

# Transcriptomic Profiles of Confirmed Pediatric Tuberculosis Patients and Household Contacts Identifies Active Tuberculosis, Infection, and Treatment Response Among Indian Children

Jeffrey A. Tornheim,<sup>1,✉</sup> Anil K. Madugundu,<sup>2,3,4,5</sup> Mandar Paradkar,<sup>6</sup> Kiyoshi F. Fukutani,<sup>7,8,9</sup> Artur T. L. Queiroz,<sup>7,8</sup> Nikhil Gupte,<sup>1,6</sup> Akshay N. Gupte,<sup>1</sup> Aarti Kinikar,<sup>10</sup> Vandana Kulkarni,<sup>6</sup> Usha Balasubramanian,<sup>6</sup> Sreelakshmi Sreenivasamurthy,<sup>2,3,11</sup> Remya Raja,<sup>2,3,4</sup> Neeta Pradhan,<sup>6</sup> Shri Vijay Bala Yogendra Shivakumar,<sup>12</sup> Chhaya Valvi,<sup>10</sup> Luke Elizabeth Hanna,<sup>13</sup> Bruno B. Andrade,<sup>7,8,9,14,15,a</sup> Vidya Mave,<sup>1,6</sup> Akhilesh Pandey,<sup>2,3,5,11,a</sup> and Amita Gupta<sup>1,16,✉</sup>, for the CTRIUMPH RePORT India Study Team

<sup>1</sup>Center for Clinical Global Health Education, Division of Infectious Diseases, Johns Hopkins University School of Medicine, Baltimore, Maryland, USA, <sup>2</sup>Institute of Bioinformatics, Bangalore, Karnataka, India, <sup>3</sup>Center for Molecular Medicine, National Institute of Mental Health and Neurosciences (NIMHANS), Bangalore, India, <sup>4</sup>Manipal Academy of Higher Education (MAHE), Manipal, Karnataka, India, <sup>5</sup>Department of Laboratory Medicine and Pathology and Center for Individualized Medicine, Mayo Clinic, Rochester, Minnesota, USA, <sup>6</sup>Byramjee Jeejeebhoy Government Medical College—Johns Hopkins University Clinical Research Site, Pune, Maharashtra, India, <sup>7</sup>Instituto Gonçalo Moniz, Fundação Oswaldo Cruz, Salvador, Brazil, <sup>8</sup>Multinational Organization Network Sponsoring Translational and Epidemiological Research (MONSTER) Initiative, Salvador, Brazil, <sup>9</sup>Faculdade de Tecnologia e Ciências (FTC), Salvador, Brazil, <sup>10</sup>Byramjee Jeejeebhoy Government Medical College, Pune, Maharashtra, India, <sup>11</sup>McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland, USA, <sup>12</sup>Johns Hopkins University—India office (Center for Clinical Global Health Education), Pune, Maharashtra, India, <sup>13</sup>National Institute for Research in Tuberculosis, Chennai, Tamil Nadu, India, <sup>14</sup>Universidade Salvador (UNIFACS), Laureate Universities, Salvador, Brazil, <sup>15</sup>Escola Bahiana de Medicina e Saúde Pública (EBMSP), Salvador, Brazil, <sup>16</sup>Johns Hopkins Bloomberg School of Public Health, Baltimore, Maryland, USA

**Background.** Gene expression profiling is emerging as a tool for tuberculosis diagnosis and treatment response monitoring, but limited data specific to Indian children and incident tuberculosis infection (TBI) exist.

**Methods.** Sixteen pediatric Indian tuberculosis cases were age- and sex-matched to 32 tuberculosis-exposed controls (13 developed incident TBI without subsequent active tuberculosis). Longitudinal samples were collected for ribonucleic acid sequencing. Differential expression analysis generated gene lists that identify tuberculosis diagnosis and tuberculosis treatment response. Data were compared with published gene lists. Population-specific risk score thresholds were calculated.

**Results.** Seventy-one genes identified tuberculosis diagnosis and 25 treatment response. Within-group expression was partially explained by age, sex, and incident TBI. Transient changes in gene expression were identified after both infection and treatment. Application of 27 published gene lists to our data found variable performance for tuberculosis diagnosis (sensitivity 0.38–1.00, specificity 0.48–0.93) and treatment response (sensitivity 0.70–0.80, specificity 0.40–0.80). Our gene lists found similarly variable performance when applied to published datasets for diagnosis (sensitivity 0.56–0.85, specificity 0.50–0.85) and treatment response (sensitivity 0.49–0.86, specificity 0.50–0.84).

**Conclusions.** Gene expression profiles among Indian children with confirmed tuberculosis were distinct from adult-derived gene lists, highlighting the importance of including distinct populations in differential gene expression models.

**Keywords.** India; pediatric TB; TB diagnosis; transcriptomics.

Tuberculosis (TB) is the infectious disease responsible for the greatest number of deaths worldwide, and 1 million children develop TB every year [1]. Despite advances in diagnosis, current tests perform poorly in children, for whom sensitivity by sputum smear, culture, and Xpert MTB/RIF are unacceptably low (<15%, 30%–40%, and 66%, respectively) [2–4]. This leaves

most pediatric TB diagnoses unconfirmed, and children are frequently treated based on the combination of clinical, radiographic, and immunological criteria [4]. Even if sputum-based diagnostic tests were as reliable in children as among adults, none of them are good surrogates for treatment outcomes. Conversion of sputum smear and culture to negative after 2 months predicts treatment success with limited sensitivity (24% and 40%, respectively) and specificity (85% for each) [5]. Under these conditions, pediatric TB carries 14% mortality worldwide [6]. To improve treatment outcomes, we need better biomarkers to diagnose pediatric TB, identify appropriate responses to treatment, and determine which exposed children will develop new infections.

There is growing interest in host-based TB diagnostics, including blood messenger ribonucleic acid (RNA) transcriptional

Received 14 December 2018; editorial decision 26 November 2019; accepted 3 December 2019; published online December 4, 2019.

<sup>a</sup>A. P. and B. B. A. contributed equally to this work.

Correspondence: J. A. Tornheim, MD, MPH, Assistant Professor of Medicine, Department of Medicine, Division of Infectious Diseases, Johns Hopkins University School of Medicine, 600 N. Wolfe St., Phipps 521, Baltimore, MD 21287 (tornheim@jhu.edu).

The Journal of Infectious Diseases® 2020;221:1647–58

© The Author(s) 2019. Published by Oxford University Press for the Infectious Diseases Society of America. All rights reserved. For permissions, e-mail: journals.permissions@oup.com. DOI: 10.1093/infdis/jiz639

profiles [7–14] and techniques to simplify published gene lists into short, “universal” signatures for point-of-care polymerase chain reaction (PCR) testing. This includes sets of 3 diagnostic genes (*GBP5*, *DUSP3*, and *KLF2*) and 5 treatment response genes (*SMARCD3*, *UCP2*, *MAP7D3*, *STT3A*, and *RP11-295G20.2*) [12, 13, 15]. However, comparisons of prior studies have been limited by variable laboratory methods (microarrays, PCR, or RNA-sequencing [RNAseq]), variable statistical methods, few replication studies, and overrepresentation of specific populations. Most data come from Africa and include few children with confirmed TB, yet India represents 27% of global TB compared with 25% in Africa [1]. Before application to clinical diagnosis worldwide, transcriptional profiles need to be replicated in specific target populations, including children and Indian study participants. To evaluate the generalizability of published transcriptional profiles to pediatric TB, we assessed gene expression among Indian children with confirmed TB and exposed controls.

