

1 **Cytomegalovirus infection is a risk factor for TB disease in**  
2 **Infants**

3

4 Julius Muller<sup>1</sup>, Rachel Tanner<sup>1</sup>, Magali Matsumiya<sup>1</sup>, Margaret A. Snowden<sup>3</sup>,  
5 Bernard Landry<sup>3</sup>, Iman Satti<sup>1</sup>, Stephanie A. Harris<sup>1</sup>, Matthew K. O'Shea<sup>1</sup>, Lisa  
6 Stockdale<sup>1,2</sup>, Leanne Marsay<sup>1</sup> Agnieszka Chomka<sup>1,4</sup>, Rachel Harrington-  
7 Kandt<sup>1</sup>, Zita-Rose Manjaly Thomas<sup>1</sup>, Vivek Naranbhai<sup>5</sup>, Elena Stylianou<sup>1</sup>,  
8 Stanley Kimbung Mbandi<sup>6</sup>, Mark Hatherill<sup>6</sup>, Gregory Hussey<sup>6</sup>, Hassan  
9 Mahomed<sup>6</sup>, Michele Tameris<sup>6</sup>, J. Bruce McClain<sup>3</sup>, Thomas G. Evans<sup>3</sup>, Willem  
10 A. Hanekom<sup>6</sup>, Thomas J. Scriba<sup>6</sup>, Helen McShane<sup>1</sup> Helen A. Fletcher<sup>1,2</sup>

11

12 <sup>1</sup>The Jenner Institute, University of Oxford, Oxford, UK, OX3 7DQ; <sup>2</sup>London  
13 School of Hygiene & Tropical Medicine, London, UK, W1CE7HT; <sup>3</sup>Aeras,  
14 Rockville, USA, MD 20850; <sup>4</sup>The Kennedy Institute, University of Oxford,  
15 Oxford, UK, OX3 7LF; <sup>5</sup>Wellcome Trust Centre for Human Genetics, University  
16 of Oxford, Oxford, UK, OX37B; <sup>6</sup>South African Tuberculosis Vaccine Initiative,  
17 Institute of Infectious Disease and Molecular Medicine & Division of  
18 Immunology, Department of Pathology, University of Cape Town, South Africa,  
19 7935.

20

1

## 2 ABSTRACT

3 Immune activation is associated with increased risk of tuberculosis (TB)  
4 disease in infants. We performed a case–control analysis to identify drivers of  
5 immune activation and disease risk. Among 49 infants who developed TB  
6 disease over the first two years of life, and 129 matched controls who remained  
7 healthy, we found the cytomegalovirus (CMV) stimulated IFN- $\gamma$  response at age  
8 4-6 months to be associated with CD8+ T-cell activation (Spearman's rho,  
9  $p=6 \times 10^{-8}$ ). A CMV specific IFN- $\gamma$  response was also associated with increased  
10 risk of developing TB disease (Conditional Logistic Regression,  $p=0.043$ , OR  
11 2.2, 95% CI 1.02-4.83), and shorter time to TB diagnosis (Log Rank Mantel-  
12 Cox  $p=0.037$ ). CMV positive infants who developed TB disease had lower  
13 expression of natural killer cell associated gene signatures and a lower  
14 frequency of CD3-CD4-CD8- lymphocytes. We identified transcriptional  
15 signatures predictive of risk of TB disease among CMV ELISpot positive  
16 (AUROC 0.98, accuracy 92.57%) and negative (AUROC 0.9, accuracy 79.3%)  
17 infants; the CMV negative signature validated in an independent infant study  
18 (AUROC 0.71, accuracy 63.9%). Understanding and controlling the microbial  
19 drivers of T cell activation, such as CMV, could guide new strategies for  
20 prevention of TB disease in infants.

21

## 1 INTRODUCTION

2 There are an estimated 1 million cases of childhood tuberculosis (TB) each year  
3 and in 2015, 210,000 children died of TB <sup>5</sup>. Children with TB are difficult to  
4 diagnose and treat and are at risk of severe disease <sup>6</sup>. The need for improved  
5 strategies to control childhood TB has led to studies to identify risk factors for  
6 TB disease in children. The longitudinal data collected during infant TB vaccine  
7 efficacy trials in South Africa with the vaccines BCG and Modified Vaccinia  
8 Virus Ankara expressing Antigen 85A (MVA85A) <sup>1,7</sup> have enabled the  
9 identification of correlates of TB disease risk in infants <sup>3,8</sup>. Using samples  
10 collected at enrolment from infants into the MVA85A efficacy trial in the Western  
11 Cape Province of South Africa <sup>1</sup>, we reported that CD4+ T-cell activation in 4-6  
12 month old infants and adolescents, measured as HLA-DR expression, was  
13 associated with increased risk of TB disease over the following 3 years of life <sup>8</sup>.  
14 The consequences of chronic T-cell activation are well-described in HIV and  
15 include risk of acquisition of infection, risk of progression from infection to  
16 disease, increased risk of non-communicable disease and aging <sup>10-14</sup> .  
17 Recognition that T-cell activation is a feature of HIV immunopathogenesis has  
18 guided the development of new interventions for management of HIV <sup>15</sup>.  
19 Improved understanding of the causes and impact of T-cell activation on TB  
20 disease risk, could transform future approaches to protect against TB. Such  
21 approaches could include the use of vaccines, antibiotics or antivirals to reduce  
22 the burden of chronic microbial drivers of T cell activation.  
23  
24 Sustained T-cell activation and dysfunction of antigen specific T-cells can result  
25 from chronic exposure to antigen from persistent viral or bacterial infections <sup>9</sup>.

1 Human cytomegalovirus (CMV) and Epstein-barr virus (EBV) are known drivers  
2 of T-cell activation <sup>11,14</sup> and with more than 95% seroprevalence in adults, South  
3 Africa is amongst the highest CMV burden countries in the world (reviewed by  
4 Adland *et. al.*, <sup>16</sup>). Recent epidemiological evidence supports a role for CMV in  
5 the aetiology of TB <sup>17-21</sup>. Here, we report results from a case–control analysis  
6 of infants from the MVA85A efficacy trial <sup>1</sup> which aimed to understand the  
7 drivers of immune activation associated with TB risk in infants.

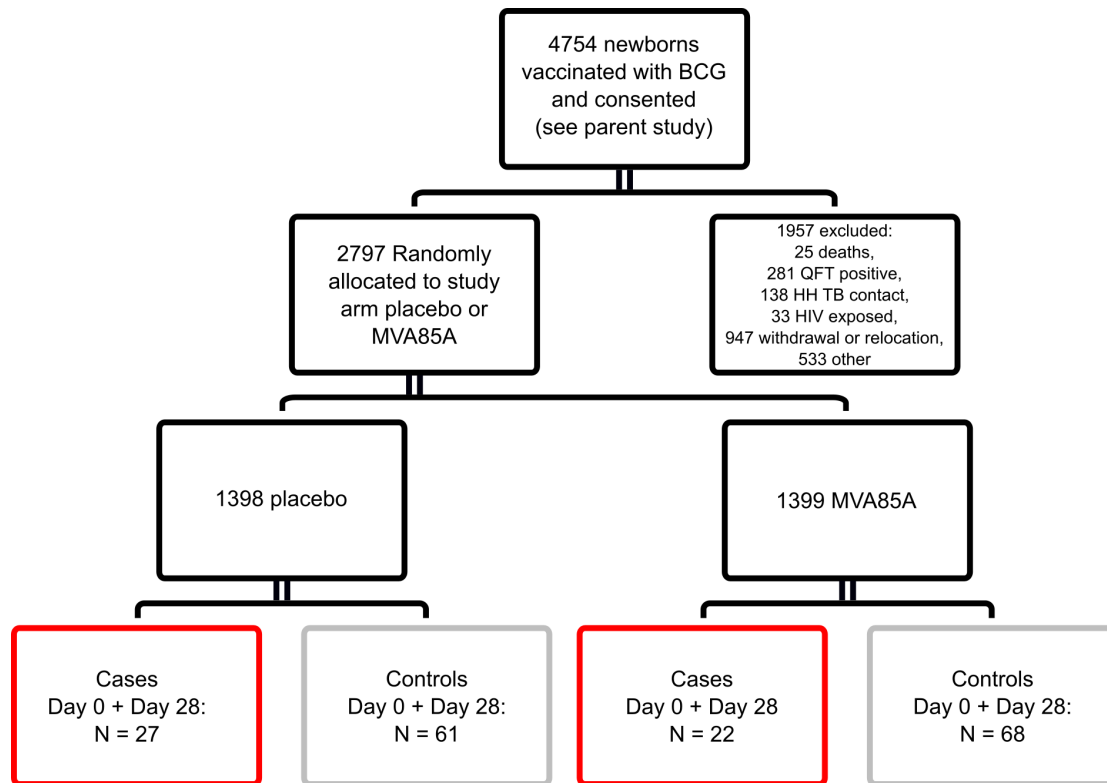
8

## 9 RESULTS

### 10 **Activated CD8+ T-cells are associated with TB disease risk and correlate** 11 **with cytomegalovirus (CMV)-specific IFN- $\gamma$ response**

12 BCG-vaccinated case and control infants who were enrolled in an efficacy trial  
13 of the candidate TB vaccine MVA85A were included in this study <sup>1</sup> (Figure 1).  
14 HIV and *M.tb* uninfected infants without known TB exposure were randomized  
15 at 16-24 weeks of age to receive a single intradermal dose of MVA85A or  
16 placebo (Candin™, a candida skin test antigen)<sup>1</sup>. In contrast to our previous  
17 analysis which focused on the sample collected at enrolment (Day 0), samples  
18 collected from two time points, Day 0 (D0) and Day 28 (D28) following MVA85A  
19 or placebo were combined for testing of both cellular parameters and  
20 transcriptional signatures associated with risk of TB disease. We have shown  
21 that the correlates of risk measured in our previous study were not affected by  
22 intervention (MVA85A or placebo) <sup>8</sup>. Only infants for whom a sample was  
23 available at both D0 and D28 were included in the analysis (Figure 1).  
24 Previously, we identified an association between TB disease risk over the first  
25 two years of life and the frequency of activated HLA-DR+ CD4+ T-cells at age

1 4-6 months <sup>8</sup>. We also found that the magnitude of BCG-specific IFN- $\gamma$   
2 expressing cells and levels of anti-Ag85A IgG were associated with reduced  
3 risk of TB disease <sup>8</sup>. In our present analysis, with combined Day 0 and Day 28  
4 samples from 49 cases and 129 controls (Figure 1), we confirmed our previous  
5 finding of TB risk associating with HLA-DR+ CD4+, BCG-specific IFN- $\gamma$  and  
6 anti-Ag85A IgG (Table 1). In addition, we found frequencies of HLA-DR+ CD8+  
7 T-cells to be associated with TB risk (conditional logistic regression (CLR) OR  
8 1.34, 95% CI 1.08-1.67,  $p = 0.008$ , FDR 0.092) (Table 1). The magnitude of  
9 PPD-stimulated IFN- $\gamma$ -expressing cells measured by ELISpot were also  
10 associated with reduced TB risk (OR 0.71, 95% CI 0.51-0.98,  $p = 0.037$ , FDR  
11 0.2) (Table 1).  
12



1

2

Figure 1. **Study design for immune correlates analysis.** Infants who were enrolled in an efficacy trial of the candidate TB vaccine MVA85A were included in this study<sup>1</sup>. Infants were randomized at 16-24 weeks of age to receive a single intradermal dose of MVA85A or placebo (Candin™, a candida skin test antigen)<sup>1</sup>. Boxes indicate the number of case infant (red) or control infant (grey) samples available for combined Day 0 (D0) and Day 28 (D28) analysis. Analysis was restricted to infants where a frozen PBMC sample was available, live cells in PBMC were >50% (or PHA IFN- $\gamma$  ELISPOT  $\geq$ 1000 SFC/million) and to infants where a sample was available for analysis from both the D0 and D28 time points. Control infants were excluded if the corresponding matched case was not in the analysis.

