched

130    for up-regulated DEG from leprosy *vs.* non-leprosy comparison. FDR, false discovery

131    rate; NL, non-leprosy; GA, granuloma annulare; non-leprosy: GA + healthy individuals.

132

133         A total of 15 genes with the largest effect size ($|log_2FC| \geq 1.5$, FDR < 0.001),

134    highest area under the curve (AUC), and plausible involvement with leprosy

135    pathogenesis (S4 Table) were then validated using a two-step RT-qPCR with a new,

136    larger, and more heterogeneous dataset including skin lesion samples from leprosy

137    patients (n=25), and other common dermatoses (n=23) (S1 Table). Other

138    dermatological diseases (ODD) included dermatitis (n=7), eczema (n=1), erythema

139    (n=4), GA (n=6), lichen planus (n=2), psoriasis (n=2) and pityriasis alba (n=1) (S1

140    Table). A total of 12 samples per group was estimated to be sufficient to attain a power

141    of 85% based on the Welch t-test (PB *vs.* ODD, MB *vs.* ODD) with alpha set at 0.03

142    to replicate the standardized effect size ($log_2FC/SD$) estimated from RNA sequencing.

143    Relative expression using the new sample set by RT-qPCR is shown in Fig 2A. Indeed,

144    the validation data are in agreement with RNA sequencing, because 11 tested genes

145    were replicated by RT-qPCR in terms of difference between mean expression (effect

146    size in $log_2FC$), except for *STAP1*, *GBP3*, *APOL3* and *CCR7* in PB *vs.* ODD

147    comparison and *CCR7* in MB *vs.* ODD (Fig 2B-C, S5 Table). As for differentiating

8

148      leprosy *per se vs.* ODD, genes *IDO1, BLK* (exon 11)*, CD38, CXCL11,* and *SLAMF7,*

149      all had an area under the curve (AUC) of at least 96% with their lower bound 97%

150      confidence intervals above 90% (Fig 2A, Fig 3C, S6 Table).

151  **Fig 2. Technical and biological validation for selected DEG discovered from RNA**

152  **sequencing.** (A) Tukey boxplots with RT-qPCR normalized (2-3 reference genes) $\log_2$

153  expression values (A.U) according to clinical and histopathological diagnosis. ODD

154  samples are colored according to *M. leprae* 16S rRNA qPCR status as positive (blue)

155  or negative (green). (B) $\log_2$FC from MB-ODD and PB-ODD comparisons estimated

156  from Bayesian linear mixed models and their 95% credible intervals. (C) Tukey boxplot

157  highlighting *IDO1* RT-qPCR normalized $\log_2$ expression values by final diagnosis

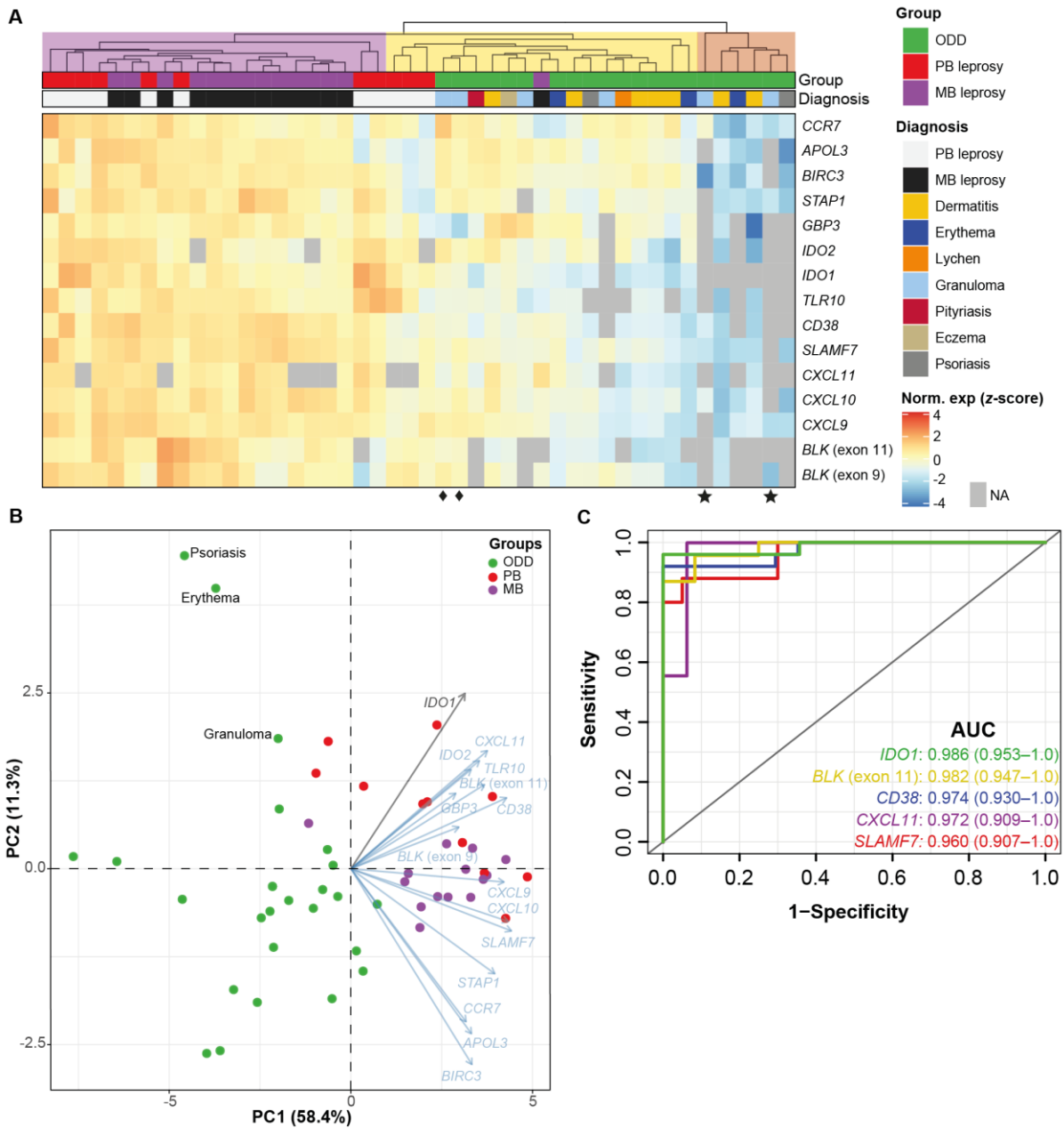158  grouped into ODD category. Missing values are omitted.

159

160       Next, hierarchical clustering with RT-qPCR data including missing values for

161  some genes (no target gene amplification by RT-qPCR) was performed to examine all

162  samples simultaneously. The analysis roughly revealed three major clusters (Fig 3A).

163  At the highest tree subdivision, one small cluster (n=6) with the dendrogram grouped

164  in light brown was composed of ODD samples with lower expression levels (Fig 3A).

165  Due to several ODD having missing values, we confirmed that these samples had

166  similar gene expression for the reference genes, thereby eliminating the possibility of

167  insufficient cDNA input. Another cluster, grouped in the light purple dendrogram,

168  included all MB and most PB samples (except four in light yellow dendrogram). GA

169  samples displayed two patterns, the first with two samples showing undetectable *IDO1*

170  expression (Fig 3A, bottom star symbols). The second set (n=4) is scattered among

171  other ODD samples (Fig 3A). It can be seen that GA and PB samples show highly

172  similar expression profiles for some genes (Fig 3A bottom diamond symbols),

173  reinforcing the difficulty in clinically discriminating between these two conditions, and

174  underlining the relevance of their inclusion in our comparisons [29–31].

10

175    Then, by applying principal component analysis (PCA) to the 15 gene signature

176    obtained with the expanded sample panel tested by RT-qPCR, we uncovered two

177    major patterns separating leprosy lesions from ODD (Fig 3B). As expected, MB

178    samples appeared more homogeneous than PB and ODD samples, while the latter

179    were more dispersed revealing heterogeneous expression patterns (Fig 3B).

180    Next, we quantified the individual classification potential of these genes in

181    distinguishing leprosy from ODD using ROC analysis on RT-qPCR data. *IDO1*

182    expression alone was found to be 98% accurate using an arbitrary threshold, followed

183    by *BLK* (exon 11), *CD38, CXCL11,* and *SLAMF7* (Fig 3C and S6 Table)*.* Finally, to

184    confirm the causal link between mycobacteria and our gene-set, we evaluated the

185    mRNA profiles induced by other live-mycobacteria using a public RNA-seq dataset

186    [32]. We observed that most gene expression signatures, including *IDO1,* could be

187    successfully replicated as induced by either *M. leprae* and/or other mycobacteria (Fig

188    1 in Appendix S1 and S7 Table). By contrast, some of the tested genes such as *BLK*,

189    *CXCL9*, *MS4A1,* and *TLR10* were not differentially expressed in any of the *in vitro*

190    assays with mycobacteria  (Fig 1 in Appendix S1 and S7 Table).

191

192

11

**Fig 3. Hierarchical clustering of RT-qPCR replicated DEG and ROC analysis.** (A)
Hierarchical clustering with scaled and centered normalized log₂ RT-qPCR expression
values (arbitrary units) and annotated according to group and specific diagnosis.
Dendrogram tree was cut arbitrarily and cluster analysis is for hypothesis generating
purposes only. Two samples had more than 13 missing expression values and were
removed from A. (B) Principal component analysis (PCA) with 15 genes measured by
RT-qPCR and using log₂ normalized scaled data. For PCA only, missing values were
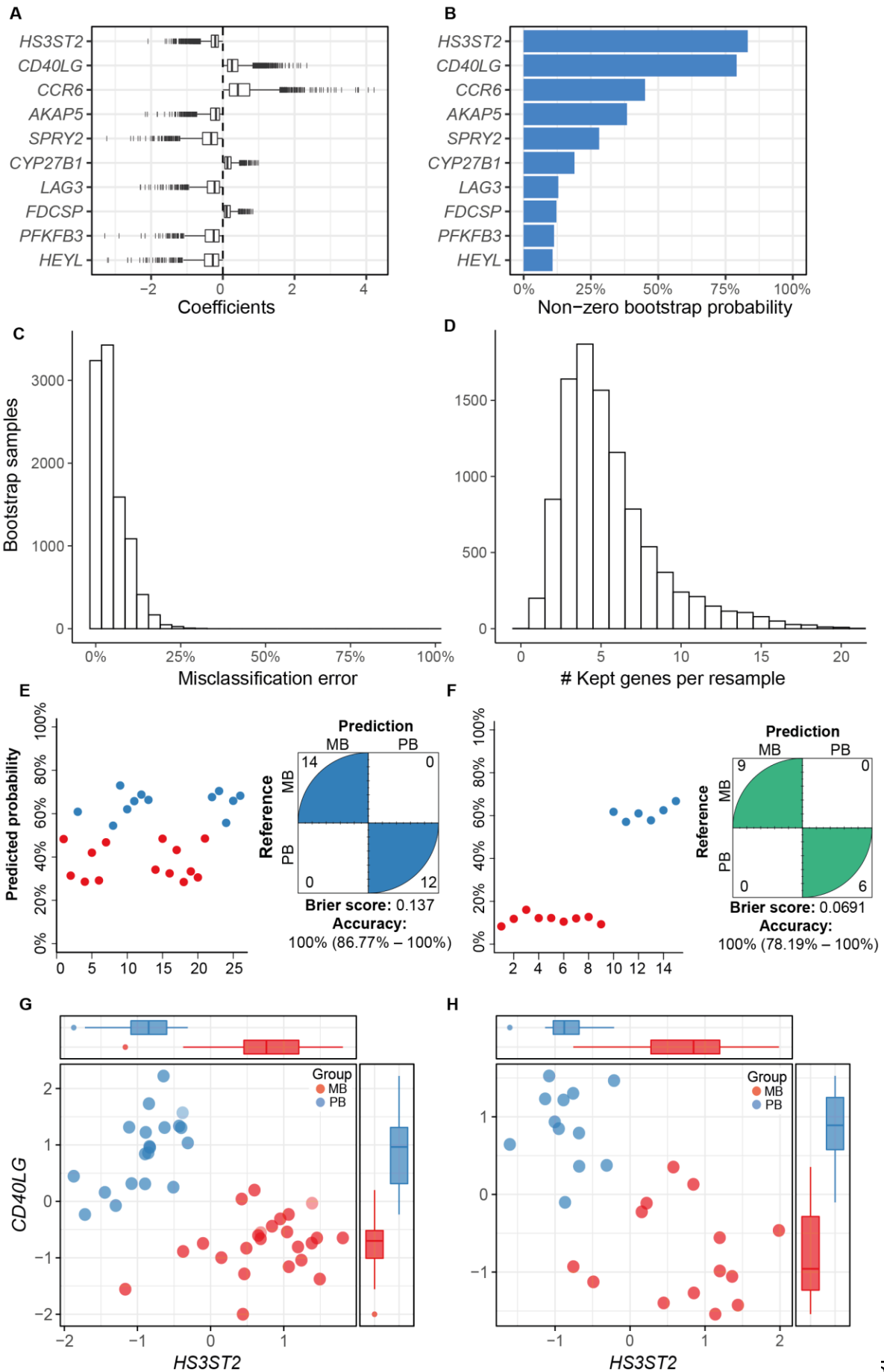
12