## METHODS

### Study Setting and Population

The Cohort for Tuberculosis Research by the Indo-US Medical Partnership (CTRIUMPh) cohort enrolled adults and children with TB and their household contacts to evaluate factors associated with transmission and progression from exposure to active TB [16]. The CTRIUMPh enrolled participants at 2 sites in India in an academic partnership with Johns Hopkins University. For this study, samples were selected from the Byramjee Jeejeebhoy Government Medical College (BJGMC), a tertiary teaching hospital in Pune, a city of >7 million people in Maharashtra, India. Whole blood was collected in PAXgene tubes from participants with TB throughout treatment and from contacts at enrollment and at 4- and 12-month follow-up visits. Contacts were also tested for TB infection (TBI) by both tuberculin skin testing (TST) and interferon gamma release assays ([IGRA] QuantiFERON-TB Gold In-Tube) at enrollment, 4, and 12 months.

### Participant and Sample Selection

Participants were <15 years of age with TB (“cases”), defined as positive by either Xpert MTB/RIF (6 participants), culture (6 participants), or caseating granulomas on histopathology of extrapulmonary samples (7 participants, not mutually exclusive). All cases completed 6 months of treatment with clinical and microbiological resolution of disease. Cases were matched by age and sex to 2 household contacts (exposed but uninfected controls) (Supplementary Data). Selected controls were negative for both active TB and TBI upon enrollment as determined by negative symptom screen, negative chest radiography, and both negative TST (<5 mm) and IGRA. Repeat TST and IGRA testing occurred at each follow-up visit, and controls were categorized as “converters” if they developed incident TBI, defined

by newly positive TST or IGRA, or “nonconverters” if they did not. Each control was followed for 12 months. None developed active TB during that time.

### Sequencing and Statistical Analysis

Samples were collected from cases at 0, 1, and 6 months of treatment and from controls at 0, 4, and 12 months for this study. Ribonucleic acid extraction was performed using commercially available PAXgene Blood RNA kits according to the manufacturer’s instructions, and RNA was sequenced by Illumina HiSeq 2500 at MedGenome in Bangalore, India to generate an average 170 million 100 base-pair paired-end reads per sample.

Raw data were aligned to the human genome (GRCh38.10) using the Spliced Transcripts Alignment to a Reference (STAR) aligner and annotated using GENCODE [17]. Samples were sequenced to a median of 168.3 million reads (interquartile range [IQR], 149.9–174.1 million reads), with medians of 81.9% (IQR, 77.3%–84.3%) uniquely mapped, 11.1% (IQR, 9.0%–13.5%) multiply mapped, and 7.0% (IQR, 6.0%–8.4%) unmapped. Count data were exported to R for subsequent analysis and subset for protein-coding genes using the biomaRt package. Heatmaps were generated using  $\log_2$ -transformed counts per million (CPM). Genes were excluded from heatmaps if mean  $\log_2$ (CPM) were <2 [18]. Hierarchical clustering used Euclidean distance as a measure of dissimilarity and average linkage for between-cluster separation. Final heatmaps presented gene expression by z-score of standard deviations from mean expression by gene across all samples using ggplot2 [19]. Differential expression analysis by DESeq2 differentiated cases from controls, cases over time, and controls over time [20]. Results were  $\log$  transformed for principal component analysis (PCA) to identify expression-associated factors. DESeq2 was repeated to differentiate cases from controls (“diagnostic gene list”) and cases before and after appropriate TB treatment (“treatment response gene list”) using design matrices including case versus control status, sex, age group (<5, 5–9, 10–14 years), pulmonary versus extrapulmonary disease, incident TBI among controls, and sample month. To minimize impact of interparticipant gene expression variation, longitudinal analyses were performed on paired samples from each participant. All genes with  $\geq 2$ -fold ( $\log_2 \geq 1$ ) differential expression and a Benjamini-Hochberg false discovery rate (FDR) of <0.05 were included in the final diagnostic and treatment response lists.

Modular analysis of functional gene pathway was performed using clusterProfiler [21]. The diagnostic gene list was annotated using Reactome for KEGG pathway according to study group (cases, converters, and nonconverters) with gene size parameter set between 20 and 500. Weighted gene coexpression network analysis was performed on  $\log_2$ -transformed CPM using the WGCNA package. A signed weighted correlation matrix of pairwise Pearson correlations between all genes and all

samples was computed using a soft threshold of  $\beta = 14$  to reach a scale-free topology from which overlap was calculated. The WGCNA dynamic hybrid tree-cut algorithm was used to detect the network modules of coexpressed genes with a minimum module size of 20. Modules were titled with an arbitrary color for module distinction and annotated with the KEGG database and enrichment analysis of each module. Fold enrichment was calculated using the quantitative set analysis for gene expression (qusage) package [22]. Modules with FDR  $P < .05$  were considered significant.

#### Concordance With Published Gene Lists

Diagnostic and treatment response gene lists were compared with other publications to determine concordance of final gene lists between studies. We assessed performance of other gene lists in our data by gene set variation analysis (GSVA) using TBSignatureProfiler [23]. This assessed the discriminating ability of 24 published gene lists for TB diagnosis and 3 published gene lists for treatment response within our dataset. Raw data from the most concordant publications were downloaded from the GEO database (GSE19491 [8], GSE42834 [9], GSE37250 [11], and GSE79362 [13]) as was a large pediatric TB diagnosis dataset (GSE39941 [7]), and 4 treatment response datasets were from 3 studies (GSE40553 [10], GSE89403 [15], GSE31348 [24], and GSE36238 [24]). In each dataset, the area under the receiver operator characteristic curve (AUC), sensitivity, and specificity of our diagnosis and treatment response gene lists were calculated using the pROC package [25]. Discriminating thresholds were identified by Youden's method, and significance was determined by rank-sum test. Two publications provided risk score calculators to generate a single, per-sample value [12, 13]. We calculated both scores for each sample and determined optimal thresholds for each score among Indian children. The Sweeney score was calculated by log transforming counts for each of 3 genes as follows:  $(GBP5 + DUSP3)/2 - KLF2$ . The Zak score was calculated using a published formula to normalize transcripts after conversion from GRCh38.10 to hg19. The AUC, sensitivity, and specificity for these scores was calculated as described above (Supplementary Data).

#### Study Approval

The CTRIUMPh was approved by the institutional review boards of BJGMC, the National Institute of Research in Tuberculosis, and Johns Hopkins University School of Medicine. All participants <18 years old had written informed consent provided by their legal guardians. Written assent for participation was provided by participants who were 8–18 years old.

#### Data Availability

Raw sequencing data analyzed in this manuscript are available from the NCBI sequence read archive under accession code PRJNA588242. Data supporting the findings of this study are

available within the article and its [Supplementary Files](#) or are available from the authors upon request.

## RESULTS

#### Study Participants

The CTRIUMPh enrolled 511 participants from 2014 to 2018 with active TB and 499 household contacts at BJGMC. Of these, 141 participants with active TB (27.6%) were <15 years of age, 16 (11.3%) of whom had confirmed active TB and sufficient follow-up for inclusion (cases). Median age for cases was 9.5 years (IQR, 7–14; range, 3–14), 8 (50.0%) were male, and 8 (50.0%) had extrapulmonary disease (Supplementary Data). Of the 32 age- and sex-matched pediatric household contacts (controls), 13 (40.6%) developed incident TBI. No controls developed active TB disease within 1 year of enrollment. One control (3.1%) with incident TBI progressed to active TB within 2 years of enrollment. All cases and controls were human immunodeficiency virus (HIV) uninfected.

#### Principal Component Analysis

Principal component analysis was performed to determine the factors most associated with differential gene expression among study participants. Visual inspection confirmed sex and age to be significant factors for gene expression (Supplementary Data). Exclusion of sex-associated genes significantly altered both PCA plots and resulting gene lists. Subsequent models included age and sex as covariates along with case versus control status (for diagnosis) and month of treatment (for treatment response).

#### Transcriptome Profiling and Gene Lists

Transcriptome profiling identified significant within-group heterogeneity in gene expression, even after controlling for age and sex (Supplementary Data). Mean within-group expression identified distinct sets of differentially expressed genes between cases at each month during treatment, between controls with incident TBI before and after conversion, and between nonconverting controls at exposure (first sample) and after 12 months (last sample). Mean gene expression levels from enrollment samples were more similar between cases at month 0 and converters before conversion than between cases at month 0 and converters after conversion, suggesting transient changes in gene expression after exposure.