3

4

1 Table 1. Conditional Logistic Regression of combined Day 0 plus Day 28 infant

2 samples

Estimated Odds-Ratio (OR) of TB disease from a Conditional Logistic Regression of D0 + D28 immunological variable							
Quantitative variable	N	Case	Est OR*	95% CI	P value	FDR value	AUROC
CD3+ T cell	329	88	0.98	0.75-1.28	0.86	0.94	
CD4+ T cell	329	88	0.79	0.6-1.03	0.087	0.29	
HLA-DR+ CD4 T cell	329	88	1.59	1.3-2	<0.0001	0.002	0.64
CD4+CD25+CD127-	329	88	0.84	0.62-1.12	0.23	0.59	
CD8+ T cell	329	88	1.23	0.98-1.53	0.074	0.28	
HLA-DR+ CD8+ T cell	329	88	1.34	1.08-1.67	0.008	0.092	0.59
CD14+CD16+	329	88	0.93	0.73-1.18	0.57	0.82	
CD14+CD16-	329	88	0.99	0.78-1.27	0.96	0.96	
CD19+ B cell	328	88	1.13	0.87-1.4	0.39	0.69	
BCG MGIA	146	50	0.82	0.52-1.29	0.39	0.69	
85A ELISpot	324	91	0.87	0.63-1.2	0.39	0.69	
BCG ELISpot	333	94	0.75	0.57-0.98	0.036	0.2	0.56
PPD ELISpot	270	79	0.71	0.51-0.98	0.037	0.2	0.57
TB10.3/10.4 ELISpot	270	79	0.96	0.74-1.23	0.72	0.83	
EBV ELISpot	290	84	1.05	0.82-1.35	0.7	0.83	

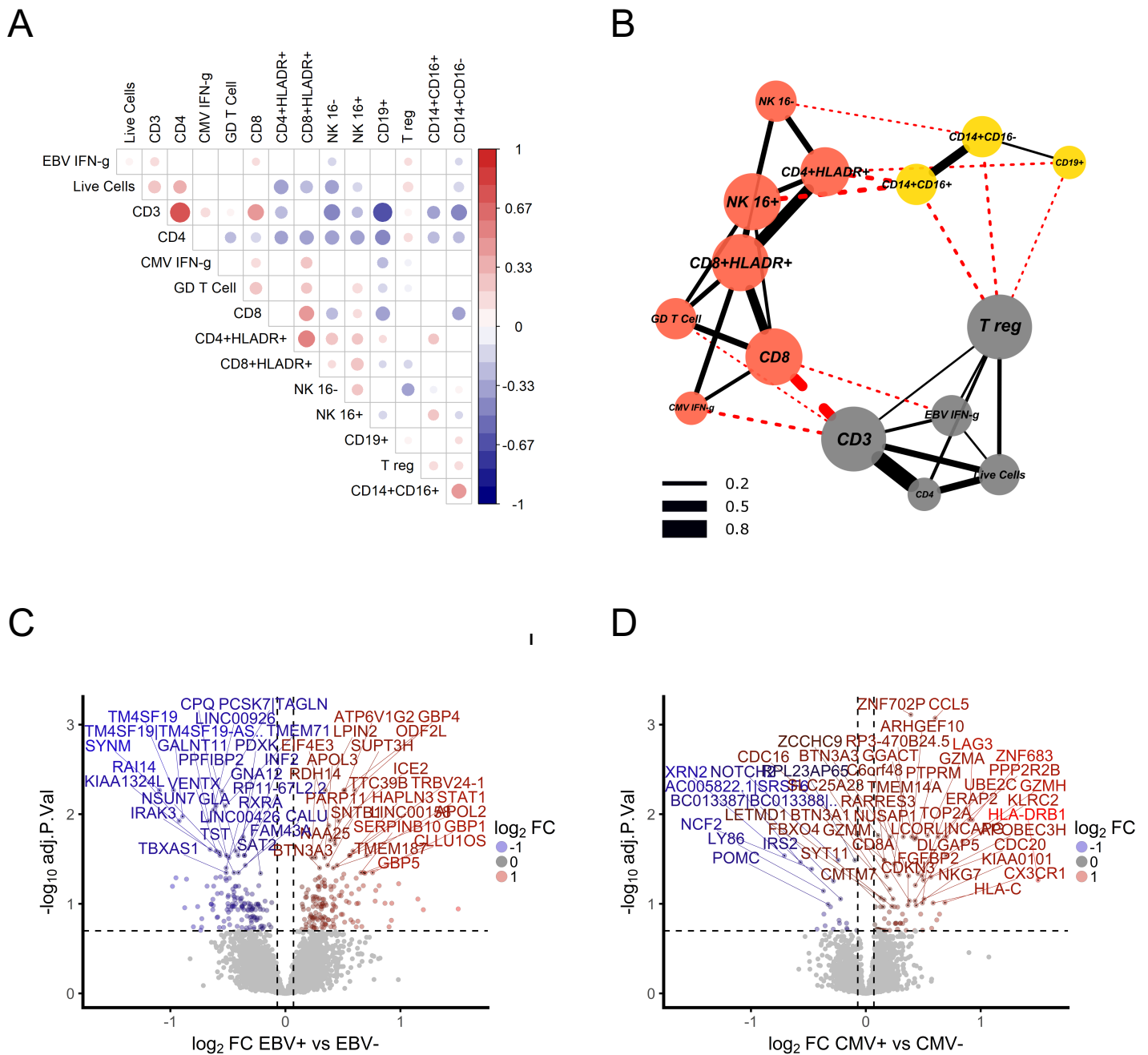
CMV ELISpot	319	91	1.01	0.8-1.23	0.9	0.94	
FLU ELISpot	289	84	1.08	0.85-1.37	0.55	0.82	
GAM.DEL (putative)	329	88	1.16	0.91-1.5	0.23	0.59	
NK 16neg (putative)	329	88	0.96	0.63-1.4	0.63	0.82	
NK 16pos (putative)	329	88	0.86	0.66-1.12	0.26	0.6	
CD14+CD16+ /CD3+	329	88	0.92	0.72-1.18	0.52	0.82	
CD14+CD16- /CD3+	329	88	0.96	0.74-1.22	0.68	0.83	
Anti-Ag85A IgG	297	87	0.79	0.63-0.99	0.043	0.2	0.59

1

2



1 We analyzed CMV and EBV-specific IFN- $\gamma$  ELISpot responses for evidence of  
2 an association between viral infection and T-cell activation in our infant cohort.  
3 Frequencies of activated CD8+ T-cells correlated with the magnitude of the  
4 CMV-specific IFN- $\gamma$  ELISpot response (Spearman's rho  $p = 6 \times 10^{-8}$ , Figure  
5 2A), suggesting that CMV infection is associated with CD8+ T-cell activation in  
6 this infant cohort. A network representation of positively correlating cell  
7 populations <sup>4</sup> (spearman rho p-value <0.05) revealed 3 clusters dominated by  
8 activated CD8+ and activated CD4+ T-cells with CMV-specific IFN- $\gamma$  ELISpot  
9 response, CD3+ T-cells with Epstein–Barr virus (EBV)-specific IFN- $\gamma$  ELISpot  
10 response and monocytes with B-cells (Figure 2B).  
11



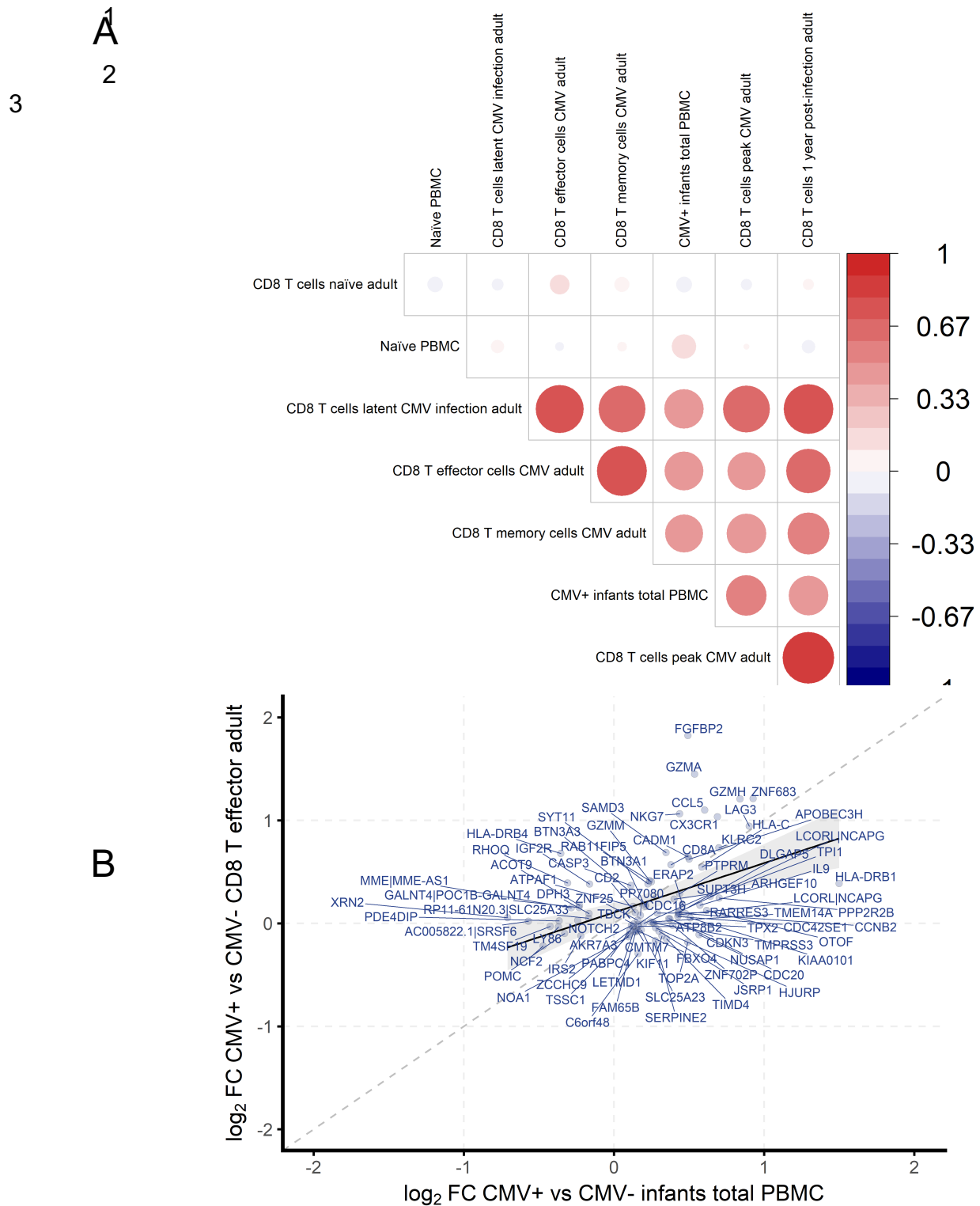
**Figure 2. CMV is associated with CD8 T cell activation in South African infants.** A) Correlation matrix of significantly (spearman rho  $p < 0.05$ ) correlated cell populations and IFN- $\gamma$  ELISpot responses to EBV and CMV. The magnitude of the CMV specific IFN- $\gamma$  ELISpot response correlated with the frequency of activated CD8+ T cells. B) Network of positively correlating cell populations (spearman rho  $p < 0.05$ ) showing 3 clusters dominated by activated T cells with CMV, CD3+ T cells with EBV and monocytes with B cells (node colour indicates cluster membership<sup>4</sup>). Red lines indicate between cluster correlations and black lines within cluster correlations. Line width indicates the correlation coefficient. C) Volcano plot showing magnitude and significance of differential expression between EBV+ and EBV- infants and D) CMV strongly positive (ELISpot  $> 100$  SFC/million) and CMV negative infants. The top 50 significant genes are labelled and horizontal and vertical dashed lines indicate 20% FDR and 5% change in gene expression respectively.