200    imputed by the gene arithmetic mean. NA, not amplified, i.e., Cp > 40. In this regard,

201    there were two outliers (psoriasis and erythema), which are samples with high

202    numbers of NA values and that were imputed using the gene arithmetic mean. (C)

203    Receiver operating characteristic analysis for genes with largest AUC (97%

204    confidence intervals) from RT-qPCR replication samples (complete data are shown in

205    S6 Table). See also S1 Appendix and S1 Fig.

## MB and PB gene expression profiling and mRNA-based classifier

208            To define a small subset of genes with high classificatory potential (i.e. with

209    non-overlapping expression values) to distinguish MB from PB lesions, we performed

210    a penalized logistic regression (LASSO) model with k-fold cross-validation trained on

211    the public microarray dataset [24]. This dataset was chosen because of the higher

212    number of PB/MB samples compared to our RNA-seq dataset. As a result, three genes

213    with non-zero coefficients were selected by the cross-validated LASSO model:

214    *HS3ST2, CD40LG,* and *CCR6*, but only the first two genes were most frequently

215    (~80%) selected across 10,000 bootstrapped samples within the training dataset (Fig

216    4A-B). The median misclassification error estimated by the resampling was about 4%

217    (±5.4% median absolute deviation), ranging from 0% to 32% (Fig 4C). Instability

218    assessment in the number of selected genes by LASSO (Fig 4D) showed that most

219    iterations resulted in four non-zero genes (range, 1-20). The final model containing the

220    three genes (*HS3ST2, CD40LG,* and *CCR6*) was evaluated on two test RNA-seq

221    datasets: our dataset and the one from Montoya *et al.* including MB (n=9) and PB

222    (n=6) groups [28]. Penalized logistic regression demonstrated an accuracy of 100%

13

223    (lower 95% CIs: 86.8% and 78.2%, respectively) in classifying MB from PB samples

224    in both test RNA-seq datasets; yet, the Brier score indicated a better performance in

225    Montoya's et al. dataset, probably due to a more homogenous sampling (Fig 4E-F).

226    The *HS3ST2* gene was consistently more expressed in MB leprosy lesions compared

227    to PB, whereas the opposite was observed for *CD40LG* (Fig 4E-H) and *CCR6* (S2

228    Fig). In both datasets, the combined expression levels of *HS3ST2* and *CD40LG*

229    showed good discrimination between the two groups (Fig 4E-H). However, given the

230    sample size and the bootstrapped estimates, it is not currently possible to exclude

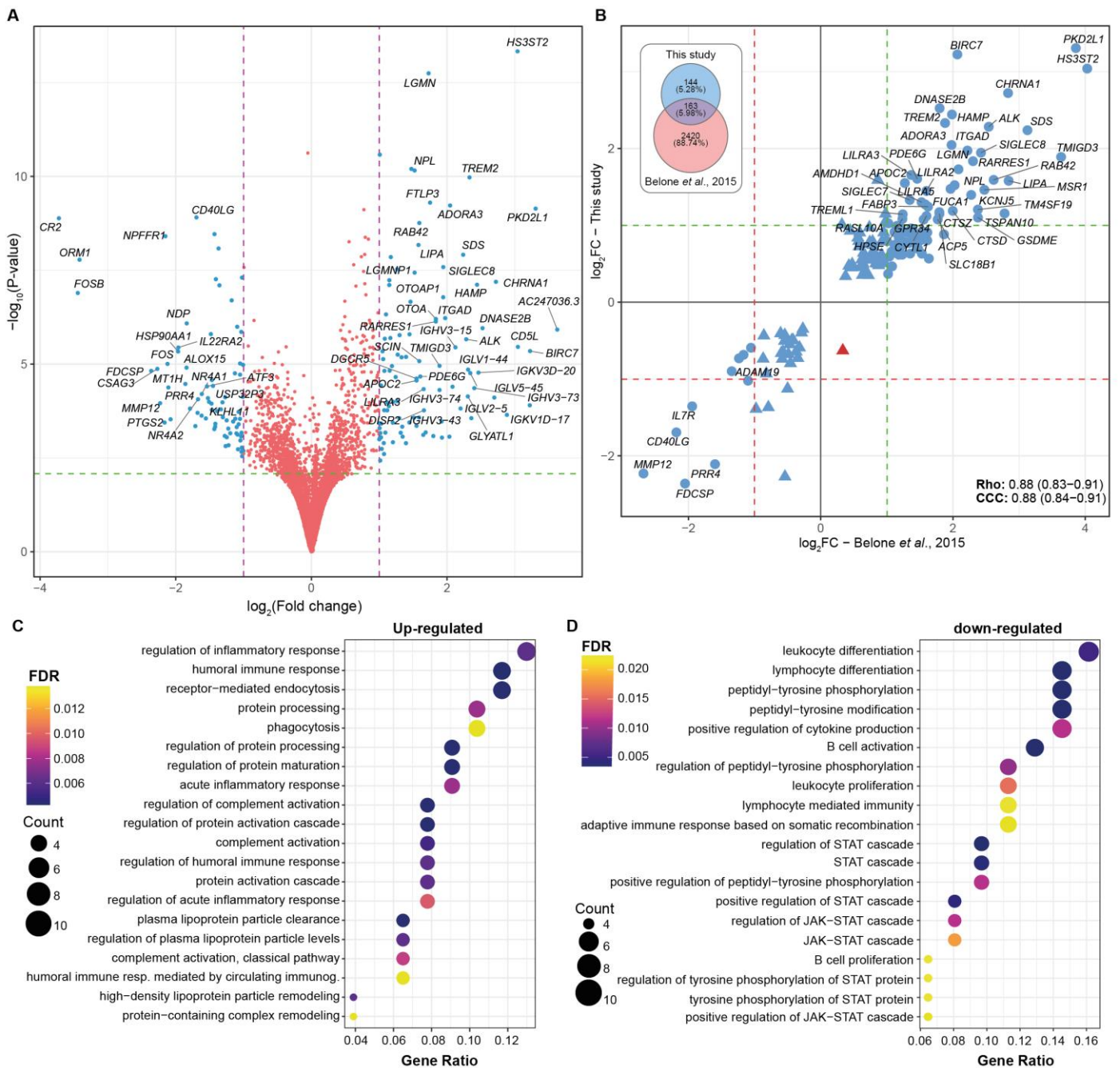231    *CCR6* from the model without additional replication.

232

233    **Fig 4. Gene candidates identified with the penalized logistic regression (LASSO)**

234    **model as the most important to distinguish PB and MB leprosy lesions.** (A)

235    Coefficients (log odds) from the top 10 most selected genes (i.e., non-zero) across

236    10,000 bootstrap samples using the microarray from Belone *et al.* as training dataset.

237    (B) Frequency of non-zero coefficients across all bootstrap samples. (C)

238    Misclassification error distribution estimated from 4-fold cross-validation (k-) across

239    10,000 bootstrap samples, with median error of 3.70% (±5.4% median absolute

240    deviation). (D) Number of genes kept across all resamples. Predicted probability from

241    the final model performance on this study test RNA-seq (E) and Montoya *et al.* RNA-

242    seq (F). Normalized $\log_2$ gene expression (z-score) of the two most frequently selected

243    variables for distinguishing MB from PB samples in the (G) microarray training dataset

244    and (H) this study test RNA-seq. PB, paucibacillary leprosy; MB, multibacillary leprosy.

245    Tukey box plots with 1st, 2nd and 3rd quartiles ± 1.5 × inter quartile range (IQR)

246    whiskers. See also S2 Fig.

247

248         Next, to assess the dichotomy beyond cellular *vs.* humoral response in leprosy

249    lesions [33,34], a comparison of gene expression in MB leprosy (LL+BL+BB) *vs.* PB

250    (TT+BT) skin lesions was performed. Differential expression analysis with $|\log_2 FC| \geq$

251    1 and FDR ≤ 0.01 resulted in 112 DEGs; 69 up-regulated and 43 down-regulated (Fig

252    5A and S8 Table). In addition, we compared DEG to the public microarray data

253    available in Gene Expression Omnibus (GEO) from Belone *et al.* [24,35] using only

254    the FDR cutoff. With an FDR < 0.01, 161 DEGs were common to both studies, all

255    except one showed concordant modulation characterized by an overall high

256    correlation coefficient and concordance index, irrespective of the technology used, the

257    sample processing, and the data analysis methods (Fig 5B). Functional enrichment

258    analysis of the RNA-seq up-regulated genes (i.e., more expressed in MB than PB)

259    revealed processes involved with regulation of immune response, humoral immunity,

260    phagocytosis, cholesterol metabolism, complement activation among others (Fig 5C

261    and S9 Table). On the contrary, enrichment analysis of genes more expressed in PB

262    revealed biological processes such as leukocyte differentiation, lymphocyte

263    differentiation, lymphocyte-mediated immunity, B cell activation, STAT cascade

264    activation/regulation, and JAK-STAT cascade activation (Fig 5D and S10 Table),

265    which are consistent with exacerbated responses in granulomatous diseases.

266    Localized clinical forms, i.e., BT and TT, show a gene expression pattern indicative of

267    differentiation towards epithelioid transformation and granuloma assembly, which is

268    also observed in cutaneous or pulmonary sarcoidosis [36,37].

269

**Fig 5. Differentially expressed genes from multibacillary (MB) *vs.* paucibacillary (PB) leprosy lesions.** (A) Volcano plot showing DEG from the MB *vs.* PB comparison, where blue points are DE with |log$_2$FC| ≥1 and FDR < 0.1. (B) Scatter plots with the 161 DEG common between this study and Belone *et al.* (24) microarray for the same comparison. Red and green dashed lines indicate log$_2$FC of -1 and 1, respectively. Blue points are genes with the same modulation signal and red indicates discordancy. Rho, Spearman's rank correlation coefficient. CCC, Lin's concordance correlation

18

277    coefficient. Venn diagram on the right displays the number of DEG in each study

278    according to FDR < 0.01. (C) Biological processes from GO enriched from up-

279    regulated and (D) down-regulated DEG. FDR, false discovery rate.