We identified 71 differentially expressed genes between cases and controls with an FDR <0.05 and  $\geq 2$ -fold difference in counts (69 upregulated, 2 downregulated) (Table 1a). Among cases, 1829 genes were differentially expressed between months 0 and 1 of treatment (all downregulated) (Supplementary Data), and 25 were downregulated between months 0 and 6 of treatment (Table 1b). We found marked heterogeneity of differentially expressed genes between comparison groups, with only 2 genes (*DEFA3* and *PRRG4*) common to both the diagnostic and

**Table 1a. Diagnostic Gene List of Differentially Expressed Genes Between Confirmed Pediatric Participants With TB and Exposed Household Contacts at Enrollment, 71 Genes**

<i>AIM2</i>	<i>CD177</i>	<i>ELOVL3</i>	<i>LIPM</i>	<i>PRTN3</i>
<i>ANKRD22</i>	<i>CDCP1</i>	<i>EPB42</i>	<i>LYPD5</i>	<i>PXT1</i>
<i>ANXA8</i>	<i>CLDN18</i>	<i>ERG</i>	<i>MAP6</i>	<i>RAB31L1</i>
<i>APOL4</i>	<i>CLEC5A</i>	<i>FAM26F</i>	<i>METTL7B</i>	<i>SCG5</i>
<i>ARHGEF17</i>	<i>CORIN</i>	<i>FCGR1A</i>	<i>MPO</i>	<i>SEPT4</i>
<i>ARHGEF37</i>	<i>CTSE</i>	<i>FCGR1B</i>	<i>MYZAP<sup>a</sup></i>	<i>SPATC1</i>
<i>AZU1</i>	<i>CTSG</i>	<i>GBP1</i>	<i>NEURL3</i>	<i>TCTEX1D1</i>
<i>BEST3</i>	<i>CXorf36</i>	<i>GBP5</i>	<i>NFXL1</i>	<i>TGM2</i>
<i>BPI</i>	<i>DDIT4L</i>	<i>GBP6</i>	<i>NRCAM<sup>a</sup></i>	<i>TMEM52B</i>
<i>C17orf105</i>	<i>DEFA3</i>	<i>GLDN</i>	<i>NRN1</i>	<i>VNN1</i>
<i>C1QA</i>	<i>DHRS9</i>	<i>GPR84</i>	<i>OLR1</i>	<i>WNT9A</i>
<i>C1QB</i>	<i>DOC2B</i>	<i>HBD</i>	<i>PCOLCE2</i>	
<i>C1QC</i>	<i>DUSP13</i>	<i>HTRA3</i>	<i>PDCC1LG2</i>	
<i>C3orf84</i>	<i>DZIP1L</i>	<i>KLHDC8A</i>	<i>PRRG4</i>	
<i>CAMP</i>	<i>ELANE</i>	<i>LCN2</i>	<i>PRRT4</i>	

Abbreviations: TB, tuberculosis.

<sup>a</sup>Downregulated genes.

**Table 1b. Treatment Response Gene List of Differentially Expressed Genes Between Confirmed Pediatric Participants With TB at the Beginning and End of Successful treatment, 25 Genes**

<i>ACSL1<sup>a</sup></i>	<i>CHI3L1<sup>a</sup></i>	<i>IGF2R<sup>a</sup></i>	<i>MGAM<sup>a</sup></i>	<i>REPS2<sup>a</sup></i>
<i>ADAMTSL4<sup>a</sup></i>	<i>CYP4F3<sup>a</sup></i>	<i>ITGAM<sup>a</sup></i>	<i>NAIP<sup>a</sup></i>	<i>SIPA1L2<sup>a</sup></i>
<i>ATP8B4<sup>a</sup></i>	<i>DEFA3<sup>a</sup></i>	<i>KCNH7<sup>a</sup></i>	<i>NLRC4<sup>a</sup></i>	<i>SORT1<sup>a</sup></i>
<i>B4GALT5<sup>a</sup></i>	<i>DYSF<sup>a</sup></i>	<i>KCNJ15<sup>a</sup></i>	<i>PADI4<sup>a</sup></i>	<i>TLR2<sup>a</sup></i>
<i>CC2D2A<sup>a</sup></i>	<i>FCAF<sup>a</sup></i>	<i>KREMEN1<sup>a</sup></i>	<i>PRRG4<sup>a</sup></i>	<i>WDFY3<sup>a</sup></i>

Abbreviations: TB, tuberculosis.

<sup>a</sup>Downregulated genes.

treatment response gene lists. Both were upregulated among cases compared with controls and downregulated between the start and end of TB treatment.

#### Modular Analysis by Tuberculosis Infection Status

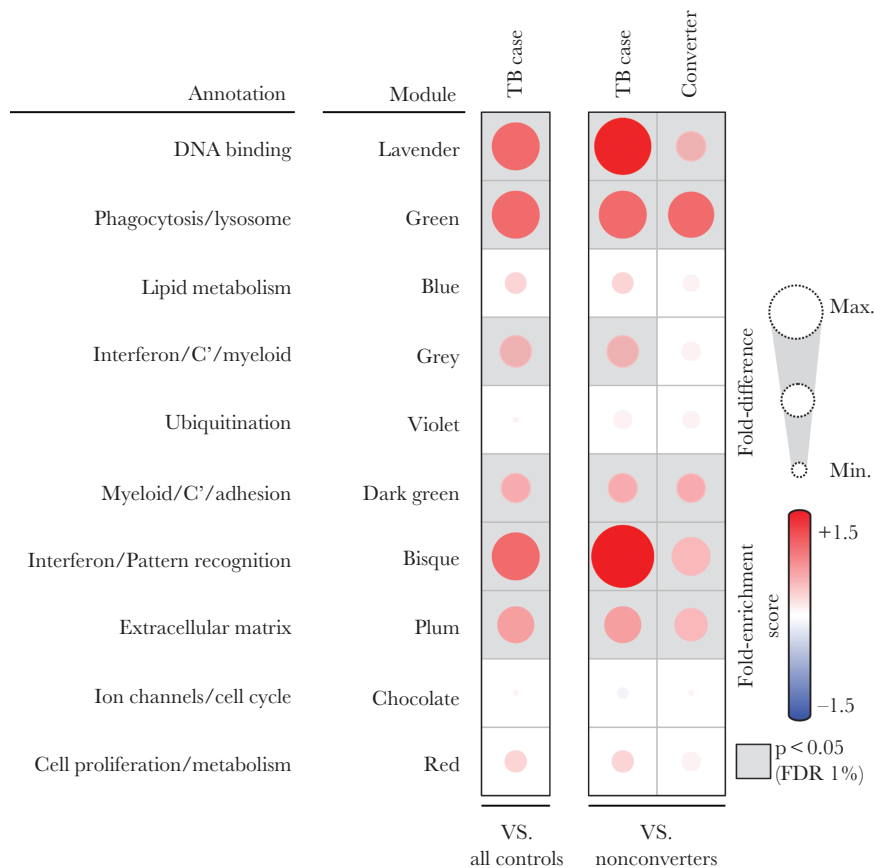
Between-group heterogeneity was confirmed by modular analysis using ingenuity pathway analysis. Specifically, expression differed by disease status for genes in the interferon pattern recognition, deoxyribonucleic acid (DNA) binding, phagocytosis and lysosome, extracellular matrix, interferon myeloid, myeloid cell adhesion, cell proliferation and metabolism, lipid metabolism, ubiquitination, and ion channel modules (Figure 1). Stratification according to incident TBI status revealed more prominent roles for genes in the DNA binding, interferon and pattern recognition, and ubiquitination modules for cases compared with nonconverting controls. In addition, comparison of controls by conversion status found that the interferon and DNA binding modules had lower fold enrichment scores, with increased enrichment of the ubiquitination module.