1 EBV had a strong effect on the blood transcriptome, with 296 genes  
2 differentially expressed between infants with positive and negative EBV  
3 responses (Figure 2C). CCL8, CXCL10 and IFIT3 had the greatest fold  
4 increase in expression in EBV+ infants indicating strong induction of a Type I/II  
5 IFN response, although we did not see a correlation between EBV ELISpot  
6 response and T cell activation in this study (Supp Table 1).

7 The impact of CMV on the blood transcriptome was smaller, with only 14 genes  
8 differentially expressed between CMV positive and negative infants (ELISpot  
9 >17 SFC/million), although there were 103 differentially expressed transcripts  
10 between CMV strongly positive (ELISpot >100 SFC/million) and negative  
11 infants (Figure 2D and Supp Tables 2A and B). Differentially expressed  
12 transcripts included HLA-DRB1, ZNF683, LAG3 and KLRC2 (NKG2C). ZNF683  
13 is an important paralog of PRDM1 and together with CCL5 has been shown to  
14 be highly induced in human CMV-specific CD8+ T-cells when compared to  
15 naïve CD8+ T-cells<sup>22</sup>. LAG3 is expressed on activated CD8+ and CD4+ T-cells  
16 and LAG3+CD8+ T-cells are found in CMV infection<sup>22,23</sup>. NKG2C+ natural killer  
17 (NK) cells are expanded in response to CMV infection although CMV infection  
18 can also induce the expression of NKG2C on CD8+ T-cells<sup>24-26</sup>.

19 Expression levels of transcripts from total peripheral blood mononuclear cells  
20 (PBMC) of infants with a CMV ELISpot response (>100 SFC/million) were  
21 correlated with expression levels of transcripts published from CD8+ T-cells  
22 isolated from adults at different stages of CMV infection<sup>22</sup>. The highly  
23 significant correlation of infant transcripts with CMV virus specific CD8+ T-cell  
24 transcripts from adults confirms the prominence of an activated CD8+ T-cell

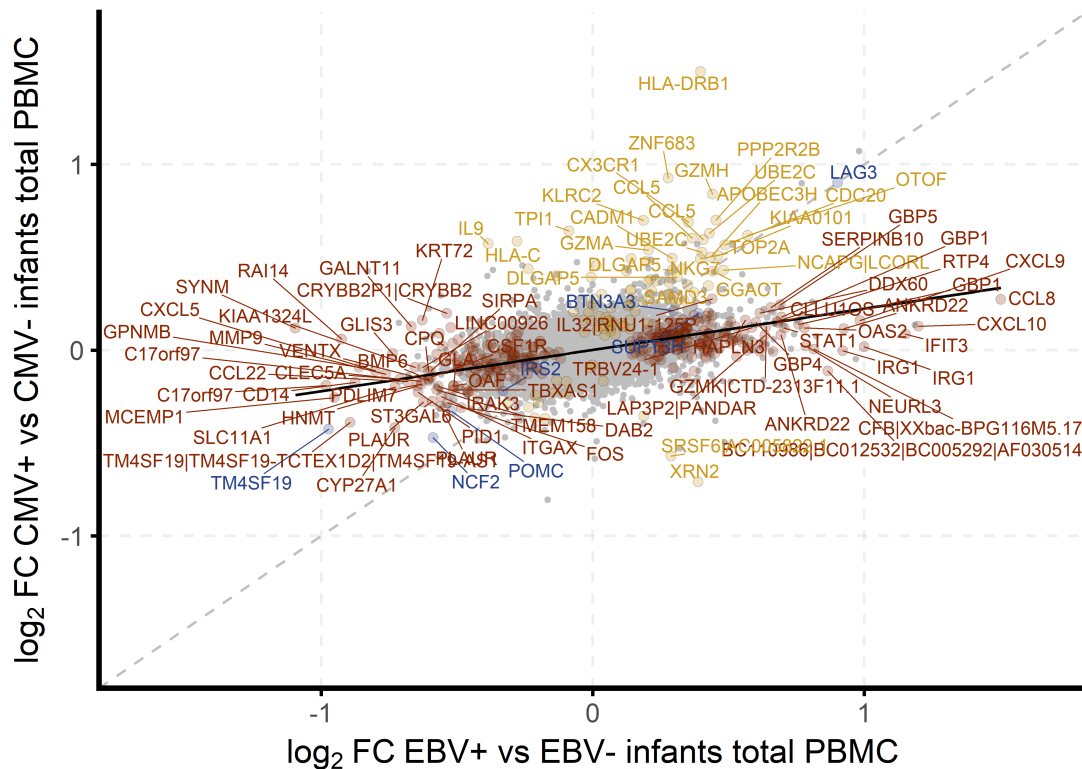
- 1 response in unstimulated PBMC of CMV ELISpot positive infants (Additional
- 2 Figure 1 A and B).
- 3



Additional Figure 1. **Prominence of CD8 T cell specific transcripts from total PBMC of CMV ELISpot positive infants.** A) Spearman's rho from correlations of gene expression from total PBMC of infants with a strong CMV ELISpot response >100 SFC/million with CD8+ T cells isolated from adults who were naïve or infected with CMV (transcripts with <20% FDR). B) Scatter plot of 84 overlapping transcripts between PBMC of infants with a strong CMV ELISpot response >100 SFC/million and transcripts from CMV virus specific effector CD8 T cells from CMV infected adults (GSE24151).

1 Interestingly, there was almost no overlap in differentially expressed genes  
2 between CMV and EBV positive infants with only seven genes significantly  
3 differentially expressed in both comparisons at an FDR of 20% (Additional  
4 Figure 2).

5



6

7

Additional Figure 2. **Little overlap in expression in infants who are either EBV or CMV ELISpot positive.** Scatterplot of fold changes in response to EBV infection (x-axis) plotted against changes in response to CMV infection (y axis). The almost horizontal linear regression line indicates stronger changes in response to EBV. At an FDR of 20%, seven genes are significantly differentially expressed in both comparisons.

8

1 Our cellular data suggests that CMV IFN- $\gamma$  responses are associated with  
2 activated HLA-DR<sup>+</sup> CD8<sup>+</sup> T-cells and transcriptional analysis supports the  
3 presence of an activated CD8<sup>+</sup> T-cell phenotype amongst total, unstimulated  
4 PBMC from CMV IFN- $\gamma$  ELISpot positive infants.

5

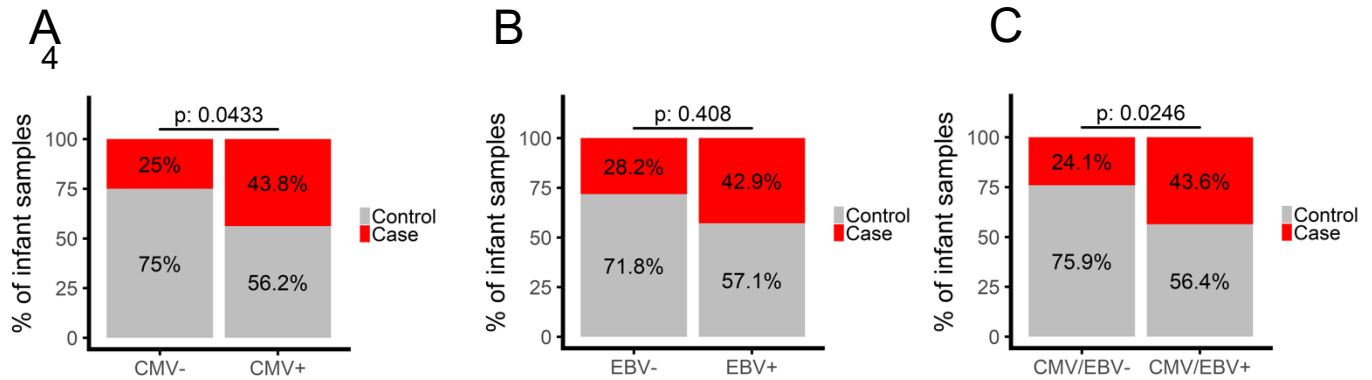
6 **CMV antigen-specific T-cell responses, measured up to 3 years prior to**  
7 **detection of TB, associate with risk of TB disease in infants**

8 CMV is associated with increased risk of HIV infection and disease progression  
9 <sup>11,12,14</sup>. To determine if viral infection is associated with TB risk in our infant  
10 cohort, we analyzed CMV-specific and EBV-specific IFN- $\gamma$  ELISpot responses  
11 in cases and controls. Infants dual positive for both EBV and CMV (n=3) were  
12 excluded from the analysis. A CMV specific IFN- $\gamma$  ELISpot response (response  
13 >17 SFC/million at either time point), measured up to 3 years prior to TB  
14 detection, was associated with increased risk of TB disease (CLR,  $p=0.043$ , OR  
15 2.22 95% CI 1.02-4.83) (Figure 3A). EBV alone was not associated with risk  
16 (Figure 3B), although a combined CMV and EBV response was associated with  
17 increased risk of TB disease (CLR,  $p=0.025$ , OR 2.3 95% CI 1.11-4.79) (Figure  
18 3C). To further explore this association, we analyzed time to TB diagnosis in  
19 TB cases. Infants with a positive CMV or CMV/EBV<sup>+</sup> ELISpot response  
20 developed TB earlier than negative infants (Log Rank Mantel-Cox  $p=0.037$ ,  
21 Figure 3D and 3E).