280

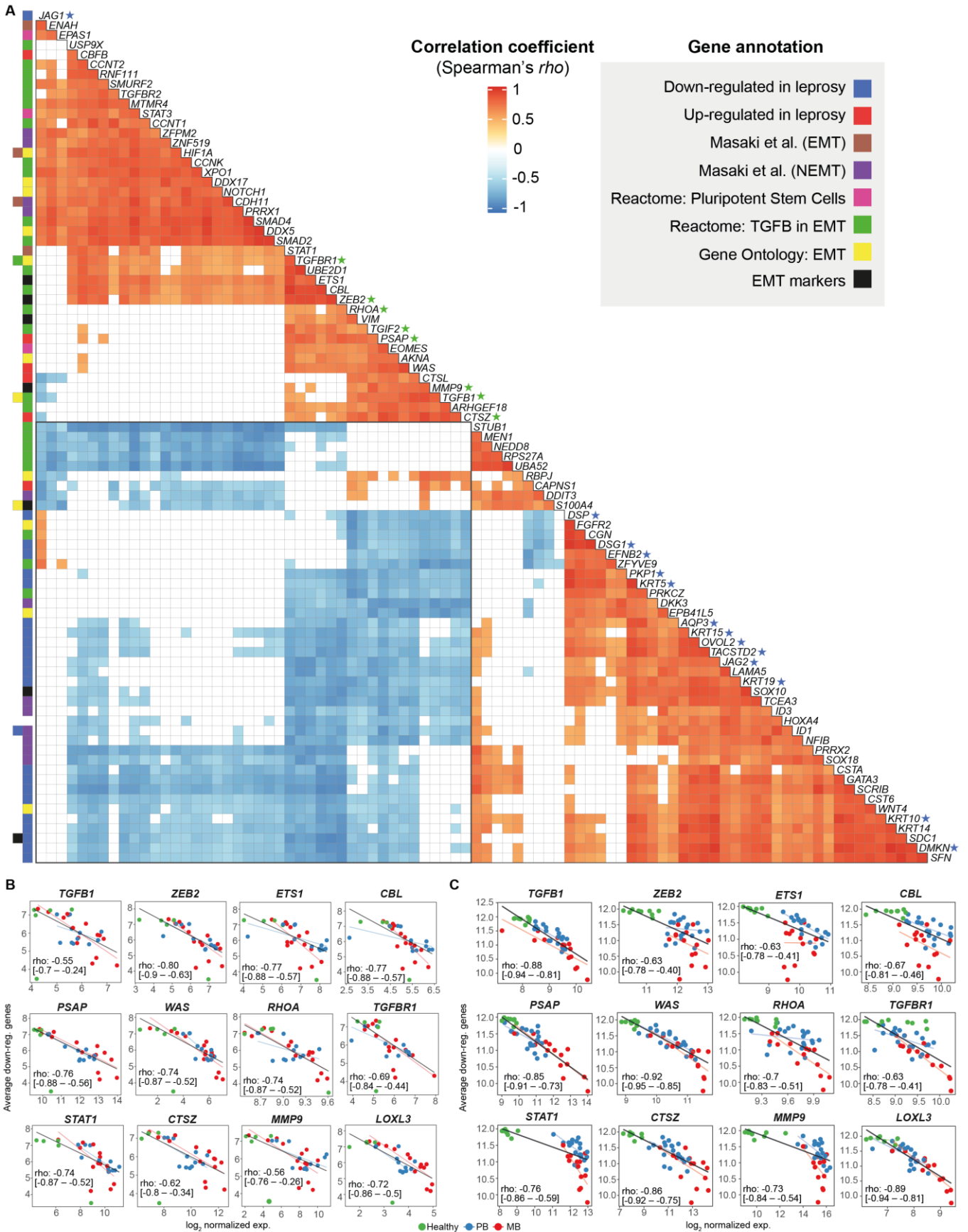## Epithelial-mesenchymal transition (EMT) in the skin of multibacillary leprosy patients

283    To make the most of our dataset, we sought to test a previous hypothesis

284    generated from our group's microarray meta-analysis results, in which we have

285    identified a consistent down-regulation of cornification, keratinocyte differentiation,

286    and epidermal development-related genes in leprosy lesions, predominantly in MB

287    [35]. We first hypothesized that such regulation could result from *M. leprae* inducing

288    dedifferentiation of keratinocytes, similar to the phenomenon described previously in

289    infected Schwann cells [38], and also seen in skin cancer by a process known as

290    epithelial-mesenchymal transition (EMT) [39,40]. To test the hypothesis that such

291    modulation was involved with EMT, we correlated the expression of the previously

292    identified down-regulated genes in leprosy [35] with a collection of genes involved with

293    previously Schwann cell dedifferentiation by *M. leprae* (Masaki *et al.* [38] signatures

294    for EMT and non-EMT genes), positive markers of EMT (from literature), as well as

295    annotated EMT and mesenchymal-related genes from Reactome (R.HSA.452723,

296    R.HSA.5619507.3, R.HSA.2173791) and Gene Ontology (GO0001837) databases.

297    Briefly, the normalized $\log_2$ expression matrices were filtered to retain only genes of

298    interest. Then, the pairwise expression correlation for all genes was calculated using

299    the Spearman's rank correlation procedure. Finally, after adjusting the P-values for

19

300    multiple testing, the genes with any pairwise correlation passing FDR ≤ 1 × 10$^{-4}$ and

301    *rho* ≤ -0.8 were visualized using a heat plot. As result, with this study's RNA-seq, we

302    found a consistent moderate negative correlation between keratinization, cornification,

303    and epidermal development genes (Fig 6A, blue stars, *AQP3, DMKN, DSG1, DSP,*

304    *EFNB2, JAG1, JAG2, KRT5, KRT10, KRT15, KRT19, OVOL2, PKP1, TACSTD2*) with

305    those involved with canonical/alternative EMT and mesenchymal phenotypes (Fig 6A,

306    green stars, *CTSZ, MMP9, PSAP, RHOA, TGFBR1, TGIF2, ZEB2, TGFB1*).

307    Interestingly, the strongest correlations with epidermal/keratinocyte genes was with

308    TGFβ-EMT-related genes (Fig. 6A blue block), as opposed to Masaki et al. non-EMT

309    and other mesenchymal/pluripotency pathways. Next, we replicated these

310    observations with Belone *et al.* microarray [24] and Montoya *et al.* RNA-seq datasets

311    [28], respectively. In Fig 6BC the strongest and representative correlations from

312    TGF☐-EMT-related pathway and a keratinocyte/epidermal gene signature are shown

313    in detail, while  the remaining are available in Fig. S3-4.

314        Overall, these results showed a decreased expression pattern of EMT-related

315    genes in healthy skin samples, and a linear expression increase in PB and MB

316    patients, especially with the microarray dataset, except for *MMP9* (Fig 6C). This was

317    accompanied by the previously reduced expression of cytokeratins and epidermal

318    development genes observed in leprosy. From these results, we hypothesize that in

319    addition to TGF☐-dependent immunosuppression in MB patients, activation of this

320    pathway could be slowing or arresting keratinocyte cornification processes in leprosy

321    lesions thereby both facilitating survival and/or spread of *M. leprae*. If not involved with

322    dedifferentiation of keratinocytes or other epithelial cells, an alternative explanation

323    would be loss of epithelial barrier in MB patients, possibly enlightening a new *M. leprae*

324    transmission route. Further mechanistic experiments ought to determine the causality

325    of our observations and test these findings in light of our hypothetical explanations of

326    the phenomenon.

327

328    **Fig 6. Strongest correlations between keratinocyte and EMT-related genes in**

329     **leprosy lesions.** (A) Heat plot with Spearman's *rho* correlation coefficient of the

330     strongest correlations after multiple testing adjustment with at least one gene-pair

331     passing FDR ≤ 0.0001 and rho ≤ -0.8. Correlations with FDR > 0.1 are filled with white.

332     Row colored squares identify gene annotations. Scatter plots of average $log_2$

333     expression calculated with keratinocyte/epidermal development-related genes

334     previously documented as down-regulated in leprosy skin against dedifferentiation-

335     related genes using either (B) this study RNA-seq dataset or (C) Belone *et al.*

336     microarray (GSE74481). Lines were drawn based on intercept and beta parameters

337     estimated from robust linear regression for all samples (black line) or separately for

338     PB (blue line), and MB (red line). Spearman's *rho* coefficient along with 95% nominal

339     confidence intervals are shown inside scatter plots calculated from all samples. See

340     also S3 Fig and S4 Fig.

# Discussion

341

342     One of the priorities in leprosy research is the development of reliable and

343     accurate laboratory diagnosis tools for all leprosy forms to provide efficient treatment

344     and prevent disability [41]. This goal includes diagnosing patients with early forms of

345     the disease, those with low or mild apparent symptoms, thus assisting with ambiguous

346     differential diagnoses, and even classifying the disease for treatment (MB *vs.* PB) [4].

347     Host response to infection as measured by gene expression in skin biopsies

348     offers diagnostic, prognostic and predictive potential. By applying host transcriptomics

349     to skin lesions from leprosy patients and other common confounding dermatoses that

350     challenge clinicians and pathologists [9,30], we identified a small set of genes that

351     provide a promising expression signature capable of distinguishing PB leprosy cases

23

352  from other confounding dermatological diseases. The top candidate, *IDO1,* is a gene

353  involved in nutritional immunity and metabolism [42–45]. Alone, the expression of this

354  gene was able to differentiate leprosy from non-leprosy lesions with high accuracy in

355  our dataset and in others. According to the latest data from single-cell analysis [46],

356  *IDO1* has been shown to be differentially expressed in Langerhans cells from leprosy

357  lesions compared to healthy skin, corroborating our findings. However, *IDO1*

358  expression is also increased in other mycobacterial diseases such as tuberculosis

359  [47,48], which might decrease its specificity. The accuracy of classification could be

360  improved by combining measurement of *IDO1* expression with that of four other

361  biomarker genes *BLK, CXCL11, CD38, TLR10 and SLAMF7,* which also showed high

362  classification accuracy in the replication dataset*.* In parallel, the penalized logistic

363  regression model, evaluated on two independent datasets, demonstrated that

364  *HS3ST2* and *CD40LG* hold potential to differentiate between MB and PB lesions. In

365  parallel, the penalized logistic regression model, evaluated on two independent

366  datasets, demonstrated that *HS3ST2* and *CD40LG* hold potential to differentiate

367  between MB and PB lesions. We recognize that there is no clinical utility in classifying

368  MB from PB lesions with laboratory assays because this can be done during

369  anamnesis alone. Hence, we aimed at identifying molecular features differing not only

370  in the measure of effect ($log_2FC$) but also having little overlap between the lesion

371  types, as this may point to previously unexplored genes and pathways relevant to

372  future investigation. Considering the functional evidence for *HS3ST2* [49], it is possible

373  that this gene may be involved with granuloma disassembly, tissue permeability, and

374  cellular migration in leprosy, which would explain its overexpression in MB lesions. On

375  the contrary, *CD40LG* (also known as CD154) is more expressed in PB patients when

376  compared to MB with a predominant role in the activation of the microbicidal *Th1*

24

377  response associated with PB lesions [50]. After mechanistic validation of our findings,

378  quantifying expression levels of *HS3ST2* and *CD40LG* from leprosy lesions could be

379  useful to assess immune responsiveness against *M. leprae*, help patient stratification

380  and/or provide a basis for host-based adjuvant treatment for leprosy lesions.

381      One of the challenges in translating gene expression signatures into medical

382  diagnosis is the cost of measuring a large number of genes and transforming these

383  values into a unique continuous or binary classifier. So far, we were able to reproduce

384  the findings using both bulk RNA-sequencing and relative RT-qPCR, with the latter

385  being more accessible to clinicians at least in reference centers or central hospitals.

386  Although there are successful approved RT-qPCR relative gene expression-based

387  diagnostic tests for diagnosing sepsis [12], clinical support for prostate [22], and breast

388  cancer [18], there is a need for alternatives to reduce the cost and complexity of such

389  assays. Quantification of mRNA based on isothermal amplification either with NASBA

390  [51,52], RT-LAMP [53,54] or CRISPR-Cas12 [55] is conceivable for less specialized

391  settings without high-end equipment. Besides, combining a multi-target expression-

392  based diagnostic test with qPCR detection of *M. leprae* DNA could increase the

393  specificity and sensitivity of leprosy diagnosis [56]. Alternatively, an ELISA assay

394  measuring the levels of IDO1 protein from skin interstitial fluid, for example, could be

395  proven useful [57]. Further studies ought to be done selecting tangible diagnostic

396  thresholds and devising a proper classification system to allow the biomarker to

397  function unsupervised.