#### Comparison With Published Gene Lists

Considering the limited overlap of specific genes in our gene lists with prior publications (Figure 2), we compared our gene

lists to 24 published gene lists that identified TB, inflammation, and incident TBI, and to 3 lists that identified treatment response (Figure 3). Gene set variation analysis demonstrated variable performance of published gene lists within our pediatric dataset compared with our gene lists for diagnosis (sensitivity range, 0.38–1.00; specificity range, 0.48–0.93) and treatment response (sensitivity range, 0.70–0.80; specificity range, 0.40–0.80) (Supplementary Data). Despite limited overlap of gene lists between studies, several studies could differentiate cases versus controls and treatment response with statistical significance, although none achieved an AUC of  $\geq 0.8$  (Supplementary Data).

The lists published by Zak et al [13], Sweeney et al [12], and Kaforou et al [11] showed the highest concordance with our diagnostic gene list (Figure 3), with 37.5%, 33.3%, and 25.9%, respectively, of those gene lists differentially expressed in this study. We assessed the performance of our gene lists in those 3 studies, as well as a large pediatric dataset [7] and 3 studies of treatment response (Figure 4, Table 2) [10, 15, 24]. Our 71-gene list differentiated TB cases from controls with an AUC  $> 0.80$  in 2 adult datasets and achieved an AUC of 0.70 (95% confidence interval [CI], 0.63–0.76) in a large African pediatric dataset [7].



**Figure 1.** Modular transcriptional assessment of pediatric tuberculosis from India, according to case and control group. Modular functional enrichment analyses of differentially expressed genes between TB cases and all controls, as well as between TB cases and converters and between converters and nonconverters, were performed as described in the Methods section. Module titles (arbitrary unique colors) are applied to represent specific pathways. The top 10 significant modules are shown along with fold enrichment scores derived from QuSAGE, with red and blue shading indicating modules over- or under-expressed compared with indicated controls. Color intensity and circle size represent the degree of enrichment, compared to the controls. Only modules with fold enrichment scores with  $p < 0.05$  with FDR  $< 1\%$  were considered significant and illustrated.

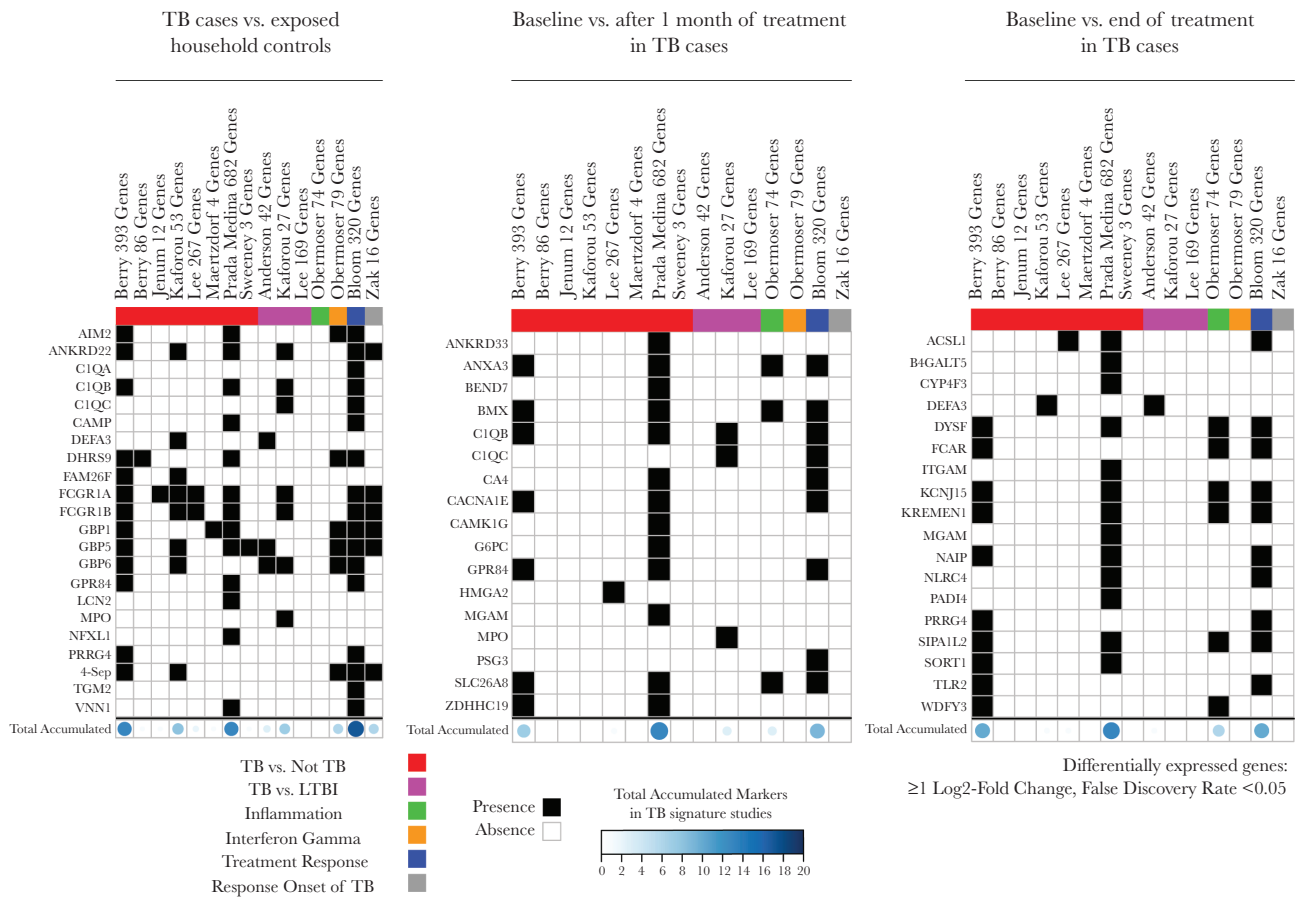
In addition, our 25-gene list identified treatment response with an AUC of 0.90 (95% CI, 0.74–1.00) in one adult dataset and 0.72 (95% CI, 0.85–0.94) in a second adult dataset, but it performed less well in the other 4 datasets assessed.

Two manuscripts provided risk score calculators that we applied to our data (Figure 5, Supplementary Data). Using an optimal cutoff score of 0.33, Sweeney et al [12] scores differentiated pediatric cases from controls in our data with an AUC of 0.76 (95% CI, 0.58–0.93), sensitivity of 0.81 (95% CI, 0.50–0.94), and specificity of 0.67 (95% CI, 0.04–0.93). Sweeney et al [12] also identified treatment response using a cutoff of 0.59 with an AUC of 0.77 (95% CI, 0.54–1.00), sensitivity of 0.80 (95% CI, 0.00–1.00), and specificity of 0.80 (95% CI, 0.20–1.00). Zak et al [13] scores differentiated cases from controls using a cutoff of 0.08 with an AUC of 0.71 (95% CI, 0.54–0.89), sensitivity of 0.63 (95% CI, 0.25–0.81), and specificity of 0.78 (95% CI, 0.33–0.96). Zak et al [13] scores also identified incident TBI among controls using a cutoff of 0.16 with an AUC of 0.83 (95% CI, 0.60–1.00), sensitivity of 0.86 (95% CI, 0.19–1.00), and specificity of 0.82 (95% CI, 0.09–1.00), but not

cases before treatment from after treatment ( $P = .13$ ). Both scores differentiated TBI from TB disease.