22 We analyzed data by vaccine group and saw no evidence that CMV positive or  
23 CMV/EBV positive infants immunized with MVA85A were at greater or lower  
24 risk of TB disease when compared to CMV negative infants immunized with  
25 MVA85A ( $p = 0.8355$  for CMV and  $p = 0.9177$  for CMV/EBV).

1  
2

3  
4



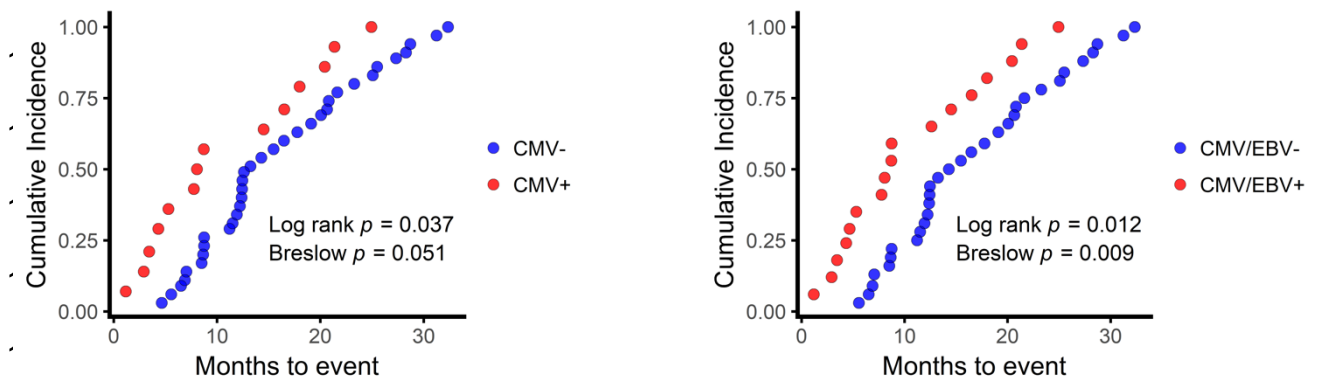
9

10

11

12

13



19

20

### Figure 3. **CMV+ infants are at increased risk of developing TB disease**

A) We saw a higher proportion of case (red) infants among CMV+ ( $n = 18/32$ ) when compared to CMV- infants ( $n = 35/140$ ). B) There was no significant enrichment for cases among EBV positive infants ( $n = 3/7$  compared to  $n = 46/162$ ) although Infants positive for either CMV or EBV (C) were at increased risk D) CMV+ infants (red) develop TB disease earlier in follow-up when compared to CMV- infants (blue) and E) Infants positive for either CMV or EBV develop TB disease earlier than CMV/EBV- infants.



1 A CMV specific IFN- $\gamma$  response measured at 4-6 months of age, up to 3 years  
2 before disease is detected, was a risk factor for the development of TB disease  
3 in South African infants and this risk was greatest during the first 10 months of  
4 follow-up.

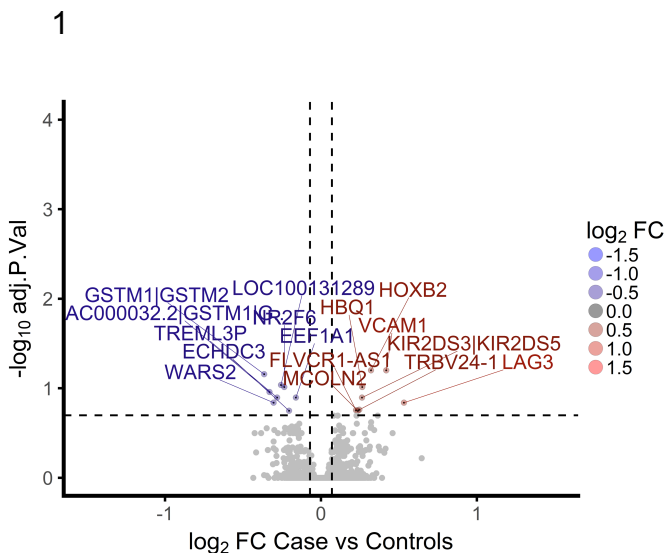
5

6 **Transcriptional evidence of activated T-cells, Type I IFN responses and**  
7 **NK cells in infants up to 3 years prior to detection of TB**

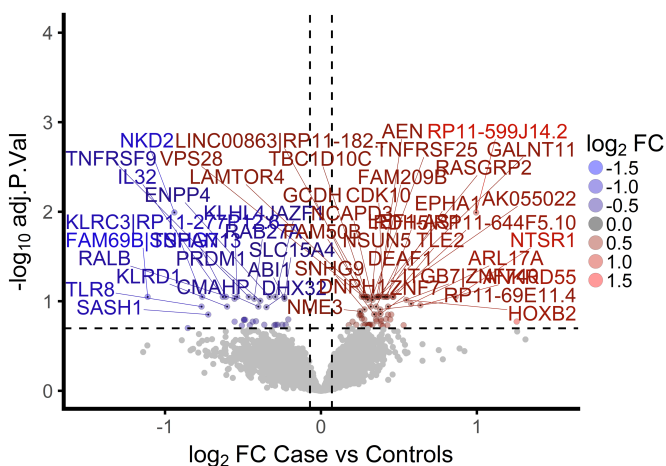
8 Because CMV ELISpot positive infants were at greater risk of disease, we  
9 stratified transcriptome data by CMV status to identify transcripts able to  
10 classify case and control infants. When CMV+ and CMV- infants were analyzed  
11 together, 16 genes were significantly differentially expressed between cases  
12 and controls (Figure 4A and Supp Table 3). LAG3 and VCAM1, known markers  
13 of T-cell activation<sup>23,27,28</sup>, had the greatest fold increase in expression in case  
14 infants. To test our ability to classify infants into cases or controls we split  
15 samples randomly into a 70% training set and 30% test set and trained an  
16 artificial neural network (Additional Figure 3A). The process was repeated fifty  
17 times with random splits into training and test set (bootstrapping) and AUROC  
18 and accuracy were recorded to assess predictive stability (Additional Figure 4).  
19 When all infants were included in the analysis we could classify infants into  
20 cases and controls with an average accuracy of 67% (AUROC 0.77, 95% CI  
21 0.69-0.85, Additional Figure 6 A). Prediction accuracies were comparable when  
22 alternative, widely used classification algorithms were used (Additional Figure  
23 3B).