398      In parallel with poor diagnosis, lack of fundamental understanding of leprosy

399  pathogenesis has misled scientists for centuries [5,6]. Herein, we also compared the

400  two leprosy poles, MB and PB, and identified several pathways already known to be

25

401    associated with leprosy, such as the humoral immune response, phagocytosis, and

402    complement activation. Genes involved with cholesterol and fatty acids were more

403    expressed in MB lesions, as already reported [58–60]. Interestingly, B-cell-related

404    genes were more expressed in PB than MB. In fact, it seems that both poles modulate

405    this pathway by a distinct set of genes. Involvement of B lymphocytes in PB leprosy

406    pathogenesis has been described by a few groups, which may indicate differential

407    involvement of such cells depending on the disease pole [61,62].

408         *M. leprae* subverts host cell metabolism [63] by inducing lipid biosynthesis,

409    while avoiding type II (IFN-gamma) responses through a type I IFNs mechanism,

410    following the phagolysosomal breach that releases DNA into the cytosol [64].

411    However, exactly how the bacilli spread throughout the body and bypass the

412    microbicidal immune response remains unknown. Here, we provide robust evidence

413    indicating that *M. leprae* may induce EMT in the skin within keratinocytes and

414    macrophages, as described in Schwann cells [38]. Indeed, *M. leprae* induced

415    dedifferentiation of infected Schwann cells into an immature stage resembling

416    progenitor/stem-like phenotype [38]. These reprogramming events induced by long-

417    term infection with *M. leprae* resulted in mesenchymal cells capable of migratory and

418    immune-permissive behavior, which in turn facilitated *M. leprae* spread to skeletal and

419    smooth muscles and furthered macrophage recruitment [38,65]. In our previous work,

420    we identified a down-regulated signature of keratinocyte differentiation and

421    cornification gene markers in MB skin lesions [35]. Here, we showed that such genes

422    are inversely correlated with genes involved with EMT, especially the members of the

423    TGFβ-EMT pathway, such as *TGFB1*, *TGFBR1*, *TGIF2, PSAP, ZEB2* [66,67]. Some

424    of these genes are directly or indirectly associated with EMT, such as a *PSAP* [68],

26

425     *WAS* [69], *RHOA* [70–73], *CTSZ* [74], *MMP9* [75], *LOXL3* [76], *HIF1A* [77,78] among

426     others.

427         Our hypothesis that *M. leprae* is inducing dedifferentiation or slowing the

428     cornification process in keratinocytes is plausible, given the evidence in Schwann cells

429     and a few reports of infection in this cell type (Fig 7) [79,80]. Nevertheless, other

430     phenomena could explain EMT's role in leprosy pathogenesis, such as wound healing

431     or loss of the epithelial barrier. Although, given its obligatory intracellular lifestyle, *M.*

432     *leprae* induces dedifferentiation in other cell types, either directly as in Schwann cells

433     or indirectly via chemokine and cytokine production in lesions. Besides inducing

434     keratinocyte dedifferentiation to mesenchymal cells, *M. leprae* might benefit from a

435     decreased or alternative immune activation of these cells [81,82]. Further functional

436     confirmatory experiments should elucidate the causality of this correlation and provide

437     definitive evidence of the relationship between the bacilli and other cell types, such as

438     keratinocytes, fibroblasts, and epithelial cells.

439         Our preliminary data also showed that the enriched pathways among PB skin

440     lesions were consistent with profiles observed in other granulomatous diseases, such

441     as noninfectious sarcoidosis and granuloma annulare, or chronic infectious diseases

442     like tuberculosis [37,83–85]. Our findings revealed that PB (TT/BT) lesions have,

443     among others, JAK-STAT cascade activation, which has been implicated in

444     sarcoidosis and GA. Remarkably, the JAK-STAT specific biological inhibitor,

445     tofacitinib, has a potent effect promoting rebalance of exacerbated immunity among

446     sarcoidosis and granuloma annulare patients reestablishing homeostasis [83].

447     Another compound, everolimus, has been shown in experimental models to achieve

448    the same response [37] suggesting that these drugs could be useful to treat PB, but

449    not MB, leprosy.

450        To conclude, our combined findings provide highly discriminatory mRNA

451    signatures from skin lesions that could distinguish leprosy from other dermatological

452    diseases and allow disease classification by monitoring only a handful of genes. In

453    addition, we report new genes and pathways that are likely informative regarding how

454    *M. leprae* interacts with and subverts host cells to promote its spread within the body

455    and subsequent transmission.

456

457

**Fig 7. Hypothetical hourglass model contextualizing the observed findings for leprosy clinical outcomes.** The host-pathogen interaction in the skin leads to opposing leprosy clinical forms. Upon infection, *M. leprae* induces baseline metabolic alterations such as an increase in glucose uptake, modulation of lipid biosynthesis, reduction of mitochondrial metabolism, and upregulation of IDO-1 and type I IFN. Eventually, progression towards an unspecified inflammatory state can be observed where three ways could be anticipated: I) self-healing; II) progression towards the tuberculoid pole; or III) progression to lepromatous pole. These outcomes are driven

29

466 by specific environmental and host genetic factors. It is expected that lower (or shorter)

467 *M. leprae* exposure, food shortage, BCG vaccination, and polymorphisms in genes

468 controlling autophagy/granuloma formation (*NOD2*, *LRRK2*, *PRKN*) all contribute to

469 developing leprosy per se. Excessive inflammation is one phenotype observed, that is

470 also seen in other granulomatous diseases (e.g., cutaneous sarcoidosis, granuloma

471 annulare), especially in paucibacillary lesions. On the other pole, epithelial-

472 mesenchymal transition and local immunosuppression are present due to a probably

473 higher (and/or longer) *M. leprae* exposure, combined with host single-nucleotide

474 polymorphisms (SNPs) at key genes, like lipid biogenesis (*APOE*) and central

475 metabolism (*HIF1A, LACC1/FAMIN*), culminating in disease progression.

# Materials and Methods

## Patient cohort

478 All patients were enrolled after informed written consent was obtained with

479 approval from the Ethics Committee of the Oswaldo Cruz Foundation, number 151/01.

480 Leprosy clinical forms were classified according to the criteria of Ridley and Jopling

481 [2]. Leprosy patients were treated according to the operational criteria established by

482 the World Health Organization [4]. Leprosy and patients with other dermatological

483 diseases were eligible if their diagnosis was confirmed by clinical and histopathological

484 findings. Additionally, detection of *M. leprae* DNA by qPCR routinely performed in our

485 laboratory could be employed to support diagnosis [56,86]. HIV and hepatitis B

486 positive patients were not included in this study, in addition, we excluded individuals

487 with a current or previous history of tuberculosis. No other comorbidities were used to

488 exclude patients and further individual information is available in S1 Table. Skin biopsy

489 specimens containing both epidermis and dermis were obtained with 3 mm (diameter)

490 sterile punches following local anesthesia from the lesion site. Skin biopsies were

491 immediately stored in one milliliter of RNALater (Ambion, Thermo Fisher Scientific Inc.,

492 MA, USA) according to the manufacturer's instructions and stored in liquid nitrogen

493 until RNA isolation. Healthy skin biopsies were from lesion-free sites of patients

494 diagnosed with indeterminate or pure neural leprosy.

## Study Design

496 The main objective of this research was to identify host gene expression

497 patterns capable of distinguishing leprosy (including the PB forms) from other

498 differential diagnosis of skin lesions. Our working hypothesis was that leprosy lesions,

499 despite their morphological and histopathological similarity to other skin diseases, may

500 induce distinct patterns of gene expression in at a small subset. We predefined the

501 comparison of leprosy (PB+MB) from non-leprosy including GA in addition to healthy

502 patients for RNA sequencing experiment. In addition, we predetermined comparisons

503 between leprosy poles: MB *vs.* PB. Our samples are representative of a population of

504 individuals attending the Sousa Araujo Outpatient Clinic based in Rio de Janeiro,

505 Brazil, which also receives patients from surrounding municipalities.

## RNA isolation

507 Snap frozen skin biopsies were thawed in wet ice and processed using TRIzol

508 Reagent (Ambion, Thermo Fisher Scientific Inc., MA, USA) according to the

509 manufacturer's instructions with the help of Polytron Homogenizer PT3100

510    (Kinematica AG, Switzerland). RNA was treated with DNAse using the DNAfree kit

511    (Thermo Fisher Scientific Inc., MA, USA) according to the standard manufacturer's

512    protocol, prior to use for library preparation and RT-qPCR. RNA integrity was

513    assessed in 1% agarose gel electrophoresis or TapeStation RNA ScreenTape (Agilent

514    Technology, CA, USA). During RNA isolation, samples were randomly assigned to

515    extraction batches and freeze-thaw cycles to minimize batch effects and the

516    introduction of technical artifacts. All procedures applied to samples were carried out

517    using reagents from the same lot. The first author conducted the experiments aware

518    of each sample group during the entire process, therefore, no blinding scheme was

519    used, although we do not rely on perceptual/abstract measurements or analyses nor

520    did we purposefully exclude samples.

## Library preparation and Illumina RNA sequencing

522    RNA-seq libraries were prepared with 1 µg of total RNA for each sample using

523    the Illumina TruSeq mRNA kit (Illumina, USA) as recommended by the manufacturer

524    using the Illumina CD RNA indexes (Illumina, USA). Libraries were quantified and

525    qualified using a qPCR quantification protocol guide (KAPA Library Quantification Kits

526    for Illumina Sequencing platforms) and TapeStation D1000 ScreenTape (Agilent

527    Technologies, USA), respectively. The resulting libraries (fragment size 200-350bp)

528    were multiplexed (17, 17, and 19 libraries, respectively) and sequenced using the

529    NextSeq 500 platform (Illumina, USA), generating approximately 520 million single-

530    end reads of 75 nucleotides in length.

## RNA-sequencing analysis

531

532     RAW bcl files were converted into .fastq using Illumina's bcl2fastq script. Then,

533     read quality was assessed using FastQC version 0.11.8 [87]. Next, transcript counts

534     were estimated using Salmon (v.1.13.0) quasi-mapping (human transcriptome

535     GRCh38_cdna sourced from Ensembl/RefGenie plus pre-computed salmon index,

536     http://refgenomes.databio.org/#hg38_cdna) with default settings and --seqBias flag

537     set [88]. Transcript counts were summarized into ENSEMBL gene counts using the R

538     v.3.6.1 package tximport v.1.12.0 [89,90] and biomaRt v.2.40.5 [91]. The expression

539     of sex-chromosome-specific genes, such as *UTY* and *XIST,* was used to rule out

540     sample mislabeling. Differential expression was estimated using DESEq2 v.1.24.0,

541     after filtering out weakly expressed genes with less than 10 counts per million and less

542     than 15 total counts in 70% of samples  [92–94]. In addition to the patient's biological

543     sex, extraction batch and sequencing run, three surrogate variables estimated with

544     RUVseq v.1.18.0 were included in DESeq2's generalized linear model [95,96].

545     Nominal P-values were inspected with histograms and adjusted for multiple testing

546     according to the method [97] proposed for controlling the false discovery rate (FDR).

547     All $\log_2$ fold-changes were shrunken prior to DE filtering with the apeglm [94] or normal

548     algorithms. For visualization, counts per million (CPM) were computed with edgeR's

549     cpm function v.3.26.1 and variance stabilized with the parametric method [92]. Then,

550     surrogate variables and covariates were regressed out from the expression matrix

551     using limma's removeBatchEffect [98–100] before being visualized with ggplot2

552     v.3.3.0 [101]. Hierarchical clustering, heatmaps, and ROC analysis were all performed

553     with the previously processed expression matrix. Heatmap with hierarchical clustering

554     was drawn with ComplexHeatmap v.2.0.0 [102] or pheatmap v.1.0.12 [103] using

555    gene-wise scaled and centered matrix with Euclidean distance and average

556    agglomeration method. Overrepresentation analysis (ORA) was used to test for Gene

557    Ontology Biological Process (GO BP) enrichment with clusterProfiler v.3.12.0 [104]

558    and org.Hs.eg.db v.3.8.2 annotations [105]. Up and down-regulated lists were used as

559    inputs and the background list was composed of all genes subjected to differential

560    expression. P-values were adjusted for multiple testing using the Benjamini-Hochberg

561    method [97]. Raw and normalized RNA sequencing data are available in EMBL-EBI's

562    ENA and ArrayExpress under accessions ERP128243 and E-MTAB-10318,

563    respectively.