## DISCUSSION

Our study of transcriptomic profiles of Indian children with and without TB disease or TBI had several important findings. First, our 71-gene diagnostic gene list and our 25-gene treatment response gene list were derived from a uniquely characterized cohort of “true positive” TB cases and “true negative” age- and sex-matched controls for whom both active and TBI were ruled out. Second, despite significant differences between study groups, we found marked within-group heterogeneity. This highlights the importance of relevant covariates in differential gene expression models, including age, sex, and incident TBI. Third, the genes in our lists had limited overlap with genes reported in previous publications, largely derived from African adults. Gene lists from other publications did not achieve high AUCs among Indian children, but our gene lists achieved high AUCs in 3 other studies. Finally, we calculated thresholds for



**Figure 2.** Concordance between differentially expressed genes found in the present study and gene lists reported by other publications. (A–C) Lists of the differentially expressed genes (DEGs) identified in each indicated comparison are displayed on the Y-axis. Definition of DEGs was  $\geq 2 \log_2$ -fold differential expression and a false discovery rate of  $< 0.05$  between the TB cases (pre-treatment) and exposed household contacts (A), TB cases at baseline (pre-treatment) and at month 1 of therapy (B) and TB cases at baseline and at the end of therapy (C). In the graphs, X-axis indicates the distinct gene lists published previously; colored squares on the top indicate the study design (comparisons tested). Black squares indicate the presence of a given DEG in the published signature. In the bubble plots on the bottom, circle size and color intensity are proportional to the total number of the DEGs detected in a publication. The gene lists with the greatest number of shared genes were the Berry diagnostic list of 393 genes, the Bloom treatment response list of 320 genes, and the Prada-Medina diagnostic list of 682 genes. The published lists with the greatest proportion of shared genes were the Zak risk of progression list of 16 transcripts, the Sweeney diagnostic list of 3 genes, and the Kaforou diagnostic list of 27 genes.

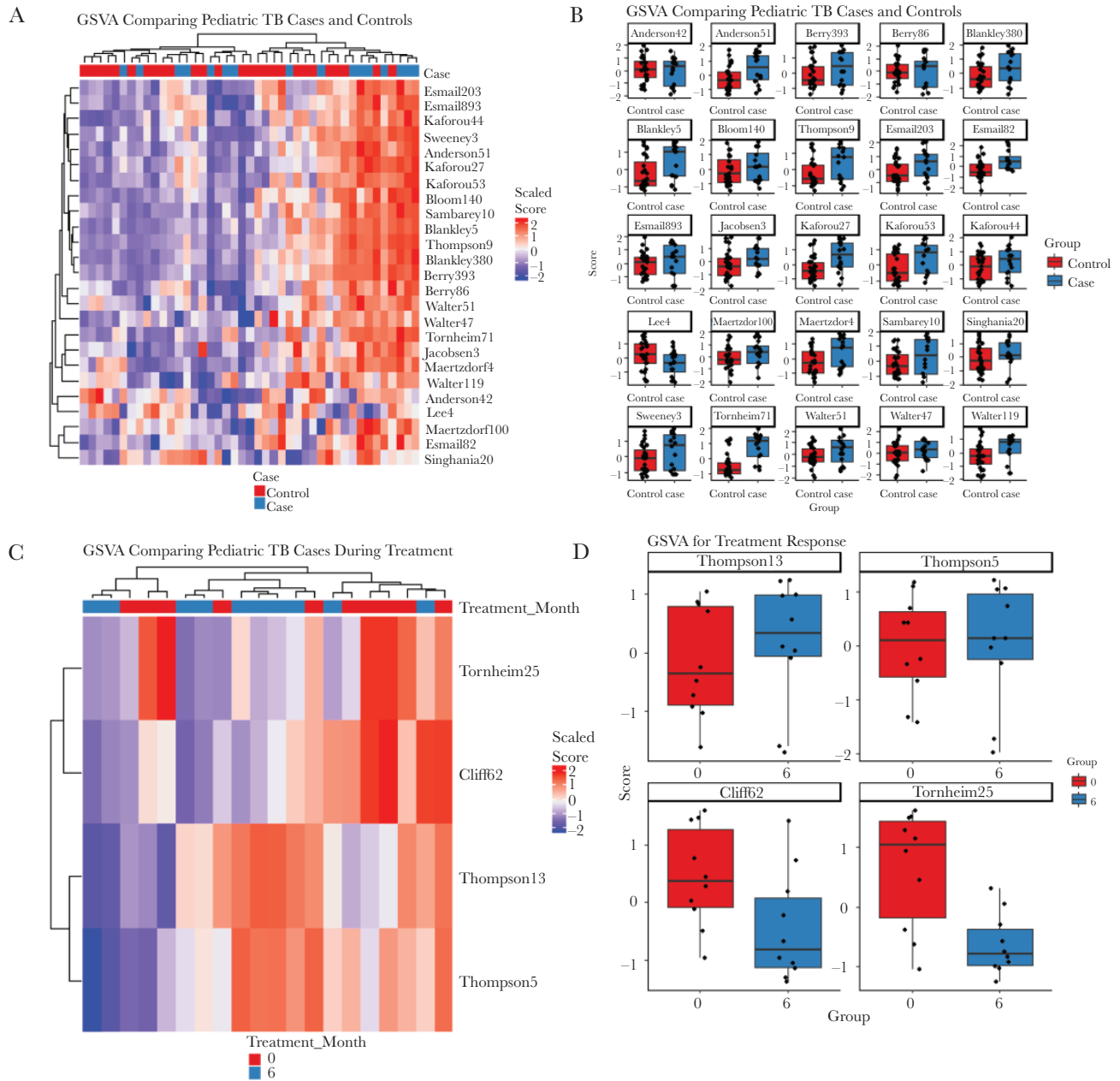
published risk scores to identify TB disease, treatment response, and incident TBI after exposure among Indian children.

Microbiological confirmation is difficult for pediatric TB patients due to lower bacterial burden, inability to produce sputum, and frequent extrapulmonary infection. This creates a dire need for non-sputum-based TB tests that perform better among children than current assays ( $>66\%$  sensitive and  $98\%$  specific) [2–4, 26]. We evaluated a “best-case” scenario for RNAseq among Indian children with TB, using high-depth sequencing among only confirmed cases and TBI-negative controls to maximize accuracy. The limited overlap found with prior publications is likely due in part to the multiplicity of epidemiologic, laboratory, and statistical methods used, including microarrays, RNAseq, and PCR-based tests with variable fold-changes and FDR thresholds between studies. The use of statistical methods to collapse results into shorter, more predictive gene lists for translation

to point-of-care also contributes to the limited overlap of final published gene lists.

It is important to note that age and sex significantly influenced our results. Several genes have known sex-biased expression [27, 28], which impacts neuropsychiatric conditions [29, 30], brain development [31, 32], lung injury [33], inflammatory responses to trauma [34], and immune function, including Fc receptor and interleukin 8 genes [35, 36]. Sex-bias has been reported for TB incidence and mortality, including susceptibility and severity differences between castrated or testosterone-treated animals [37]. X-linked polymorphisms have been described, as have sex-specific responses across immune cell types affecting interferon, interleukin, tumor necrosis factor, Toll-like receptor, and Bacillus Calmette-Guérin (BCG)-specific responses [37–40]. The transition through adolescence into early adulthood, with resulting hormone changes, has also been associated with TB [41, 42], but it has

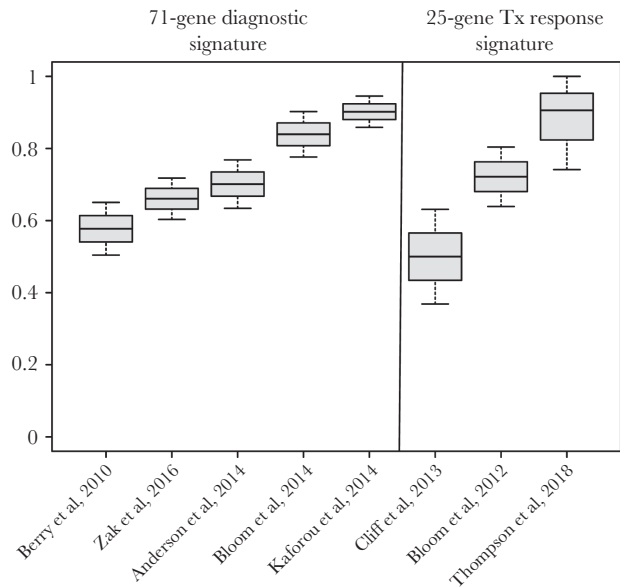
GSVA Comparing Performance of Gene List for Pediatric TB Diagnosis and Treatment Response Compared to Other Published Gene Lists



**Figure 3.** Gene set variation analysis (GSVA) demonstrating performance of gene lists for pediatric tuberculosis diagnosis and treatment response compared with published gene lists. Gene set variation analysis (GSVA), the relative performance of gene lists derived from this Indian pediatric dataset and those from other publications when tested among our samples. The gene lists derived from this study for diagnosis of tuberculosis ("Tornheim71" in A and B) and treatment response ("Tornheim25" in C and D) show the greatest discriminating ability in unit-free GSVA analysis for both diagnosis and treatment. No gene lists achieved the sensitivity and specificity of existing diagnostic tools in our dataset.

not been widely discussed in transcriptomic studies of TB. Given the significant role of sex-specific gene transcription in our data and in other diseases, we propose that age and sex should be included in discovery models of differential gene expression analysis, particularly among children and adolescents.