24

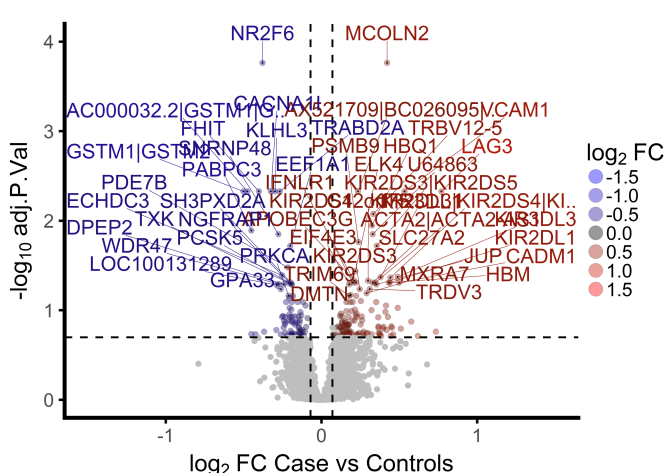
**A**



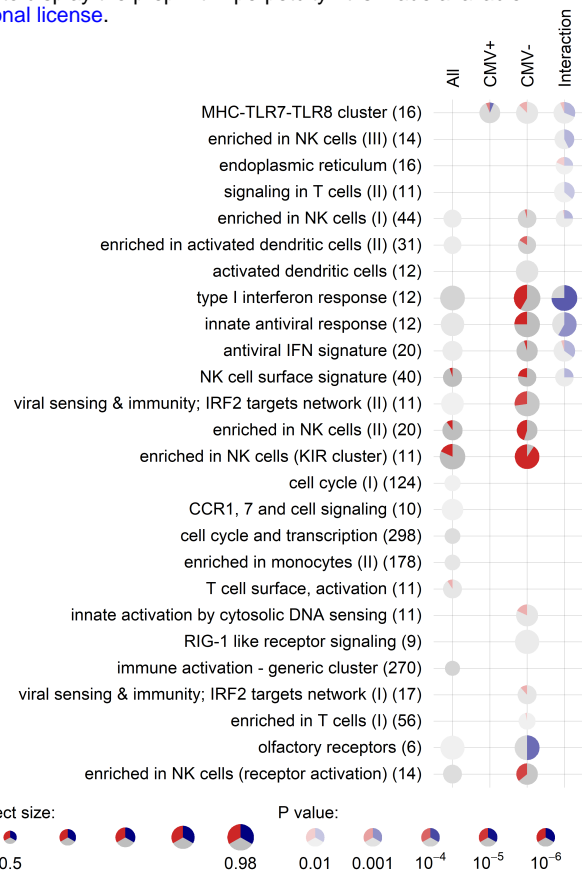
**B**



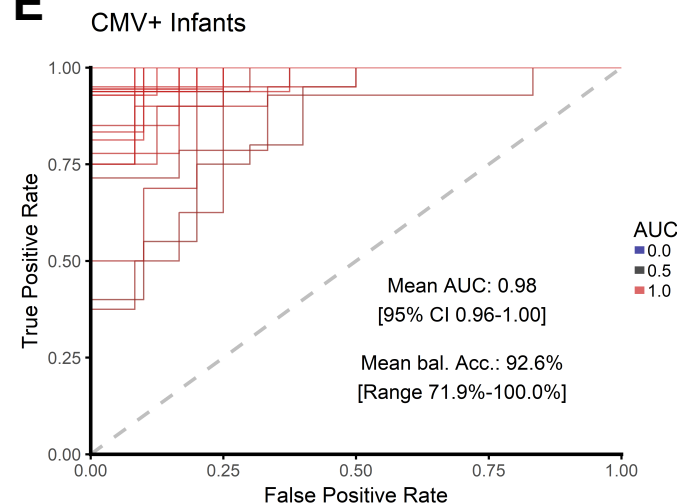
**C**



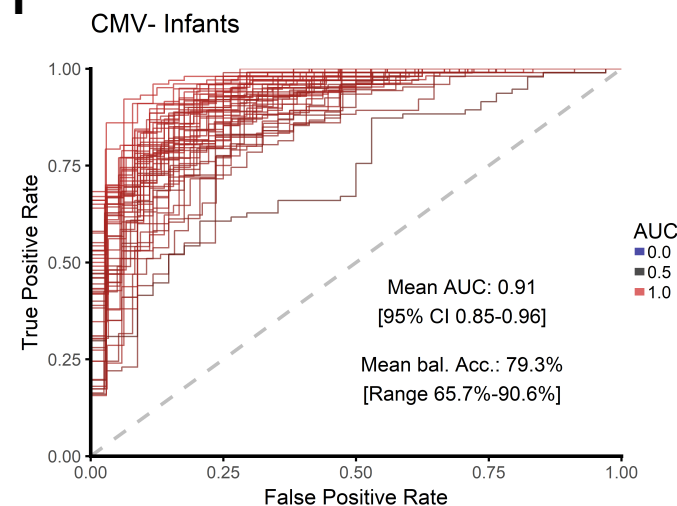
**D**



**E**



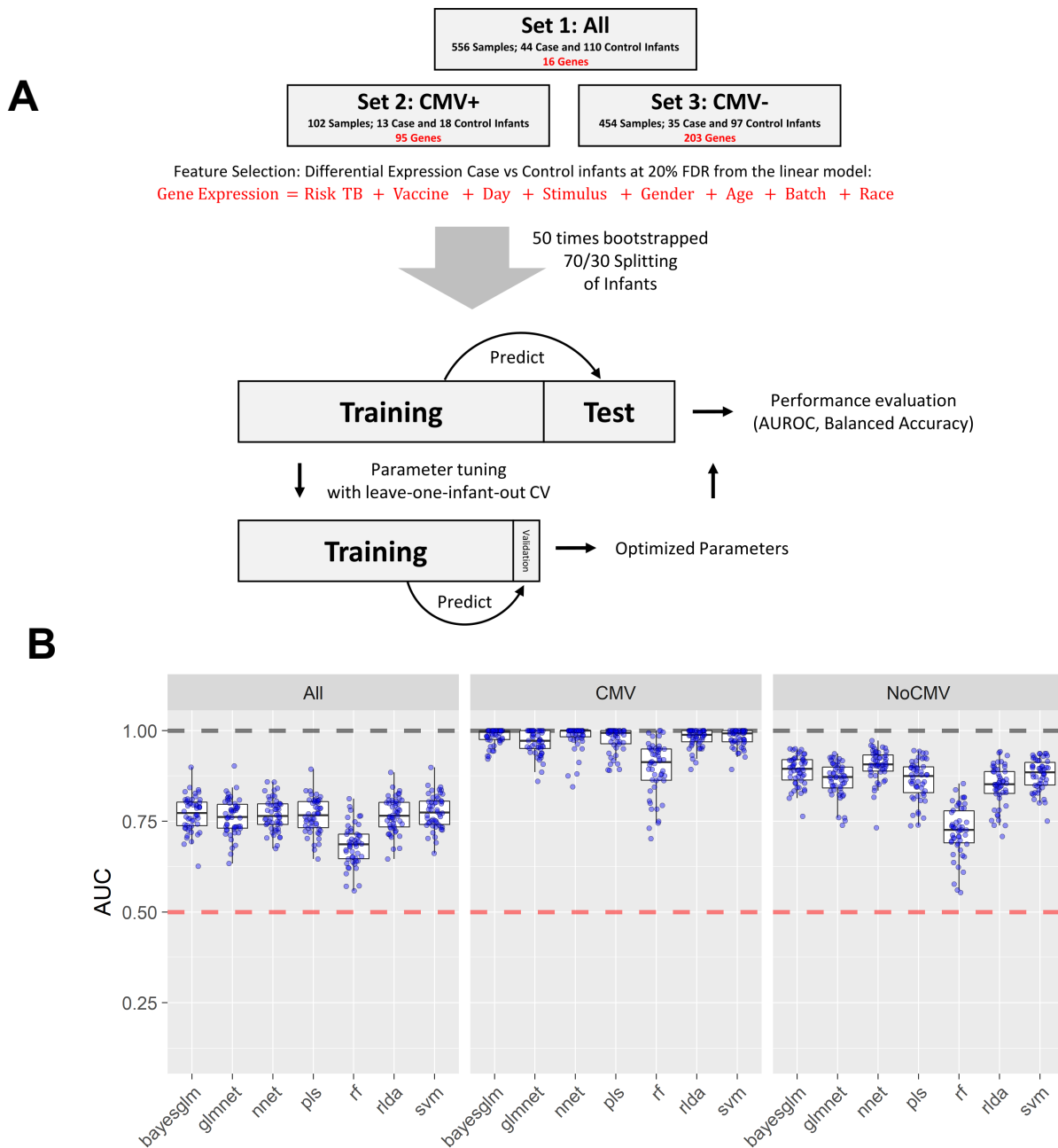
**F**



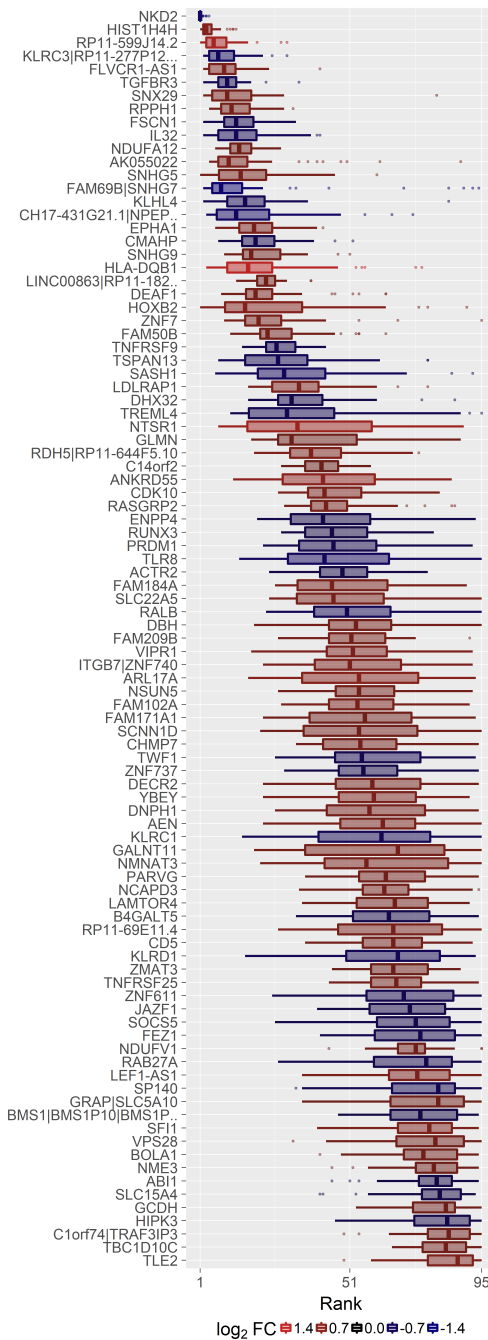
1

**Figure 4. Transcriptomic correlates of risk of TB disease are different in CMV+ and CMV- infants**

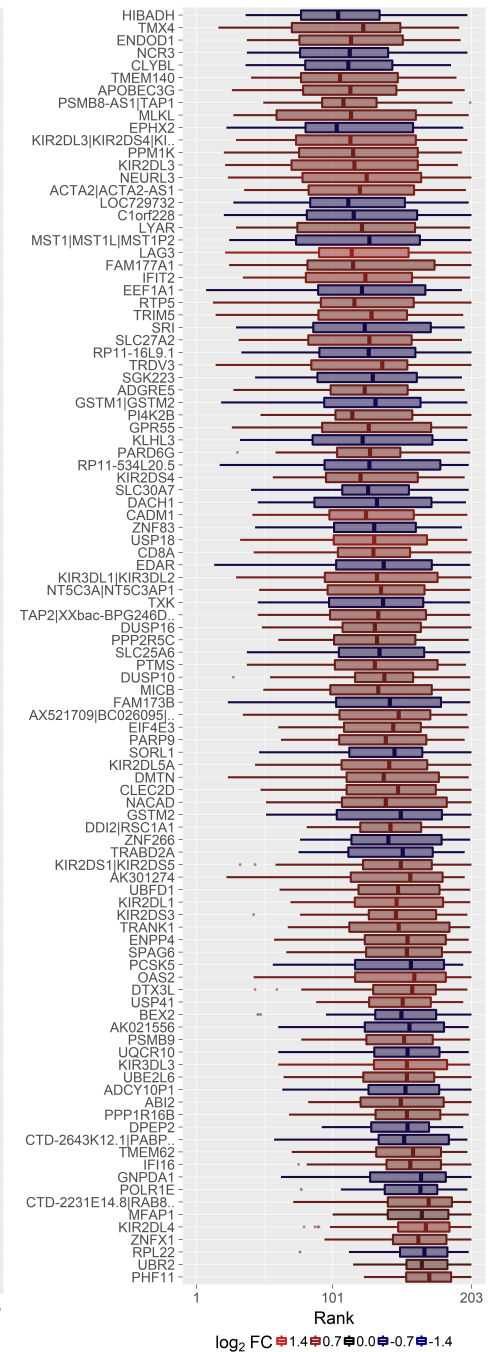
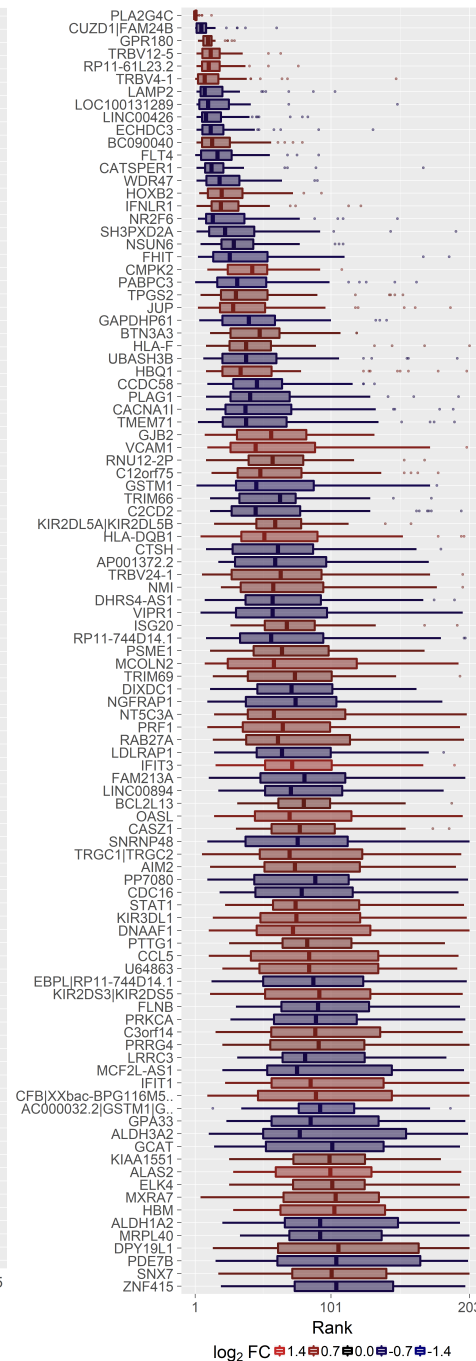
Volcano plots showing magnitude and significance of differential expression between all case and control infants A), CMV+ case and control infants B) and C) CMV- case and control infants. The top 50 significant genes are labelled and horizontal and vertical dashed lines indicate 20% FDR and 5% change in gene expression respectively. D) Enriched modules for differential expression in case and control infants amongst all, CMV+ and CMV- infants. Each row contains one module with the number of genes indicated. Each significantly enriched module at a p-value < 0.05 is shown as a pie chart. The size of the pie corresponds to the AUROC in the cerno test, and intensity of the colour corresponds to the enrichment q-value. The red and blue colour indicates the amount of significant up and down regulated genes respectively, grey color indicates the remaining not significant genes within the category. The interaction term evaluates the statistical difference between changes in CMV+ and CMV- infants. E) AUROC of the classification performance of the artificial neural network model which was trained using approximately 70% of the data and risk of TB was predicted on the withheld 30% of the data (Additional Figure 3 A). The process was repeated fifty times with random splits into training and test set (bootstrapping) and the AUROC was recorded for each round for E) CMV+ and F) CMV- infants respectively.