## RT-qPCR

565    A total of 2.5 µg of RNA was reversed transcribed into cDNA using 4 µL of Vilo

566    Master Mix (Thermo Fisher Scientific Inc., USA) according to the manufacturer's

567    instructions. Then, cDNA was diluted to a final concentration of 5 ng/µL using TE buffer

568    (10 mM Tris-HCL and 0.1 mM EDTA in RNAse-free water). RT-qPCR was performed

569    using Fast Sybr Master Mix (Thermo Fisher Scientific Inc., USA) in a final reaction

570    volume of 10 µL. For each reaction, performed in duplicate, 5 µL of Fast Sybr Green

571    were combined with 200 nM of each primer, 10 ng of cDNA, and q.s.p of injection-

572    grade water. Thermal cycling and data acquisition were performed on Viia7 with 384

573    well block (Applied Biosystems, Thermo Fisher Scientific Inc., USA) following the

574    master mix manufacturer cycling preset with a final melting curve analysis (65 °C to

575    95 °C, captured at every 0.5 °C). All primers were designed with NCBI Primer-Blast

576    [106–109] to either flank intron(s) or span exon-exon junction(s) to avoid gDNA

577    amplification (S11 Table). Further, primers were quality checked for specificity, dimers

578 and hairpin with MFEPrimer v.3.0 [110,111] and IDT's oligoAnalyzer

579 (https://www.idtdna.com/calc/analyzer). Data were exported from QuantStudio

580 software v.1.3 in RDML format, which was imported to LinRegPCR v.2020.0 for RT-

581 qPCR efficiency determination and calculation of the $N_0$ value [112,113]. Finally, $N_0$

582 values were imported to R and normalized using as the denominator the normalization

583 factor (NF) calculated from the geometric mean of at least three reference genes

584 (*RPS16*, *RPL35* and *QRICH1*), which were previously tested for stability [114]. These

585 $N_0$ normalized values were used for visualization in Fig 2A. For mean difference

586 estimation between groups, RT-qPCR data were analyzed in a Bayesian framework

587 (Markov Chain Monte Carlo sampling, MCMC) using generalized linear mixed effect

588 models under lognormal-Poisson error with MCMC.qpcr v.1.2.4 [115,116]. Per-gene

589 efficiency estimates from LinRegPCR were used in conjunction with Cp (crossing

590 point) calculated in QuantStudio software v.1.3 to generate the counts table. Then, the

591 generalized linear mixed-effect model was fitted using three reference genes (allowing

592 up to 20% between-group variation) with 550,000 iterations, thin = 100, and burn-in of

593 50,000. The model specification included the sample (factor with 51 levels) as a

594 random effect and the diagnosis group (factor with 3 levels) as a fixed effect. MCMC

595 diagnostics were done by inspecting chain mixing plots and linear mixed model

596 diagnostic plots. Ninety-five percent credible intervals were drawn around the posterior

597 means and MCMC equivalent P-values were also computed.

## 598 Reanalysis of public gene expression datasets

599 Belone and collaborators GSE74481 [24] and de Toledo-Pinto and cols.

600 GSE35423 [64] microarray datasets were reanalyzed as described elsewhere [35].

601 Blischak and cols. [32] RNA-seq dataset (GSE67427) was reanalyzed from counts per

602 sample file from the author's Bitbucket repository (https://bitbucket.org/jdblischak/tb-

603 data/src/master/). Briefly, a normalized $log_2$ expression matrix was regressed out for

604 RNA integrity number and extraction batch variables. Then, differences in gene

605 expression (48h post-infection) for specific genes and treatments were tested using a

606 gene-wise linear mixed model with a random intercept per sample (replicate) followed

607 by Dunnet comparison against a "mock" group using emmeans v.1.5.3. Montoya and

608 collaborators' dataset was retrieved from GEO (GSE125943) already normalized

609 (DESeq2 median ratio method) and transformed with base 2 logarithm with no further

610 processing [28].

## Correlation analyses

612       For RNA-seq datasets, normalized $log_2$ counts-per-million values were used

613 and $log_2$ normalized intensities for microarray. Spearman's rank correlation method

614 was chosen because it is robust against outliers, does not rely on normality

615 assumption, and also identifies monotonic but non-linear relationships. Initially, a list

616 of keratinocyte/cornification/epidermal development genes that were DE in the meta-

617 analysis was assembled [35]. Then, lists of target genes were compiled from results

618 of Masaki *et al.* [38]: EMT and non-EMT; from Reactome: R-HSA-452723

619 (Transcriptional regulation of pluripotent stem cells), R-HAS-5619507.3 (Activation of

620 HOX genes during differentiation), R-HAS-2173791 (TGFβ receptor signaling in EMT);

621 Gene Ontology GO:0001837 (EMT), and literature for EMT canonical markers. Next

622 pairwise Spearman correlation was calculated using the Hmisc's rcorr function v.4.2-

623 0 for every pair of genes from keratinocyte/epidermal development and EMT gene

624 lists. P-values were adjusted for multiple testing using the BH method for FDR control

625     for all tests [97]. Additionally, 95% nominal confidence intervals were calculated using

626     the Fieller method implemented by correlation R package v.0.5.0 [117,118]. To

627     visualize the results, only genes with at least one pairwise correlation with Spearman's

628     rho coefficient ≤ -0.8 and FDR ≤ 0.0001 were selected. Additionally, the average $\log_2$

629     expression from genes involved with keratinocyte/epidermal development was

630     calculated and used in scatter plots against the expression of the EMT genes. Scatter

631     plots were drawn with ggplot2 v.3.3.3 showing lines from coefficients estimated using

632     default robust regression (MASS::rlm v.7.3-51.4) either for all samples or stratified by

633     group. No outliers were omitted.

## Regularized (LASSO) logistic regression classification

635         Normalized $\log_2$ expression matrices regressed out for covariates and batches

636     were used as input predictors. The model was trained using the microarray dataset

637     from Belone et al. [24] with penalized regression (L1-norm, LASSO) and 4-fold cross-

638     validation (k-fold CV) with the negative binomial log-likelihood link function, glmnet

639     v.4.1 [119–121]. Predictors were standardized to have mean zero and unit variance

640     inside the cv.glmnet function. We opted for L1-norm because it results in a smaller

641     number of genes (#features ≤ n) with non-zero coefficients, as compared to elastic-

642     net or ridge regression counterparts. In addition, this model is suitable for high-

643     dimensional data as it combines feature selection during model tuning and training,

644     mitigating the effects of predictors' collinearity and reducing overfitting. To assess the

645     coefficients' error, misclassification error rate, feature stability and model size we used

646     non-parametric bootstrap (boot v.1.3.25) with 10,000 samples, with 4-fold cross-

647     validation inside each loop [122,123]. The final LASSO model selected by 4-fold cross-

37

648 validation contained three non-zero genes. Finally, independent RNA-seq test

649 datasets were used to compute the accuracy of the final model. Alternatively, the

650 whole process was repeated with leave-one-out cross-validation instead of k-fold. The

651 results were practically indistinguishable, especially regarding the feature stability

652 (data not shown).

## 653 Sample sizes

654 The sample size for RNA sequencing was selected based on previous leprosy

655 work with microarrays, aiming at detecting genes with at least a differential fold-change

656 of two. For RT-qPCR validation, sample size calculation was performed using the per-

657 gene standardized effect size estimated from the RNA-seq data, aiming at a power of

658 85% and alpha = 0.03. No samples were discarded after successful data collection

659 (i.e. outliers). In the end, the sample sizes per group for RT-qPCR were: MB = 14,

660 PB=11, ODD = 23. All RT-qPCR reactions were conducted in duplicate for each

661 biological unit (here, a fragment of a skin biopsy derived from an individual).

## 662 RT-qPCR and ROC statistical analyses

663 Normalized RT-qPCR gene expression data were $\log_2$ transformed before use

664 in data visualization. Additionally, we checked if the Bayesian results remained

665 consistent using a more common procedure (data not shown). For this, the mean

666 normalized expression (from $N_0$) was compared pairwise for the prior stipulated groups

667 using Welch's t-test implemented in R language, using the predetermined alpha of

668 0.03. Normality assumption was verified with normal quantile-quantile plots (qqplots,

669    car v. 3.0-2). In cases where quantile-quantile plots showed huge deviation from

670    theoretical normal distribution, the Wilcoxon Rank Sum was used to verify results.

671    Receiver Operating Curve (ROC) analysis was used to determine the accuracy

672    (measured by the area under the curve, AUC) and its respective best classification

673    threshold, aiming at maximizing AUC with equal importance for sensitivity and

674    specificity. Confidence intervals (95%) for AUC were calculated using the Delong non-

675    parametric method as implemented in pROC v.1.15.3 [124–126].

## Data and code reporting

677    Raw .fastq data are available in EMBL-EBI European Nucleotide Archive (ENA)

678    database (ERP128243). Raw Salmon counts and normalized batch cleaned

679    expression matrices are available in EMBL-EBI ArrayExpress, under E-MTAB-10318,

680    along with experimental and phenotypic metadata. R source code and accompanying

681    intermediate data used in all analyses in this manuscript are also readily available

682    through Zenodo, doi.org/10.5281/zenodo.4682010.

## Acknowledgements

39

# References

691

692    1.    Britton WJ, Lockwood DN. Leprosy. The Lancet. 2004;363: 1209–1219.

693    doi:10.1016/S0140-6736(04)15952-7

694    2.    Ridley DS, Jopling WH. Classification of leprosy according to immunity. A five-

695    group system. Int J Lepr Mycobact Dis Off Organ Int Lepr Assoc. 1966;34: 255–73.

696    3.    Scollard DM, Adams LB, Gillis TP, Krahenbuhl JL, Truman W, Williams DL.

697    The Continuing Challenges of Leprosy The Continuing Challenges of Leprosy. Clin

698    Microbiol Rev. 2006;19: 338–381. doi:10.1128/CMR.19.2.338

699    4.    WHO. Guidelines for the Diagnosis, Treatment and Prevention of Leprosy.

700    Geneva: World Health Organization; 2018 p. 106.

701    5.    WHO. Global leprosy (Hansen disease) update, 2019: time to step-up

702    prevention initiatives. Wkly Epidemiol Rec. 2020;95: 417–440.

703    6.    Nath I, Saini C, Valluri VL. Immunology of leprosy and diagnostic challenges.

704    Clin Dermatol. 2015;33: 90–98. doi:10.1016/j.clindermatol.2014.07.005

705    7.    van Hooij A, Tjon Kon Fat EM, Batista da Silva M, Carvalho Bouth R, Cunha

706    Messias AC, Gobbo AR, et al. Evaluation of Immunodiagnostic Tests for Leprosy in

707    Brazil, China and Ethiopia. Sci Rep. 2018;8: 1–9. doi:10.1038/s41598-018-36323-1

708    8.    van Hooij A, van den Eeden S, Richardus R, Tjon Kon Fat E, Wilson L,

709    Franken KLMC, et al. Application of new host biomarker profiles in quantitative point-

710    of-care tests facilitates leprosy diagnosis in the field. EBioMedicine. 2019;47: 301–

711    308. doi:10.1016/j.ebiom.2019.08.009

712   9.      Manta FS de N, Leal-Calvo T, Moreira SJM, Marques BLC, Ribeiro-Alves M,

713   Rosa PS, et al. Ultra-sensitive detection of Mycobacterium leprae: DNA extraction

714   and PCR assays. Poonawala H, editor. PLoS Negl Trop Dis. 2020;14: e0008325.