After controlling for these factors, we identified substantial within- and between-group heterogeneity in gene expression (Figure 2). Some heterogeneity can be explained by genes representing biological pathways, such as GBP1, GBP5, and GBP6 in our diagnostic gene list. Previous manuscripts vary in their inclusion of such genes, but they tend to overrepresent interferon-induced



**Figure 4.** Performance of differentially expressed genes derived from pediatric gene lists among published datasets for diagnosis of pediatric tuberculosis and response to tuberculosis treatment. The list of 71 genes that were identified among Indian children in this study was assessed for its ability to discriminate between cases and controls in four diagnostic datasets for tuberculosis (Left). Area under the receiver operator characteristic curve (AUC) of our gene list exceeded 80% for diagnosis of tuberculosis in two of those datasets. The list of 25 genes that were differentially expressed in response to treatment among Indian children with tuberculosis was assessed for its ability to identify treatment response among 3 studies of adults with tuberculosis. AUC of the 25 gene treatment response list exceeded 80% for treatment response in one dataset.

genes, which may reduce specificity for TB [14]. This is particularly true with TBI diagnosis, for which current tests are frequently discordant [43, 44]. We found little overlap between our diagnostic and treatment response gene lists. In addition, modular analysis found interferon-inducible genes to vary depending on incident TBI among controls. These findings suggest that a “one-size fits all” approach may not work in all populations, particularly among children from TB-endemic areas where controls have frequent TB exposure and may not yet have positive TSTs or IGRAs.

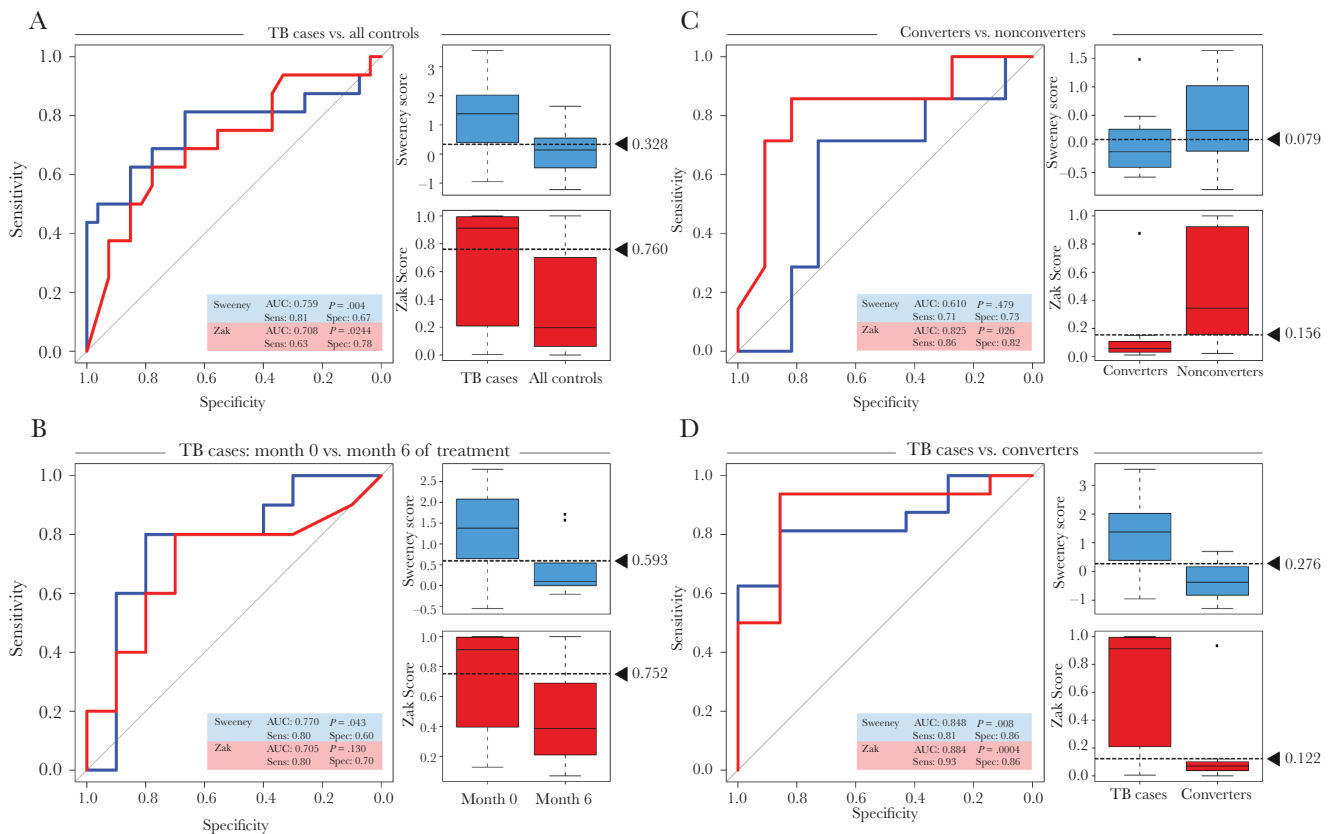
Given limited overlap of genes and variable populations, it is not surprising that lists from other publications did not achieve comparable sensitivity and specificity to existing tests for pediatric TB diagnosis within our data. Likewise, application of our gene lists to other datasets found lower sensitivity and specificity than Xpert MTB/RIF for pediatric TB diagnosis (66% sensitive and 89% specific) [45]. However, our 71-gene list for diagnosis had better sensitivity than smear (15%) and culture (30%–40%) for TB diagnosis in both adult and pediatric datasets (Table 3) [2–4]. Similarly, our 25-gene list for treatment response had better sensitivity to predict treatment outcomes than smear (24%) and culture (40%) in all adult datasets, with similar specificity (84% vs 85%) [5], suggesting a role in pediatric TB management. If these gene lists were developed into point-of-care PCR assays, they may require distinct cutoffs for use in diagnosis and treatment response assessment. We have included optimal thresholds for this population in the [Supplemental Materials](#). Ideally, future assays would improve upon these targets, because low provider confidence in new diagnostics can reduce their impact on empiric treatment rates, particularly in settings with a high pretest probability of TB [46].