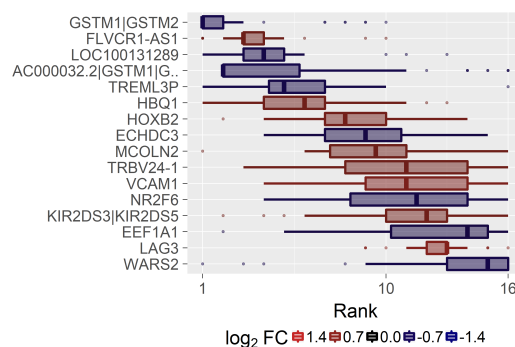
## CMV+



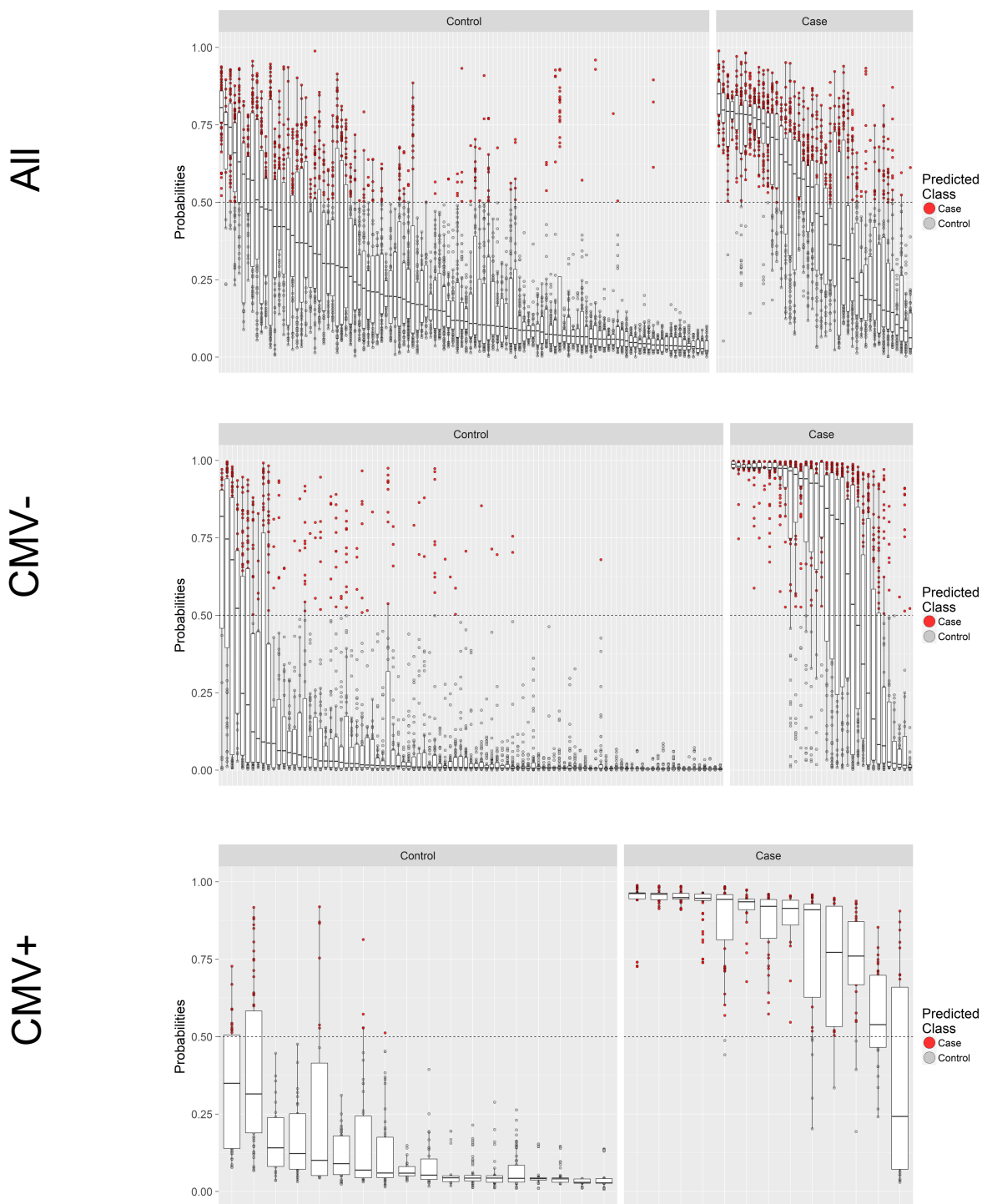
## CMV-



## All



Additional Figure 4. **Boxplots showing the relative importance of features used to train the classifier for each data set.** For each repeat, features were ranked by variable importance and the relative rank was recorded. The lower the average number, the more often a feature has been assigned a high importance by the classifier during the fifty repeated predictions. The color indicates the Case vs Control infants  $\log_2$  fold change estimate based on the differential expression analysis.<sup>21</sup>



Additional Figure 5. **Average prediction accuracies for risk of TB are very high in infants stratified by their CMV status.** Box plots of probabilities assigned to each sample by the trained neural network. Each box plot represents one infant and contains stimulated as well as unstimulated samples and the results of fifty bootstrapped repeats of the prediction. Left and right panel divide the infants into cases and controls and point color indicates the classification decision by the model.

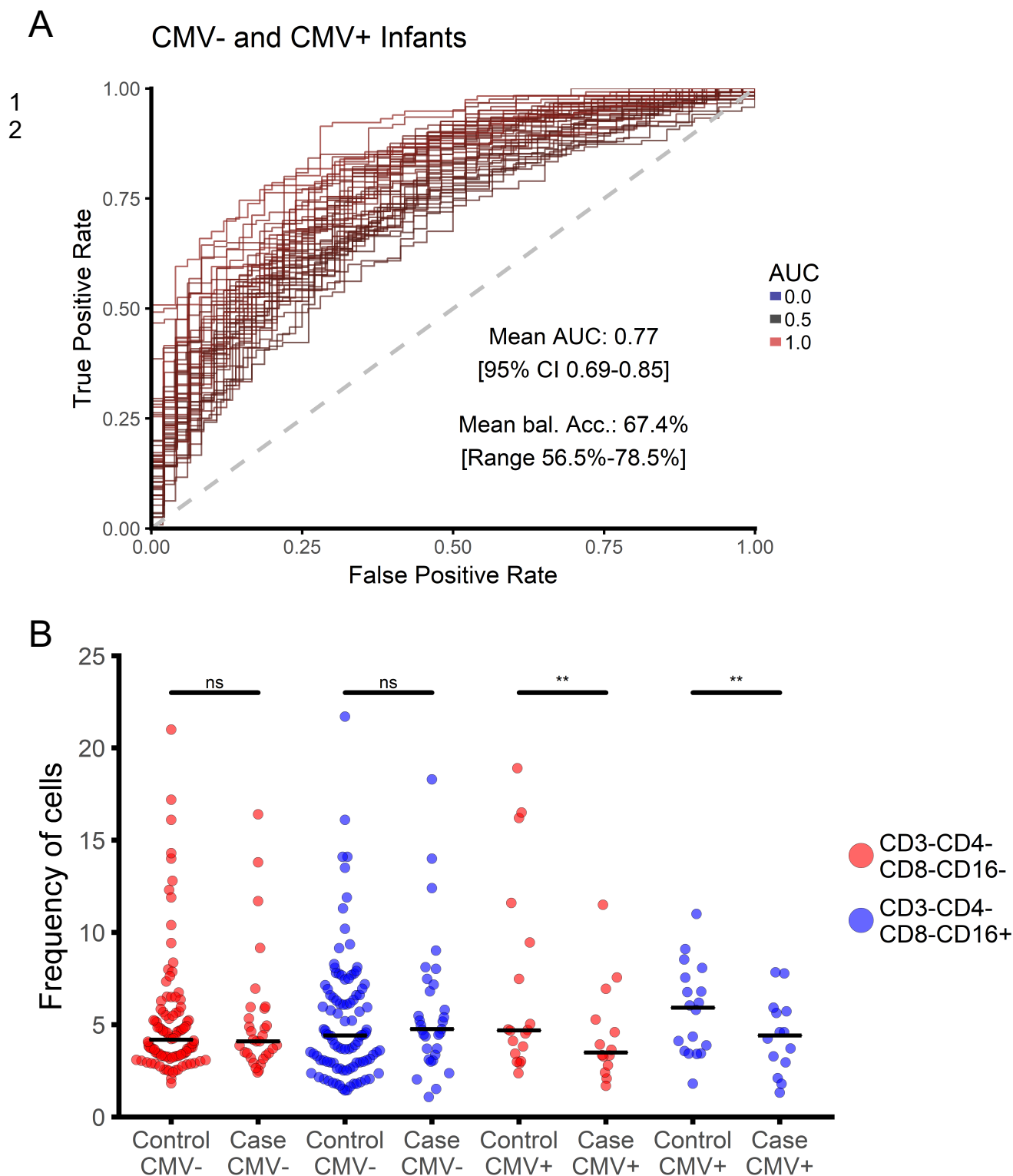
1 In CMV+ healthy infants who developed TB in the following 3 years (cases), the  
2 NK cell associated cytokine IL32 and the NK cell-specific lectin-like receptors  
3 KLRC1 and KLRC3 were among transcripts with the greatest decrease in fold-  
4 change of expression and the highest predictive power (Figure 4B, Additional  
5 Figure 4 and Supp Table 4). IL32 enhances maturation of monocytes to  
6 macrophages and has been shown to be important for protection against  
7 *Mycobacterium tuberculosis* (*M.tb*)<sup>29</sup>. In our cellular analysis we observed  
8 decreased frequencies of CD3-CD8-CD4- (triple negative) CD16- and CD16+  
9 natural killer cells in infants who develop TB disease in the next 3 years when  
10 compared to CMV+ infants who do not develop disease (Additional Figure 6B).  
11 In CMV- case infants, elevated expression of T-cell activation markers,  
12 including LAG3 and VCAM1, markers of a type I/II IFN response including  
13 IFIT3, and enhanced expression of a broad range of both activated and  
14 inhibitory KIR receptors including KIR2DL1, KIR2DL3, KIR2DL4, KIR2DL5A,  
15 KIR2DS3, KIR2DS5, KIR3DL1 and KIR3DL3 were observed (Figure 4C and  
16 Supp Table 5). Modular pathway analysis showed enrichment for NK and KIR  
17 cluster genes and a type I/II IFN antiviral immune response in CMV- case  
18 infants when compared to controls (Figure 4D). However, we saw no evidence  
19 of increased NK cell frequency in infants who develop TB disease in our cellular  
20 analysis (Additional Figure 6B).  
21 TB risk associated transcripts and immune pathways were different among  
22 CMV+ and CMV- infants, as highlighted by a modular pathway analysis which  
23 uses an interaction term to compare pathways associated with TB risk among  
24 CMV+ and CMV- infants (Figure 4D and Supp Table 6).