715   doi:10.1371/journal.pntd.0008325

716   10.     Gliddon HD, Herberg JA, Levin M, Kaforou M. Genome-wide host RNA

717   signatures of infectious diseases: discovery and clinical translation. Immunology.

718   2018;153: 171–178. doi:10.1111/imm.12841

719   11.     Ko ER, Yang WE, McClain MT, Woods CW, Ginsburg GS, Tsalik EL. What

720   was old is new again: Using the host response to diagnose infectious disease.

721   Expert Rev Mol Diagn. 2015;15: 1143–1158. doi:10.1586/14737159.2015.1059278

722   12.     Miller RR, Lopansri BK, Burke JP, Levy M, Opal S, Rothman RE, et al.

723   Validation of a host response assay, SeptiCyte LAB, for discriminating sepsis from

724   systemic inflammatory response syndrome in the ICU. Am J Respir Crit Care Med.

725   2018;198: 903–913. doi:10.1164/rccm.201712-2472OC

726   13.     Van Hooij A, Fat EMTK, Van Den Eeden SJF, Wilson L, Da Silva MB,

727   Salgado CG, et al. Field-friendly serological tests for determination of M. Leprae-

728   specific antibodies. Sci Rep. 2017;7: 1–8. doi:10.1038/s41598-017-07803-7

729   14.     Warsinske H, Vashisht R, Khatri P. Host-response-based gene signatures for

730   tuberculosis diagnosis: A systematic comparison of 16 signatures. PLoS Med.

731   2019;16. doi:10.1371/journal.pmed.1002786

732   15.     Röltgen K, Pluschke G, Spencer JS, Brennan PJ, Avanzi C. The immunology

733   of other mycobacteria: M. ulcerans, M. leprae. Semin Immunopathol. 2020;42: 333–

734   353. doi:10.1007/s00281-020-00790-4

735    16.    Mesko B, Poliska S, Nagy L. Gene expression profiles in peripheral blood for

736    the diagnosis of autoimmune diseases. Trends Mol Med. 2011;17: 223–233.

737    doi:10.1016/j.molmed.2010.12.004

738    17.    Wang B, Chen S, Zheng Q, Gao Z, Chen R, Xuan J, et al. Development and

739    initial validation of diagnostic gene signatures for systemic lupus erythematosus. Ann

740    Rheum Dis. 2019. doi:10.1136/annrheumdis-2019-216695

741    18.    Carlson JJ, Roth JA. The impact of the Oncotype Dx breast cancer assay in

742    clinical practice: A systematic review and meta-analysis. Breast Cancer Res Treat.

743    2013;141: 13–22. doi:10.1007/s10549-013-2666-z

744    19.    Gordon GJ, Jensen RV, Hsiao LL, Gullans SR, Blumenstock JE, Ramaswamy

745    S, et al. Translation of microarray data into clinically relevant cancer diagnostic tests

746    using gene expression ratios in lung cancer and mesothelioma. Cancer Res.

747    2002;62: 4963–4967.

748    20.    Narrandes S, Xu W. Gene expression detection assay for cancer clinical use.

749    J Cancer. 2018;9: 2249–2265. doi:10.7150/jca.24744

750    21.    Clark-Langone KM, Sangli C, Krishnakumar J, Watson D. Translating tumor

751    biology into personalized treatment planning: analytical performance characteristics

752    of the Oncotype DX®Colon Cancer Assay. BMC Cancer. 2010;10: 691.

753    doi:10.1186/1471-2407-10-691

754    22.    Knezevic D, Goddard AD, Natraj N, Cherbavaz DB, Clark-Langone KM,

755    Snable J, et al. Analytical validation of the Oncotype DX prostate cancer assay - a

756    clinical RT-PCR assay optimized for prostate needle biopsies. BMC Genomics.

757    2013;14: 1–12. doi:10.1186/1471-2164-14-690

758    23.    Laible M, Schlombs K, Kaiser K, Veltrup E, Herlein S, Lakis S, et al. Technical

759    validation of an RT-qPCR in vitro diagnostic test system for the determination of

760    breast cancer molecular subtypes by quantification of ERBB2 , ESR1 , PGR and

761    MKI67 mRNA levels from formalin- fixed paraffin-embedded breast tumor

762    specimens. BMC Cancer. 2016; 1–14. doi:10.1186/s12885-016-2476-x

763    24.    Belone A de FF, Rosa PS, Trombone APF, Fachin LRV, Guidella CC, Ura S,

764    et al. Genome-wide screening of mRNA expression in leprosy patients. Front Genet.

765    2015;6: 1–12. doi:10.3389/fgene.2015.00334

766    25.    Jorge KTOS, Souza RP, Assis MTA, Araújo MG, Locati M, Jesus AMR, et al.

767    Characterization of MicroRNA Expression Profiles and Identification of Potential

768    Biomarkers in Leprosy. J Clin Microbiol. 2017;55: 1516–1525.

769    doi:10.1128/JCM.02408-16

770    26.    Tió-Coma M, van Hooij A, Bobosha K, van der Ploeg-van Schip JJ, Banu S,

771    Khadge S, et al. Whole blood RNA signatures in leprosy patients identify reversal

772    reactions before clinical onset: a prospective, multicenter study. Sci Rep. 2019;9:

773    17931. doi:10.1038/s41598-019-54213-y

774    27.    Tió-Coma M, Kiełbasa SM, van den Eeden SJF, Mei H, Roy JC, Wallinga J, et

775    al. Blood RNA signature RISK4LEP predicts leprosy years before clinical onset.

776    EBioMedicine. 2021;68: 103379. doi:10.1016/j.ebiom.2021.103379

777    28.    Montoya DJ, Andrade P, Silva BJA, Teles RMB, Ma F, Bryson B, et al. Dual

778    RNA-Seq of Human Leprosy Lesions Identifies Bacterial Determinants Linked to

779    Host Immune Response. Cell Rep. 2019;26: 3574-3585.e3.

780    doi:10.1016/j.celrep.2019.02.109

781    29.    Bhatia S, Shenoi SD, Pai K, Srilatha PS. Granuloma multiforme: an

782    uncommon differential for leprosy. Trop Doct. 2019;49: 55–58.

783    doi:10.1177/0049475518803191

784    30.    Kundakci N, Erdem C. Leprosy: A great imitator. Clin Dermatol. 2019;37:

785    200–212. doi:10.1016/j.clindermatol.2019.01.002

786    31.    Zhu TH, Kamangar F, Silverstein M, Fung MA. Borderline Tuberculoid

787    Leprosy Masquerading as Granuloma Annulare: A Clinical and Histological Pitfall.

788    Am J Dermatopathol. 2017;39: 296–299. doi:10.1097/DAD.0000000000000698

789    32.    Blischak JD, Tailleux L, Mitrano A, Barreiro LB, Gilad Y. Mycobacterial

790    infection induces a specific human innate immune response. Sci Rep. 2015;5: 1–16.

791    doi:10.1038/srep16882

792    33.    Modlin RL. Th1-Th2 paradigm: insights from leprosy. J Invest Dermatol.

793    1994;102: 828–832. doi:10.1111/1523-1747.ep12381958

794    34.    Yamamura M, Uyemura K, Deans RJ, Weinberg K, Rea TH, Bloom BR, et al.

795    Defining protective responses to pathogens: Cytokine profiles in leprosy lesions.

796    Science. 1991;254: 277–279. doi:10.1126/science.1925582

797    35.    Leal-Calvo T, Moraes MO. Reanalysis and integration of public microarray

798    datasets reveals novel host genes modulated in leprosy. Mol Genet Genomics.

799    2020;295: 1355–1368. doi:10.1007/s00438-020-01705-6

800    36.    Judson MA, Marchell RM, Mascelli M, Piantone A, Barnathan ES, Petty KJ, et

801    al. Molecular profiling and gene expression analysis in cutaneous sarcoidosis: the

802    role of interleukin-12, interleukin-23, and the T-helper 17 pathway. J Am Acad

803    Dermatol. 2012;66: 901–910, 910.e1–2. doi:10.1016/j.jaad.2011.06.017

804   37.     Linke M, Pham HTT, Katholnig K, Schnöller T, Miller A, Demel F, et al.

805   Chronic signaling via the metabolic checkpoint kinase mTORC1 induces

806   macrophage granuloma formation and marks sarcoidosis progression. Nat Immunol.

807   2017;18: 293–302. doi:10.1038/ni.3655

808   38.     Masaki T, Qu J, Cholewa-Waclaw J, Burr K, Raaum R, Rambukkana A.

809   Reprogramming adult Schwann cells to stem cell-like cells by leprosy bacilli

810   promotes dissemination of infection. Cell. 2013;152: 51–67.

811   doi:10.1016/j.cell.2012.12.014

812   39.     Brabletz T, Kalluri R, Nieto MA, Weinberg RA. EMT in cancer. Nat Rev

813   Cancer. 2018;18: 128–134. doi:10.1038/nrc.2017.118

814   40.     Pastushenko I, Blanpain C. EMT Transition States during Tumor Progression

815   and Metastasis. Trends Cell Biol. 2019;29: 212–226. doi:10.1016/j.tcb.2018.12.001

816   41.     Khazai Z, Van Brakel W, Essink D, Gillis T, Kasang C, Kuipers P, et al.

817   Reviewing Research Priorities of the Leprosy Research Initiative (LRI): a

818   stakeholder's consultation. Lepr Rev. 2019;90: 3–30. doi:10.47276/lr.90.1.3

819   42.     Chen W. IDO: more than an enzyme. Nat Immunol. 2011;12: 809–811.

820   doi:10.1038/ni.2088

821   43.     Greco FA, Coletti A, Camaioni E, Carotti A, Marinozzi M, Gioiello A, et al. The

822   Janus-faced nature of IDO1 in infectious diseases: challenges and therapeutic

823   opportunities. Future Med Chem. 2016;8: 39–54. doi:10.4155/fmc.15.165

824   44.     Melé M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, et al.

825   The human transcriptome across tissues and individuals. Science. 2015;348: 660–

826   665. doi:10.1126/science.aaa0355

827    45.    Yamazaki F, Kuroiwa T, Takikawa O, Kido R. Human indolylamine 2,3-

828    dioxygenase. Its tissue distribution, and characterization of the placental enzyme.

829    Biochem J. 1985;230: 635–638. doi:10.1042/bj2300635

830    46.    Hughes TK, Wadsworth MH, Gierahn TM, Do T, Weiss D, Andrade PR, et al.

831    Second-Strand Synthesis-Based Massively Parallel scRNA-Seq Reveals Cellular

832    States and Molecular Features of Human Inflammatory Skin Pathologies. Immunity.

833    2020;53: 878-894.e7. doi:10.1016/j.immuni.2020.09.015

834    47.    Gautam US, Foreman TW, Bucsan AN, Veatch AV, Alvarez X, Adekambi T, et

835    al. In vivo inhibition of tryptophan catabolism reorganizes the tuberculoma and

836    augments immune-mediated control of Mycobacterium tuberculosis. Proc Natl Acad

837    Sci U S A. 2018;115: E62–E71. doi:10.1073/pnas.1711373114

838    48.    Yeung AWS, Terentis AC, King NJC, Thomas SR. Role of indoleamine 2,3-

839    dioxygenase in health and disease. Clin Sci. 2015;129: 601–672.

840    doi:10.1042/CS20140392

841    49.    Denys A, Allain F. The emerging roles of heparan sulfate 3-O-

842    sulfotransferases in cancer. Front Oncol. 2019;9. doi:10.3389/fonc.2019.00507

843    50.    Yamauchi PS, Bleharski JR, Uyemura K, Kim J, Sieling PA, Miller A, et al. A

844    Role for CD40-CD40 Ligand Interactions in the Generation of Type 1 Cytokine

845    Responses in Human Leprosy. J Immunol. 2000;165: 1506–1512.

846    doi:10.4049/jimmunol.165.3.1506

847    51.    Heim A. Highly sensitive detection of gene expression of an intronless gene:

848    amplification of mRNA, but not genomic DNA by nucleic acid sequence based

849    amplification (NASBA). Nucleic Acids Res. 1998;26: 2250–2251.