This study had several limitations. In Pune, only 56% of public sector patients are microbiologically confirmed, and only 4% are children [47]. To improve specificity for pediatric TB, we only included confirmed cases in this study, resulting in a small sample size. To provide a true negative control group, we selected children who, despite exposure to participants with TB, were negative by both TST and IGRA at the time of enrollment, and we assessed children who were recently exposed, rather than those who were sick with diseases that present similarly to TB. This may not reflect a “real-world” assessment, because TBI rates among Indian adolescents is common, approaching 50%–60% of household contacts in some studies [48, 49]. Tuberculosis is common in India, with a national incidence of 217 per 100 000 person-years [47]. This offers many opportunities for TB exposure,

**Table 2.** Area Under the Curve, Sensitivity, and Specificity of Pediatric Gene Lists in Published Diagnostic and Treatment Response Studies for Tuberculosis

Study	AUC of Pediatric Gene Lists in Other Study's Dataset	Sensitivity of Pediatric Gene Lists in Other Study's Dataset	Specificity of Pediatric Gene Lists in Other Study's Dataset
<b>71-Gene List for Diagnosis of Tuberculosis</b>			
Kaforou et al [11]	0.83 (0.77–0.90)	0.85 (0.80–0.88)	0.85 (0.81–0.89)
Bloom et al [9]	0.90 (0.85–0.94)	0.50 (0.50–0.50)	0.50 (0.50–0.50)
Anderson et al [7]	0.70 (0.63–0.76)	0.79 (0.71–0.85)	0.79 (0.72–0.86)
Zak et al [13]	0.66 (0.60–0.71)	0.63 (0.58–0.69)	0.64 (0.59–0.69)
Berry et al [8]	0.57 (0.50–0.65)	0.56 (0.51–0.63)	0.56 (0.51–0.62)
<b>25-Gene List for Response to Tuberculosis Treatment</b>			
Thompson et al [15]	0.90 (0.74–1.00)	0.86 (0.70–1.00)	0.84 (0.70–1.00)
Bloom et al [10]	0.72 (0.63–0.80)	0.69 (0.61–0.76)	0.69 (0.62–0.77)
Cliff et al [24]	0.50 (0.36–0.63)	0.49 (0.37–0.63)	0.50 (0.37–0.60)

Abbreviations: AUC, area under the receiver operator characteristic curve.



**Figure 5.** Performance of published risk scores among Indian children, by disease state. Receiver operator characteristic (ROC) curves and box plots assessing the performance of the Sweeney [12] and Zak [13] scores among Indian children, by disease status. Area under the ROC curve (AUC), sensitivity (“Sens”), specificity (“Spec”), and Wilcoxon rank-sum *p*-value are indicated for each score, and box plots identify the best threshold by Youden’s method to differentiate patients by each score (dashed lines). Plots indicate the ability of each score to differentiate active TB from all household contacts (A), to identify response to TB treatment (B), to predict the development of incident TBI after exposure among TST and IGRA negative children (C), and to differentiate active TB from TBI (D).

not just in the home, but also at work, school, or in transit [50]. As a result, our results may not be generalizable to low-burden regions with less-frequent TB exposure. It is also possible that our use of a composite definition of TBI (either TST or IGRA positive) identified a heterogeneous assessment of new infection. Disagreement between TST and IGRA has been well documented, particularly among people with HIV and children for whom TST is confounded by BCG [43, 44]. Finally, because it is harder to confirm TB among younger children, the median age in this study was 9.5 years, with only 6 children  $\leq 5$  years. This limits generalizability of our findings to young children and may have overrepresented the significance of age and sex on puberty-associated gene transcription.

## CONCLUSIONS

The expansion of sequencing technology worldwide offers the opportunity for host-derived biomarkers to improve insensitive microbiological techniques. Moving forward, it will be important to standardize transcriptional profiling models to incorporate age, sex, and TBI status to identify the genes

that are most specific to pediatric TB. To determine the best applications of gene expression to pediatric TB diagnosis, future studies will need to include a broad array of global populations, including children with HIV, children  $< 5$  years old, children with multidrug-resistant TB, children with “probable” and “possible” TB, and those who fail to improve with TB treatment.

## Supplementary Data

Supplementary materials are available at *The Journal of Infectious Diseases* online. Consisting of data provided by the authors to benefit the reader, the posted materials are not copyedited and are the sole responsibility of the authors, so questions or comments should be addressed to the corresponding author.

**Supplement 1.** Additional figures and tables.

**Supplement 2.** Differential gene expression, by condition.

## Notes

**Acknowledgments.** We thank Evan Johnson for the development of and his guidance using the TB Signature Profiler software package used in these analyses [23].

**Author contributions.** J. A. T. contributed to study design, sample processing, data analysis, and manuscript preparation. A. K. M. and A. N. G. contributed to study design, data analysis, and manuscript preparation. M. P. and N. G. contributed to study design, data collection, data analysis, and manuscript preparation. A. T. L. Q., K. F. F., and B. B. A. contributed to data analysis and manuscript preparation. A. K. contributed to data collection and manuscript preparation. V. K., U. B., S. S., R. R., and N. P. contributed to study design, sample processing, and data analysis. S. V. B. Y. S. and C. V. contributed to data collection and analysis. L. E. H. contributed to study design, sample management, and manuscript preparation. V. M. supervised study activities and contributed to study design, data collection, and manuscript preparation. A. P. and A. G. supervised study design, data collection, analysis, and manuscript preparation.

**Disclaimer.** The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Financial support.** This work was funded by the US National Institutes of Health (NIH)/Indian Department of Biotechnology (DBT) RePORT India Consortium. This project has been funded in whole or in part with Federal funds from the Government of India's DBT, the Indian Council of Medical Research, the NIH, National Institute of Allergy and Infectious Diseases (NIAID), Office of AIDS Research, and distributed in part by CRDF Global. This work was also supported by the Wellcome Trust/DBT India Alliance Margdarshi Fellowship (Grant Number IA/M/15/1/502023; to A. P.). J. A. T. was supported by NIH NIAID Grant K23AI135102. This work was additionally supported by the NIH Office of the Director, Fogarty International Center, Office of AIDS Research, National Cancer Center, National Heart, Blood, and Lung Institute, and the NIH Office of Research for Women's Health through the Fogarty Global Health Fellows Program Consortium comprised of the University of North Carolina, John Hopkins, Morehouse, and Tulane (R25TW009340), the Johns Hopkins University School of Medicine Clinician Scientist Career Development Award, and NIH NIAID (R21AI122922). Additional support came from NIH/NIAID (UM1AI069465), the Fogarty International Center BJGMC-JHU HIV-TB Program (D43TW009574), the Ujala Foundation, the Gilead Foundation, the Wyncote Foundation, and Persistent Systems.

**Potential conflicts of interest.** All authors: No reported conflicts of interest. All authors have submitted the ICMJE Form for Disclosure of Potential Conflicts of Interest.

## References

1. World Health Organization. Global Tuberculosis Report 2018. Geneva: World Health Organization; 2018: pp 277.
2. Swaminathan S, Rekha B. Pediatric tuberculosis: global overview and challenges. *Clin Infect Dis* 2010; 50 (Suppl 3):S184–94.
3. Starke JR, Taylor-Watts KT. Tuberculosis in the pediatric population of Houston, Texas. *Pediatrics* 1989; 84:28–35.
4. Graham SM, Cuevas LE, Jean-Philippe P, et al. Clinical case definitions for classification of intrathoracic tuberculosis in children: an update. *Clin Infect Dis* 2015; 61(Suppl 3):S179–87.
5. Horne DJ, Royce SE, Gooze L, et al. Sputum monitoring during tuberculosis treatment for predicting outcome: systematic review and meta-analysis. *Lancet Infect Dis* 2010; 10:387–94.
6. World Health Organization. Roadmap for childhood tuberculosis. Geneva, Switzerland: World Health Organization; 2013: pp 44.
7. Anderson ST, Kaforou M, Brent AJ, et al. Diagnosis of childhood tuberculosis and host RNA expression in Africa. *N Engl J Med* 2014; 370:1712–23.
8. Berry MP, Graham CM, McNab FW, et al. An interferon-inducible neutrophil-driven blood transcriptional signature in human tuberculosis. *Nature* 2010; 466:973–7.
9. Bloom CI, Graham CM, Berry MP, et al. Transcriptional blood signatures distinguish pulmonary tuberculosis, pulmonary sarcoidosis, pneumonias and lung cancers. *PLoS One* 2013; 8:e70630.
10. Bloom CI, Graham CM, Berry MP, et al. Detectable changes in the blood transcriptome are present after two weeks of antituberculosis therapy. *PLoS One* 2012; 7:e46191.
11. Kaforou M, Wright VJ, Oni T, et al. Detection of tuberculosis in HIV-infected and -uninfected African adults using whole blood RNA expression signatures: a case-control study. *PLoS Med* 2013; 10:e1001538.
12. Sweeney TE, Braviak L, Tato CM, Khatri P. Genome-wide expression for diagnosis of pulmonary tuberculosis: a multicohort analysis. *Lancet Respir Med* 2016; 4:213–24.
13. Zak DE, Penn-Nicholson A, Scriba TJ, et al.; ACS and GC6-74 cohort study groups. A blood RNA signature for tuberculosis disease risk: a prospective cohort study. *Lancet* 2016; 387:2312–22.
14. Singhania A, Verma R, Graham CM, et al. A modular transcriptional signature identifies phenotypic heterogeneity of human tuberculosis infection. *Nat Commun* 2018; 9:2308.
15. Thompson EG, Du Y, Malherbe ST, et al.; Catalysis TB-Biomarker Consortium. Host blood RNA signatures predict the outcome of tuberculosis treatment. *Tuberculosis* 2017; 107:48–58.
16. Gupte A, Padmapriyadarsini C, Mave V, et al.; CTRIUMPH Study Team. Cohort for Tuberculosis Research by the Indo-US Medical Partnership (CTRIUMPH): protocol for a multicentric prospective observational study. *BMJ Open* 2016; 6:e010542.
17. Dobin A, Davis CA, Schlesinger F, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013; 29:15–21.