1 When infants were stratified by CMV status we were able to classify cases and  
2 controls within the CMV+ cohort with 92% average balanced accuracy and  
3 within the CMV- cohort with 89% average balanced accuracy with an average  
4 AUROC of 0.98 (95% CI 0.96-1) and 0.9 (95% CI 0.85-0.96), respectively  
5 (Figure 4 E and F and Additional Figures 4).

6

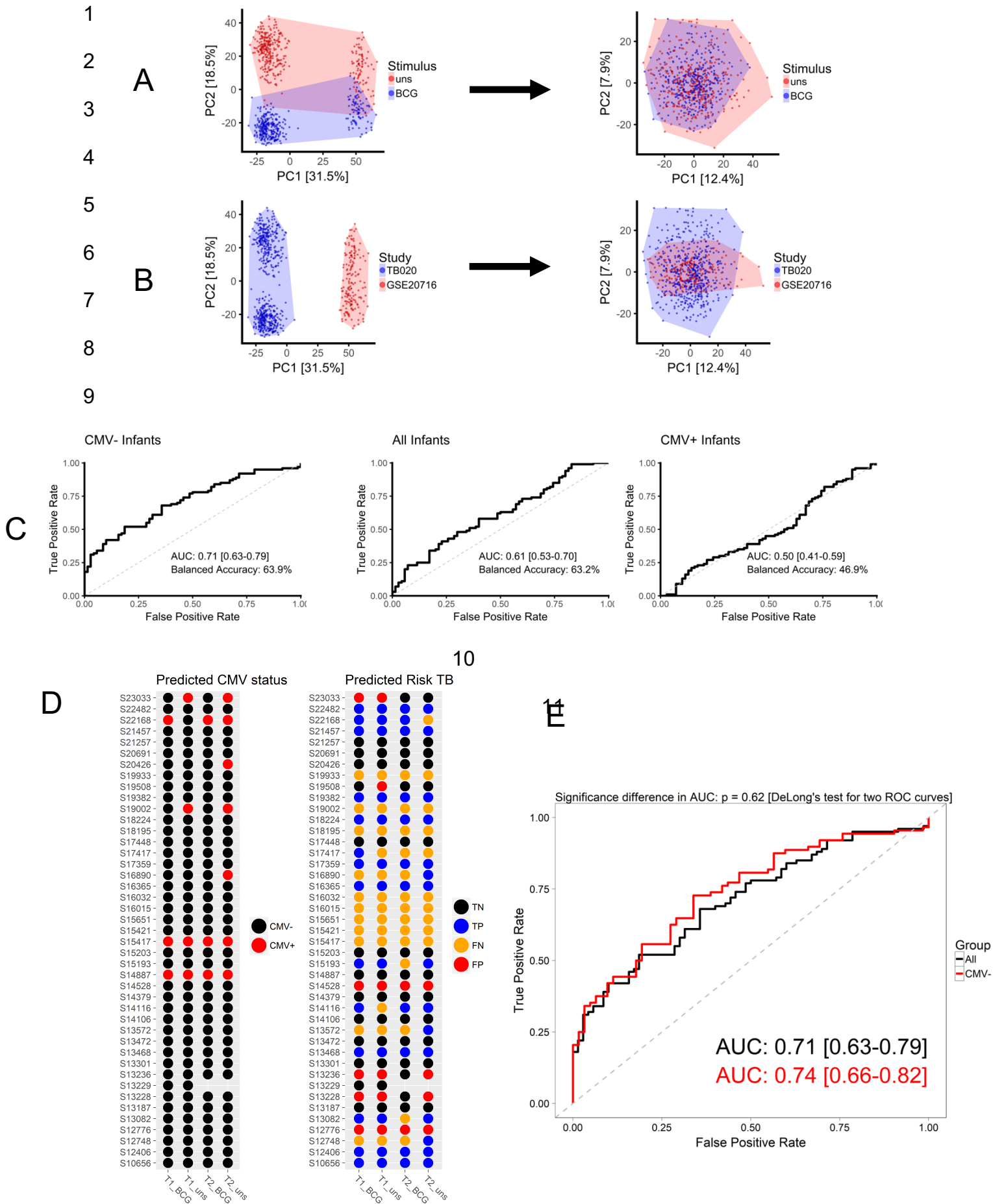
7





**Additional Figure 6. Classification into cases and controls using all samples.** A) Classification performance of the artificial neural network using CMV+ and CMV- infants. B) CD16- and CD16+ (putative) NK cell frequencies among case and control infants with and without a CMV response

1 To validate the classifier signatures, we used raw data from an independent  
2 cohort of 10-week-old South African infants vaccinated with BCG at birth and  
3 unknown CMV status (GEO gene set GSE20716, Fletcher *et. al.* <sup>3</sup>) Infants were  
4 enrolled into an efficacy trial of intradermal or percutaneous delivery of  
5 Japanese BCG at birth <sup>7</sup>. A nested correlates of risk study was performed using  
6 blood from 10-week-old healthy infants who developed TB disease within the  
7 next two years <sup>3,30</sup>. We sought to remove technical differences between the two  
8 array data sets to enable validation of classifier signatures (Additional Figure  
9 7A), however, we could not fully normalize expression between the two data  
10 sets, most likely due to the age difference between infants in the two cohorts  
11 (2-3 months in GSE20716 and 4-6 months in the MVA85A efficacy trial)  
12 (Additional Figure 7B). Despite these cohort differences and the loss of more  
13 than half of the genes in our signature due to microarray platform differences,  
14 we were able to predict risk of TB using our CMV- classifier with an accuracy  
15 of 63.9% and an AUROC of 0.71 (95% CI 0.63-0.79, Additional Figure 7C).  
16 Next, we attempted to improve accuracy by prediction and removal of infants  
17 with suspected CMV infection. The number of CMV positive infants with more  
18 than one predicted positive sample was low (n = 5, Additional Figure 7C) and  
19 excluding these infants lead to a slight but not significant improvement in  
20 balanced accuracy of 65.6% and an AUROC of 0.74 (Additional Figure 7E).



**Additional Figure 7. Validation of TB predictive biomarkers from CMV- infants in an independent cohort of BCG vaccinated infants**

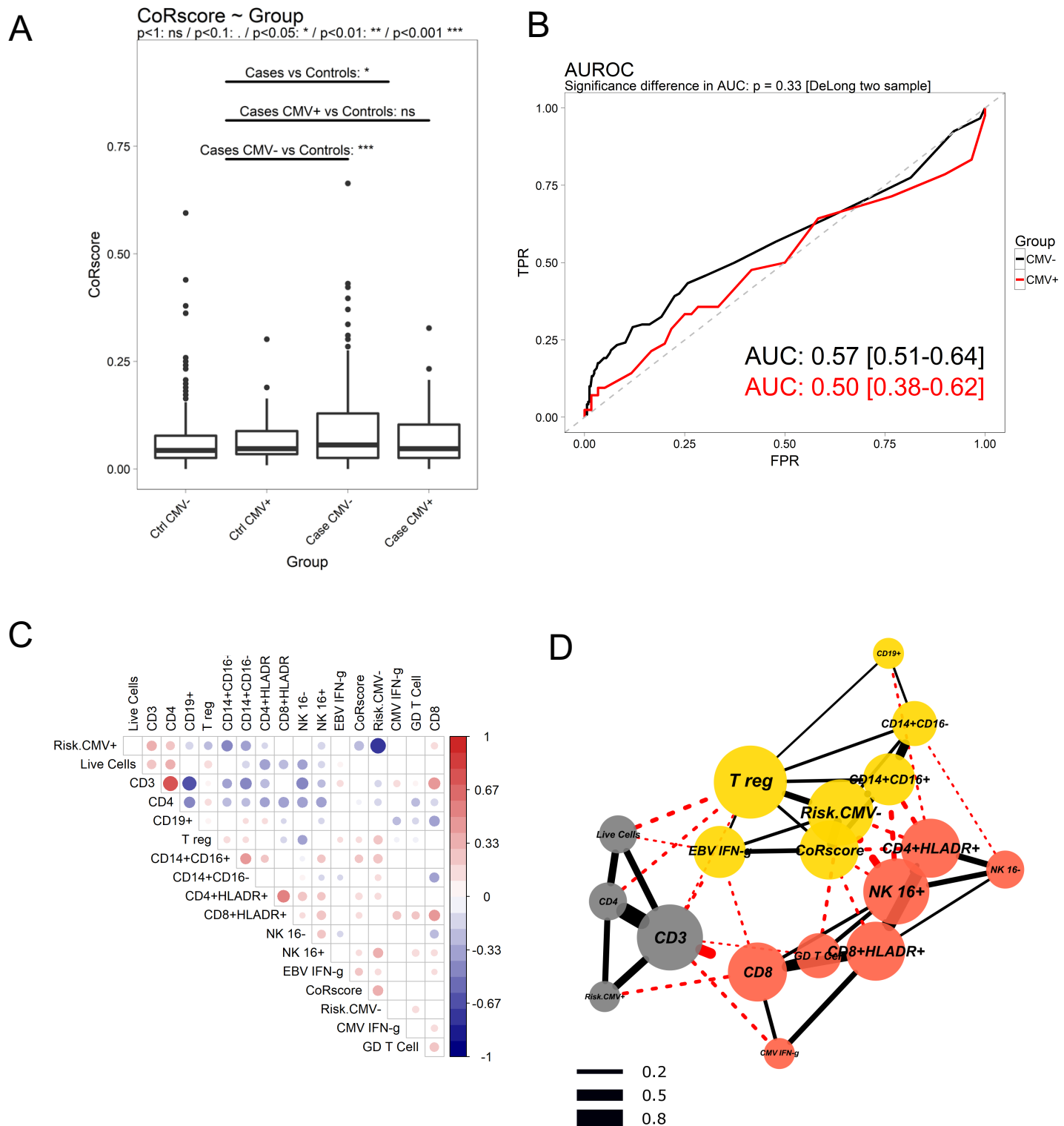
Raw data and detection p values from an independent cohort of 10 week old South African infants vaccinated with BCG at birth (GEO gene set GSE20716<sup>3</sup>) were used as a validation cohort for the classifier signatures from CMV+ and CMV- infants. 20589 Probes from 182 samples were mapped from Illumina HumanRef-seq 8 arrays to Gencode v.23 and 8 outlier samples were excluded A) The effects of time point (4 or 12 hours) and BCG Stimulation (stimulated or unstimulated) were removed and 9195 probes identified as overlapping between the two studies. B) We were only partially able to remove differences in gene expression due to study of origin. The differences observed in PC2 are driven by age, which was non-overlapping in the two studies, 2-3 months in GSE20716 and 4-6 months in the MVA85A efficacy trial.C) Using 55 genes differentially expressed between infants with very high or low CMV titres at an FDR of 5%, 5 infants were identified to have two or more CMV positive samples. D) The charts represents one dot per sample and red colour indicates samples classified as CMV+. The TB Risk dot plot shows samples with correctly predicted Case (TP) and Control (TN) status as well as samples with wrong Case (FP) and Control (FN) labels. E). Excluding all five suspected CMV+ infants we were able to improve risk of TB prediction accuracy on GSE20716<sup>3</sup>slightly, with balanced accuracy of 65.6%.