850    doi:10.1093/nar/26.9.2250

851    52.    Patterson SS, Casper ET, Garcia-Rubio L, Smith MC, Paul JH. Increased

852    precision of microbial RNA quantification using NASBA with an internal control. J

853    Microbiol Methods. 2005;60: 343–352. doi:10.1016/j.mimet.2004.10.011

854    53.    Ganguli A, Ornob A, Spegazzini N, Liu Y, Damhorst G, Ghonge T, et al.

855    Pixelated spatial gene expression analysis from tissue. Nat Commun. 2018;9.

856    doi:10.1038/s41467-017-02623-9

857    54.    Pandey M, Singh D, Onteru SK. Reverse transcription loop-mediated

858    isothermal amplification (RT-LAMP), a light for mammalian transcript analysis in low-

859    input laboratories. J Cell Biochem. 2018;119: 4334–4338. doi:10.1002/jcb.26624

860    55.    Broughton JP, Deng X, Yu G, Fasching CL, Servellita V, Singh J, et al.

861    CRISPR–Cas12-based detection of SARS-CoV-2. Nat Biotechnol. 2020;38: 870–

862    874. doi:10.1038/s41587-020-0513-4

863    56.    Barbieri RR, Manta FSN, Moreira SJM, Sales AM, Nery JAC, Nascimento

864    LPR, et al. Quantitative polymerase chain reaction in paucibacillary leprosy

865    diagnosis: A follow-up study. PLoS Negl Trop Dis. 2019;13: e0007147.

866    doi:10.1371/journal.pntd.0007147

867    57.    Strassner JP, Rashighi M, Ahmed Refat M, Richmond JM, Harris JE. Suction

868    blistering the lesional skin of vitiligo patients reveals useful biomarkers of disease

869    activity. J Am Acad Dermatol. 2017;76: 847-855.e5. doi:10.1016/j.jaad.2016.12.021

870    58.    Elamin AA, Stehr M, Singh M. Lipid Droplets and Mycobacterium leprae

871    Infection. J Pathog. 2012;10. doi:10.1155/2012/361374

872    59.    Lobato LS, Rosa PS, Ferreira J da S, Neumann A da S, da Silva MG, do

873    Nascimento DC, et al. Statins increase rifampin mycobactericidal effect. Antimicrob

874    Agents Chemother. 2014;58: 5766–74. doi:10.1128/AAC.01826-13

875    60.    Wang D, Zhang D-F, Li G-D, Bi R, Fan Y, Wu Y, et al. A pleiotropic effect of

876    the APOE gene: association of APOE polymorphisms with multibacillary leprosy in

877    Han Chinese from Southwest China. Br J Dermatol. 2018;178: 931–939.

878    doi:10.1111/bjd.16020

879    61.    Fabel A, Giovanna Brunasso AM, Schettini AP, Cota C, Puntoni M, Nunzi E,

880    et al. Pathogenesis of Leprosy. Am J Dermatopathol. 2019;41: 422–427.

881    doi:10.1097/DAD.0000000000001310

882    62.    Iyer AM, Mohanty KK, van Egmond D, Katoch K, Faber WR, Das PK, et al.

883    Leprosy-specific B-cells within cellular infiltrates in active leprosy lesions. Hum

884    Pathol. 2007;38: 1065–1073. doi:10.1016/j.humpath.2006.12.017

885    63.    Medeiros RCA, Girardi K do C de V, Cardoso FKL, Mietto B de S, Pinto TG

886    de T, Gomez LS, et al. Subversion of Schwann Cell Glucose Metabolism by

887    Mycobacterium leprae. J Biol Chem. 2016;291: 21375–21387.

888    doi:10.1074/jbc.M116.725283

889    64.    de Toledo-Pinto TG, Ferreira ABR, Ribeiro-Alves M, Rodrigues LS, Batista-

890    Silva LR, Silva BJ de A, et al. STING-Dependent 2′-5′ Oligoadenylate Synthetase–

891    Like Production Is Required for Intracellular Mycobacterium leprae Survival. J Infect

892    Dis. 2016;214: 311–320. doi:10.1093/infdis/jiw144

893    65.    Hess S, Rambukkana A. Bacterial-induced cell reprogramming to stem cell-

894    like cells: new premise in host–pathogen interactions. Curr Opin Microbiol. 2015;23:

895    179–188. doi:10.1016/j.mib.2014.11.021

896    66.    Vandewalle C, Comijn J, De Craene B, Vermassen P, Bruyneel E, Andersen

897    H, et al. SIP1/ZEB2 induces EMT by repressing genes of different epithelial cell-cell

898    junctions. Nucleic Acids Res. 2005;33: 6566–6578. doi:10.1093/nar/gki965

899    67.    DaSilva-Arnold SC, Kuo CY, Davra V, Remache Y, Kim PCW, Fisher JP, et

900    al. ZEB2, a master regulator of the epithelial-mesenchymal transition, mediates

901    trophoblast differentiation. Mol Hum Reprod. 2018;25: 61–75.

902    doi:10.1093/molehr/gay053

903    68.    Jiang Y, Zhou J, Hou D, Luo P, Gao H, Ma Y, et al. Prosaposin is a biomarker

904    of mesenchymal glioblastoma and regulates mesenchymal transition through the

905    TGF-β1/Smad signaling pathway. J Pathol. 2019;249: 26–38. doi:10.1002/path.5278

906    69.    Frugtniet BA, Martin TA, Zhang L, Jiang WG. Neural Wiskott-Aldrich

907    syndrome protein (nWASP) is implicated in human lung cancer invasion. BMC

908    Cancer. 2017;17. doi:10.1186/s12885-017-3219-3

909    70.    Bendris N, Arsic N, Lemmers B, Blanchard JM. Cyclin A2, Rho GTPases and

910    EMT. Small GTPases. 2012;3: 225–228. doi:10.4161/sgtp.20791

911    71.    Bhowmick NA, Ghiassi M, Bakin A, Aakre M, Lundquist CA, Engel ME, et al.

912    Transforming growth factor-β1 mediates epithelial to mesenchymal

913    transdifferentiation through a RhoA-dependent mechanism. Mol Biol Cell. 2001;12:

914    27–36. doi:10.1091/mbc.12.1.27

915    72.    Salvi A, Thanabalu T. WIP promotes in-vitro invasion ability, anchorage

916    independent growth and EMT progression of A549 lung adenocarcinoma cells by

917    regulating RhoA levels. Biochem Biophys Res Commun. 2017;482: 1353–1359.

918    doi:10.1016/j.bbrc.2016.12.040

919    73.    Wang Q, Yang X, Xu Y, Shen Z, Cheng H, Cheng F, et al. RhoA/Rho-kinase

920    triggers epithelial-mesenchymal transition in mesothelial cells and contributes to the

921    pathogenesis of dialysis-related peritoneal fibrosis. Oncotarget. 2018;9: 14397–

922    14412. doi:10.18632/oncotarget.24208

923    74.    Wang J, Chen L, Li Y, Guan XY. Overexpression of cathepsin Z contributes to

924    tumor metastasis by inducing epithelial-mesenchymal transition in hepatocellular

925    carcinoma. PLoS ONE. 2011;6. doi:10.1371/journal.pone.0024967

926    75.    Lin CY, Tsai PH, Kandaswami CC, Lee PP, Huang CJ, Hwang JJ, et al.

927    Matrix metalloproteinase-9 cooperates with transcription factor Snail to induce

928    epithelial-mesenchymal transition. Cancer Sci. 2011;102: 815–827.

929    doi:10.1111/j.1349-7006.2011.01861.x

930    76.    Peinado H, del Carmen Iglesias-de la Cruz M, Olmeda D, Csiszar K, Fong

931    KSK, Vega S, et al. A molecular role for lysyl oxidase-like 2 enzyme in Snail

932    regulation and tumor progression. EMBO J. 2005;24: 3446–3458.

933    doi:10.1038/sj.emboj.7600781

934    77.    Tam SY, Wu VWC, Law HKW. Hypoxia-Induced Epithelial-Mesenchymal

935    Transition in Cancers: HIF-1α and Beyond. Front Oncol. 2020;10.

936    doi:10.3389/fonc.2020.00486

937    78.    Zhu Y, Tan J, Xie H, Wang J, Meng X, Wang R. HIF-1α regulates EMT via the

938    Snail and β-catenin pathways in paraquat poisoning-induced early pulmonary

939    fibrosis. J Cell Mol Med. 2016;20: 688–697. doi:10.1111/jcmm.12769

940    79.    Lyrio ECD, Campos-Souza IC, Corrêa LCD, Lechuga GC, Verícimo M, Castro

941    HC, et al. Interaction of Mycobacterium leprae with the HaCaT human keratinocyte

942    cell line: new frontiers in the cellular immunology of leprosy. Exp Dermatol. 2015;24:

943    536–542. doi:10.1111/exd.12714

944    80.    Okada S, Komura J, Nishiura M. Mycobacterium leprae found in epidermal

945    cells by electron microscopy. IntJLeprOther MycobactDis. 1978;46: 30–34.

946    81.    Pivarcsi A, Kemény L, Dobozy A. Innate Immune Functions of the

947    Keratinocytes. Acta Microbiol Immunol Hung. 2004;51: 303–310.

948    doi:10.1556/AMicr.51.2004.3.8

949    82.    Pivarcsi A, Nagy I, Lajos K. Innate Immunity in the Skin: How Keratinocytes

950    Fight Against Pathogens. Curr Immunol Rev. 2005;1: 29–43.

951    doi:10.2174/1573395052952941

952    83.    Damsky W, Thakral D, McGeary MK, Leventhal J, Galan A, King B. Janus

953    kinase inhibition induces disease remission in cutaneous sarcoidosis and granuloma

954    annulare. J Am Acad Dermatol. 2020;82: 612–621. doi:10.1016/j.jaad.2019.05.098

955    84.    Flynn JL, Chan J, Lin PL. Macrophages and control of granulomatous

956    inflammation in tuberculosis. Mucosal Immunol. 2011;4: 271–278.

957    doi:10.1038/mi.2011.14

958    85.    Locke LW, Crouser ED, White P, Julian MW, Caceres EG, Papp AC, et al. IL-

959    13–regulated Macrophage Polarization during Granuloma Formation in an In Vitro

960    Human Sarcoidosis Model. Am J Respir Cell Mol Biol. 2019;60: 84–95.

961    doi:10.1165/rcmb.2018-0053OC

962    86.    Manta FSN, Barbieri RR, Moreira SJM, Santos PTS, Nery JAC, Duppre NC,

963    et al. Quantitative PCR for leprosy diagnosis and monitoring in household contacts:

964    A follow-up study, 2011–2018. Sci Rep. 2019;9. doi:10.1038/s41598-019-52640-5

965    87.    Brabaham Bioinformatics. FastQC: A Quality Control Tool for High

966    Throughput Sequence Data [Online]. 2015. Available:

967    http://www.bioinformatics.babraham.ac.uk/projects/fastqc/

968    88.    Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast

969    and bias-aware quantification of transcript expression. Nat Methods. 2017;14: 417–

970    419. doi:10.1038/nmeth.4197

971    89.    R Core Team. R: A language and environment for statistical computing.

972    Vienna, Austria; 2017. Available: https://www.r-project.org/

973    90.    Soneson C, Love MI, Robinson MD. Differential analyses for RNA-seq:

974    transcript-level estimates improve gene-level inferences. F1000Research. 2016;4:

975    1521. doi:10.12688/f1000research.7563.2

976    91.    Durinck S, Moreau Y, Kasprzyk A, Davis S, De Moor B, Brazma A, et al.

977    BioMart and Bioconductor: a powerful link between biological databases and

978    microarray data analysis. Bioinformatics. 2005;21: 3439–3440.

979    doi:10.1093/bioinformatics/bti525

980    92.    Anders S, Huber W. Differential expression analysis for sequence count data.