18. McCarthy DJ, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res* **2012**; 40:4288–97.
19. Wickham H. *ggplot2: elegant graphics for data analysis*. Springer-Verlag, New York: Springer; **2016**.
20. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* **2014**; 15:550.
21. Yu G, Wang LG, Han Y, He QY. clusterProfiler: an R package for comparing biological themes among gene clusters. *OMICS* **2012**; 16:284–7.
22. Yaari G, Bolen CR, Thakar J, Kleinstein SH. Quantitative set analysis for gene expression: a method to quantify gene set differential expression including gene-gene correlations. *Nucleic Acids Res* **2013**; 41:e170.
23. Odom A, Jenkins Y, Zhao WE. The tbsignatureprofiler: a novel R package for comparing tuberculosis gene expression signatures. abstract presented at the annual biomedical research conference for minority students (ABRCMS). 13–16 November 2019, Anaheim, California. <http://sites.bu.edu/briteru/files/2019/11/OdomPosterABRCMS2019-Final.pdf>. Accessed 21 December 2019.
24. Cliff JM, Lee JS, Constantinou N, et al. Distinct phases of blood gene expression pattern through tuberculosis treatment reflect modulation of the humoral immune response. *J Infect Dis* **2013**; 207:18–29.
25. Youden WJ. Index for rating diagnostic tests. *Cancer* **1950**; 3:32–5.
26. World Health Organization. High-priority target product profiles for new tuberculosis diagnostics: report of a consensus meeting. Geneva: World Health Organization; **2014**: pp 98.
27. Ung CL, Lam SH, Zhang X, Johnson WE. Inverted expression profiles of sex-biased genes in response to toxicant perturbations and diseases. *PLoS One* **2013**; 8:e56668.
28. Gershoni M, Pietrokovski S. The landscape of sex-differential transcriptome and its consequent selection in human adults. *BMC Biol* **2017**; 15:7.
29. Labonté B, Engmann O, Purushothaman I, et al. Sex-specific transcriptional signatures in human depression. *Nat Med* **2017**; 23:1102–11.
30. Daskalakis NP, Cohen H, Cai G, Buxbaum JD, Yehuda R. Expression profiling associates blood and brain glucocorticoid receptor signaling with trauma-related individual differences in both sexes. *Proc Natl Acad Sci U S A* **2014**; 111:13529–34.
31. Ziats MR, Rennert OM. Sex-biased gene expression in the developing brain: implications for autism spectrum disorders. *Mol Autism* **2013**; 4:10.
32. Xu HW, Wang F, Liu Y, Gelernter J, Zhang H. Sex-biased methylome and transcriptome in human prefrontal cortex. *Hum Mol Genet* **2014**; 23:11.
33. Lingappan K, Srinivasan C, Jiang W, Wang L, Courouclis XI, Moorthy B. Analysis of the transcriptome in hyperoxic lung injury and sex-specific alterations in gene expression. *PLoS One* **2014**; 9:e101581.
34. Lopez ME, Efron PA, Ozrazgat-Baslanti T, et al. Sex-based differences in the genomic response, innate immunity, organ dysfunction, and clinical outcomes after severe blunt traumatic injury and hemorrhagic shock. *J Trauma Acute Care Surg* **2016**; 81:8.
35. Marttila SJ, Jylhävä J, Nevalainen T, et al. Transcriptional analysis reveals gender-specific changes in the aging of the human immune system. *PLoS One* **2013**; 8:e66229.
36. Nevalainen TK, Kananen L, Marttila S, et al. Transcriptomic and epigenetic analyses reveal a gender difference in aging-associated inflammation: the Vitality 90+ study. *Age (Dordr)* **2015**; 37:9814.
37. Nhamoyebonde S, Leslie A. Biological differences between the sexes and susceptibility to tuberculosis. *J Infect Dis* **2014**; 209(Suppl 3):S100–6.
38. Fish EN. The X-files in immunity: sex-based differences predispose immune responses. *Nat Rev Immunol* **2008**; 8:737–44.
39. Muenchhoff M, Goulder PJ. Sex differences in pediatric infectious diseases. *J Infect Dis* **2014**; 209 (Suppl 3):S120–6.
40. Gabriel G, Arck PC. Sex, immunity and influenza. *J Infect Dis* **2014**; 209 (Suppl 3):S93–9.
41. Garenne M, Lafon M. Sexist diseases. *Perspect Biol Med* **1998**; 41:176–89.
42. Garenne M. Demographic evidence of sex differences in vulnerability to infectious diseases. *J Infect Dis* **2015**; 211:331–2.
43. Ayubi E, Doosti-Irani A, Sanjari Moghaddam A, Sani M, Nazarzadeh M, Mostafavi E. The clinical usefulness of tuberculin skin test versus interferon-gamma release assays for diagnosis of latent tuberculosis in HIV patients: a meta-analysis. *PLoS One* **2016**; 11:e0161983.
44. Grinsdale JA, Islam S, Tran OC, Ho CS, Kawamura LM, Higashi JM. Interferon-gamma release assays and pediatric public health tuberculosis screening: the San Francisco Program Experience 2005 to 2008. *J Pediatric Infect Dis Soc* **2016**; 5:122–30.
45. Walter ND, Reves R, Davis JL. Blood transcriptional signatures for tuberculosis diagnosis: a glass half-empty perspective. *Lancet Respir Med* **2016**; 4:e28.
46. Theron G, Zijenah L, Chanda D, et al.; TB-NEAT team. Feasibility, accuracy, and clinical effect of point-of-care Xpert MTB/RIF testing for tuberculosis in primary-care settings in Africa: a multicentre, randomised, controlled trial. *Lancet* **2014**; 383:424–35.
47. Revised National Tuberculosis Control Programme. Guidelines on programmatic management of drug-resistant Tuberculosis in India 2017. New Delhi, India: Revised National Tuberculosis Control Programme, Central TB Division, **2017**: pp 320.

48. Sharma SK, Vashishtha R, Chauhan LS, Sreenivas V, Seth D. Comparison of TST and IGRA in diagnosis of latent tuberculosis infection in a high TB-Burden setting. *PLoS One* **2017**; 12:e0169539.
49. Uppada DR, Selvam S, Jesuraj N, et al.; TB Trials Study Group. The tuberculin skin test in school going adolescents in South India: associations of socio-demographic and clinical characteristics with TST positivity and non-response. *BMC Infect Dis* **2014**; 14:571.
50. Andrews JR, Morrow C, Walensky RP, Wood R. Integrating social contact and environmental data in evaluating tuberculosis transmission in a South African township. *J Infect Dis* **2014**; 210:597–603.