981    Genome Biol. 2010;11: R106. doi:10.1186/gb-2010-11-10-r106

982    93.    Love MI, Huber W, Anders S. Moderated estimation of fold change and

983    dispersion for RNA-seq data with DESeq2. Genome Biol. 2014;15: 550.

984    doi:10.1186/s13059-014-0550-8

985    94.    Zhu A, Ibrahim JG, Love MI. Heavy-Tailed prior distributions for sequence

986    count data: Removing the noise and preserving large differences. Bioinformatics.

987    2019;35: 2084–2092. doi:10.1093/bioinformatics/bty895

988    95.    Gagnon-Bartsch JA, Speed TP. Using control genes to correct for unwanted

989    variation in microarray data. Biostatistics. 2012;13: 539–552.

990    doi:10.1093/biostatistics/kxr034

991    96.    Risso D, Ngai J, Speed TP, Dudoit S. Normalization of RNA-seq data using

992    factor analysis of control genes or samples. Nat Biotechnol. 2014;32: 896–902.

993    doi:10.1038/nbt.2931

994    97.    Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical

995    and Powerful Approach to Multiple Testing. Journal of the Royal Statistical Society.

996    Series B (Methodological). WileyRoyal Statistical Society; 1995.

997    doi:10.2307/2346101

998    98.    Phipson B, Lee S, Majewski IJ, Alexander WS, Smyth GK. Robust

999    hyperparameter estimation protects against hypervariable genes and improves

1000    power to detect differential expression. Ann Appl Stat. 2016;10: 946–963.

1001    doi:10.1214/16-AOAS920

1002    99.    Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, et al. limma powers

1003    differential expression analyses for RNA-sequencing and microarray studies. Nucleic

1004    Acids Res. 2015;43: e47. doi:10.1093/nar/gkv007

1005    100.    Smyth GK. Linear Models and Empirical Bayes Methods for Assessing

1006    Differential Expression in Microarray Experiments Linear Models and Empirical

1007    Bayes Methods for Assessing Differential Expression in Microarray Experiments.

1008    Stat Appl Genet Mol Biol. 2004;3: 1–26. doi:10.2202/1544-6115.1027

1009    101.    Wickham H. ggplot2-Elegant Graphics for Data Analysis. 1st ed. New York,

1010    NY: Springer New York; 2009. doi:10.1007/978-0-387-98141-3

1011    102.    Gu Z, Eils R, Schlesner M. Complex heatmaps reveal patterns and

1012    correlations in multidimensional genomic data. Bioinformatics. 2016;32: 2847–2849.

1013    doi:10.1093/bioinformatics/btw313

1014    103.    Kolde R. pheatmap: Pretty Heatmaps. 2015. Available: https://cran.r-

1015    project.org/package=pheatmap

1016    104.    Yu G, Wang L-G, Han Y, He Q-Y. clusterProfiler: an R Package for

1017    Comparing Biological Themes Among Gene Clusters. OMICS J Integr Biol. 2012;16:

1018    284–287. doi:10.1089/omi.2011.0118

1019    105.    Carlson M. org.Hs.eg.db: Genome wide annotation for Human. 2019.

1020    Available: 10.18129/B9.bioc.org.Hs.eg.db

1021    106.    Koressaar T, Remm M. Enhancements and modifications of primer design

1022    program Primer3. Bioinformatics. 2007;23: 1289–1291.

1023    doi:10.1093/bioinformatics/btm091

1024    107.    Kõressaar T, Lepamets M, Kaplinski L, Raime K, Andreson R, Remm M.

1025    Primer3_masker: integrating masking of template sequence with primer design

1026    software. Bioinformatics. 2018;34: 1937–1938. doi:10.1093/bioinformatics/bty036

1027    108.    Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, et al.

1028    Primer3—new capabilities and interfaces. Nucleic Acids Res. 2012;40: e115–e115.

1029    doi:10.1093/nar/gks596

1030    109.    Ye J, Coulouris G, Zaretskaya I, Cutcutache I, Rozen S, Madden TL. Primer-

1031    BLAST: a tool to design target-specific primers for polymerase chain reaction. BMC

1032    Bioinformatics. 2012;13: 134. doi:10.1186/1471-2105-13-134

1033    110.    Qu W, Shen Z, Zhao D, Yang Y, Zhang C. MFEprimer: Multiple factor

1034    evaluation of the specificity of PCR primers. Bioinformatics. 2009;25: 276–278.

1035    doi:10.1093/bioinformatics/btn614

1036    111.    Wang K, Li H, Xu Y, Shao Q, Yi J, Wang R, et al. MFEprimer-3.0: Quality

1037    control for PCR primers. Nucleic Acids Res. 2019;47: W610–W613.

1038    doi:10.1093/nar/gkz351

1039    112.    Ramakers C, Ruijter JM, Lekanne Deprez RH, Moorman AFM. Assumption-

1040    free analysis of quantitative real-time polymerase chain reaction (PCR) data.

1041    Neurosci Lett. 2003;339: 62–66. doi:10.1016/S0304-3940(02)01423-4

1042    113.    Ruijter JM, Ramakers C, Hoogaars WMH, Karlen Y, Bakker O, Van den hoff

1043    MJB, et al. Amplification efficiency: Linking baseline and bias in the analysis of

1044    quantitative PCR data. Nucleic Acids Res. 2009;37. doi:10.1093/nar/gkp045

1045    114.    Vandesompele J, De Preter K, Pattyn ilip, Poppe B, Van Roy N, De Paepe A,

1046    et al. Accurate normalization of real-time quantitative RT-PCR data by geometric

1047    averaging of multiple internal control genes. Genome Biol. 2002;3: 34–1.

1048    doi:10.1186/gb-2002-3-7-research0034

1049    115.    Matz MV, Wright RM, Scott JG. No control genes required: Bayesian analysis

1050    of qRT-PCR data. PloS One. 2013;8: 1–12. doi:10.1371/journal.pone.0071448

1051    116.    Steibel JP, Poletto R, Coussens PM, Rosa GJM. A powerful and flexible linear

1052    mixed model framework for the analysis of relative quantification RT-PCR data.

1053    Genomics. 2009;94: 146–152. doi:10.1016/j.ygeno.2009.04.008

1054    117.    Fieller EC, Hartley HO, Pearson ES. TESTS FOR RANK CORRELATION

1055    COEFFICIENTS I. Biometrika. 1957;44: 470–481. doi:10.1093/biomet/44.3-4.470

1056    118.    Makowski D, Ben-Shachar MS, Patil I, Lüdecke D. Methods and Algorithms

1057    for Correlation Analysis in R. J Open Source Softw. 2020;5: 2306.

1058    doi:10.21105/joss.02306

1059    119.    Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear

1060    models via coordinate descent. J Stat Softw. 2010;33: 1–22.

1061    doi:10.18637/jss.v033.i01

1062    120.    Simon N, Friedman JH, Hastie T, Tibshirani R. Regularization Paths for Cox's

1063    Proportional Hazards Model via Coordinate Descent. J Stat Softw. 2011;39: 1–13.

1064    doi:10.18637/jss.v039.i05

1065    121.    Tibshirani R. Regression Shrinkage and Selection via the Lasso. J R Stat Soc

1066    Ser B Methodol. 1996;58: 267–288.

1067    122.    Hastie T, Tibshirani R, Wainwright M. Statistical Learning with Sparsity. 1st

1068    ed. Chapman and Hall/CRC; 2015.

1069    123.    Davison AC, Hinley DV. Bootstrap Methods and Their Application. Cambrige

1070    University Press; 1997. Available: http://statwww.epfl.ch/davison/BMA/

1071    124.    DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the Areas under

1072    Two or More Correlated Receiver Operating Characteristic Curves: A Nonparametric

1073    Approach. Biometrics. 1988;44: 837. doi:10.2307/2531595

1074    125.    Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez JC, et al. pROC:

1075    An open-source package for R and S+ to analyze and compare ROC curves. BMC

1076    Bioinformatics. 2011;12: 77. doi:10.1186/1471-2105-12-77

1077    126.    Sun X, Xu W. Fast implementation of DeLong's algorithm for comparing the

1078    areas under correlated receiver operating characteristic curves. IEEE Signal Process

1079    Lett. 2014;21: 1389–1393. doi:10.1109/LSP.2014.2337313

# Supporting Information

1080

1081    **S1 Appendix. Linking expression profiles to mycobacteria species.**

1082    **S1 Fig. Gene expression in MB and PB groups from test and training datasets.**

1083    Normalized $\log_2$ expression values per group from (A) this study RNA-seq dataset or

1084    (B) Belone *et al.* (GSE74481) [24]. The genes shown were selected in 25%–50% of

1085    the LASSO models (Fig 4B) according to the bootstrap. MB, multibacillary leprosy; PB,

1086    paucibacillary leprosy; TT, tuberculoid leprosy; BT, borderline-tuberculoid; BB,

1087    borderline-borderline; BL, borderline-lepromatous; LL, lepromatous. Each point

1088    represents an independent skin biopsy from a patient. Y-axis values are not

1089    comparable between panels A and B.

1090    **S2 Fig. Strongest correlations between the average expression of genes**

1091    **associated with keratinocyte/cornification against dedifferentiation-related**

1092    **genes using Montoya *et al.* RNA-seq dataset** [28]**.** Scatter plots of scores (average

1093   normalized $\log_2$ expression) calculated from genes with previously documented down-

1094   regulation in leprosy skin lesions against dedifferentiation-related genes with Montoya

1095   *et al.* RNA-seq dataset (GSE125943) [28]. Lines were drawn based on intercept and

1096   beta estimates from robust linear regression for all samples (black) or separately for

1097   TL (tuberculoid leprosy, blue), and LL (lepromatous leprosy, red). X-axis shows $\log_2$

1098   normalized expression values. Spearman's rho are shown along with nominal 95%

1099   confidence intervals inside the plots. Most genes shown have FDR < 0.1 and rho ≤ -

1100   0.6. Related to figure 6.

1101   **S3 Fig. Strongest correlations between modulated genes from**

1102   **keratinocyte/cornification and dedifferentiation-related genes using Belone et**

1103   **al. microarray dataset (GSE74481)** [24]**.** Heat plot with Spearman's rho correlation

1104   coefficient of the strongest correlations from all ontologies screened after multiple

1105   testing adjustment (BH-FDR). Most genes shown have FDR ≤ 0.0001 and rho ≤ -0.7.

1106   Bottom colored rectangles indicate which category the gene was present (some genes

1107   co-occur). Related to figure 6.

1108   **S1 Table. Demographic and clinical metadata from human participants.**

1109   **S2 Table. Genes differentially expressed from leprosy *vs.* non-leprosy with**

1110   **|log$_2$FC| ≥ 1 and FDR ≤ 0.01.**

1111   **S3 Table. Over-representation analysis (ORA) for leprosy *vs.* non-leprosy (up-**

1112   **regulated) differentially expressed genes.**

1113   **S4 Table. ROC analysis from RNA-seq dataset using leprosy *vs.* non-leprosy**

1114   **samples.**

1115    **S5 Table. Posterior log$_2$FC estimates, 95% credible intervals and MCMC P-**

1116    **values from PB-OD and MB-OD comparisons**.

1117    **S6 Table. ROC analysis results using RT-qPCR with the validation dataset**

1118    **(Related to Fig 3).** 95% confidence intervals are shown, except for AUCs of 1.0. The

1119    table is sorted from highest to lowest AUC.

1120    **S7 Table. Log$_2$FC estimates, confidence intervals, and Dunnet *P*-values from**

1121    **distinct mycobacterial stimuli in human macrophages *in vitro*.**

1122    **S8 Table. Genes differentially expressed from multibacillary paucibacillary**

1123    **leprosy with |log$_2$FC| ≥ 1 and FDR ≤ 0.01.**

1124    **S9 Table. Over-representation analysis (ORA) for MB *vs.* PB (up-regulated)**

1125    **differentially expressed genes.**

1126    **S10 Table. Over-representation analysis (ORA) for MB *vs.* PB (down-regulated)**

1127    **differentially expressed genes.**

1128    **S11 Table. Oligonucleotide sequences.**