1

1 Finally, we looked for enrichment of transcripts from the 16- peripheral blood  
2 gene correlate of risk (CoR) signature associated with progression to TB  
3 disease in *M.tb* infected adolescents<sup>2</sup> among our CMV+ and CMV- infants. This  
4 adolescent CoR signature was significantly enriched among case infants when  
5 compared to controls and the enrichment was strongest amongst CMV- infants  
6 (Additional Figure 8A). However, the adolescent CoR signature was not able  
7 to accurately classify infants into cases and controls (Additional Figure 8B).  
8 Including the CoR score in our infant network analysis we show that the CoR  
9 score correlates with the frequency of inflammatory monocytes and with our  
10 CMV- classifier signature (Additional Figure 8C). Our CMV+ classifier signature  
11 is positively correlated with T cell frequency and negatively correlated with  
12 monocyte and CD16+ NK cell frequency.

13 A network representation of positively correlating cell populations (Spearman's  
14 rho p-value <0.05) revealed 3 major clusters with the adolescent CoR and  
15 infant CMV- classifier signature clustering together with inflammatory and  
16 classical monocytes (Additional Figure 8D).

17

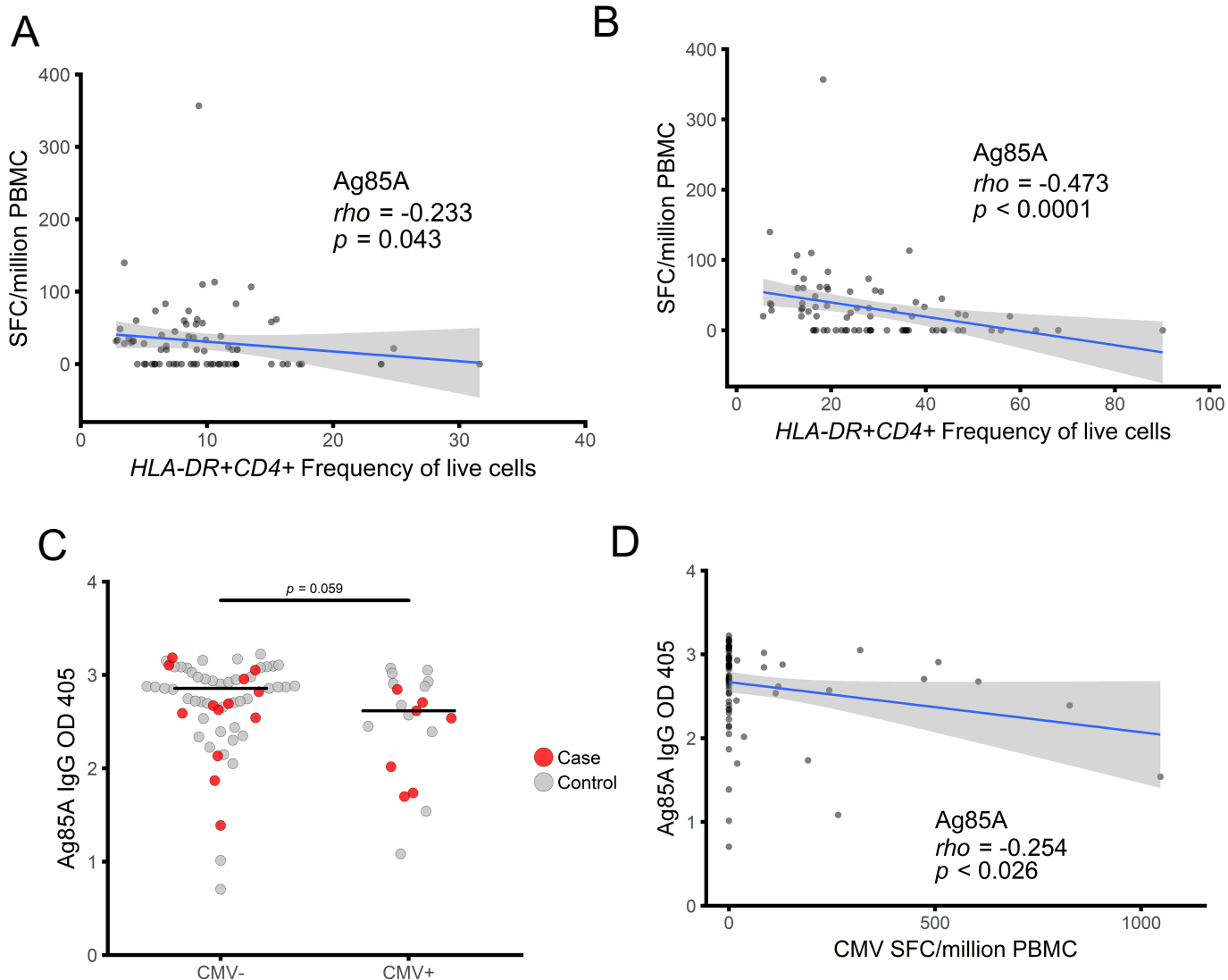


**Additional Figure 8. Enrichment of TB predictive Correlate of TB Risk score (CoR) from adolescents among CMV- case infants when compared to controls, but no accuracy for prediction of TB in infants.** Using a locked-down Illumina model published by Zak et al<sup>2</sup> for prediction of TB progression, a blinded CoR score was derived for each infant. A) After unblinding, the CoR score was found to be significantly higher in case infants when compared to controls and the greatest difference in CoR score was between CMV- case and control infants B) However, despite enrichment in case infants, the CoR score was unable to accurately classify either CMV+ or CMV- infants as cases or controls. C) The CoR score correlates with the frequency of inflammatory monocytes and the CMV-infant classifier signature. D) Network of positively correlating cell populations (spearman rho p-value  $< 0.05$ ) showing a cluster containing the adolescent CoR score with monocytes, EBV ELISpot and the infant CMV- classifier score. Node colour indicates community membership and red and black edges are drawn between and within communities respectively (see methods). Edge width indicates the correlation coefficient.

These findings further support distinct immunological correlates of risk of TB disease in CMV+ and CMV- infants. Taken together our data show increased T-cell activation, KIR receptor signaling and type I IFN response in CMV- case infants and decreased NK cell associated transcripts in CMV+ case infants when compared to their respective controls. Furthermore, we identified transcriptomic signatures, which can identify infants with risk of TB disease with high accuracy in CMV+ and CMV- infants and were able to validate the CMV- biomarker signature in an independent study with moderate accuracy.

### **T-cell activation associates with lower mycobacterial antigen specific immune response following immunization with MVA85A and BCG**

To assess the impact of T-cell activation and CMV on the ability to mount an antigen specific response following immunization with MVA85A, we examined mycobacterial antigen specific IFN- $\gamma$  responses and anti-Ag85A IgG in MVA85A immunized infants at D28 following immunization with MVA85A. IFN- $\gamma$  responses to Ag85A, measured on D28, were inversely correlated with activated CD4+ and CD8+ T-cell frequencies (Figure 5A and B). There was a trend towards lower mycobacterial antigen specific IFN- $\gamma$  responses and lower anti-Ag85A IgG in CMV+ when compared to CMV- infants ( $p = 0.058$ , Mann-Whitney U test) and Anti-Ag85A IgG was inversely correlated with CMV ELISpot response in MVA85A immunized infants (Spearman's Rho,  $p = 0.026$ , Figure 5 C and D).



**Figure 5. T cell activation is associated with lower mycobacterial antigen specific immune response following immunization with MVA85A** The D28 IFN- $\gamma$  ELISpot response to Ag85A was inversely correlated with both A) activated CD4 T cell and B) activated CD8 T cell frequency C) There was a trend for lower anti-Ag85A IgG in CMV+ compared to CMV- infants following immunization with MVA85A. D) anti-Ag85A IgG was inversely correlated with CMV ELISpot response. Mann Whitney U test and Spearman Spearman's Rho correlation. Red = cases and grey = controls.



These data show that T-cell activation and CMV infection influence MVA85A boosting of an antigen specific immune response.

## DISCUSSION

We demonstrate, in healthy infants, that CD8+ T-cell activation and prior or sub-clinical infection with CMV (defined by a positive CMV-specific T-cell response) are associated with increased risk of developing TB disease over the next 3 years of life. We also show that CMV+ infants acquire TB disease earlier than CMV- infants. These data complement our previous finding that CD4+ T-cell activation was associated with TB disease risk in infants and adolescents from a South African community with very high TB incidence<sup>8</sup>. Previous studies have reported that infection with CMV enhances the risk of HIV acquisition and disease progression, through expansion of activated CD8+ T-cells, depletion of naïve T-cells and T-cell senescence<sup>11-14</sup>. More recently, CMV has been implicated in the aetiology of TB, supported by epidemiological associations between the two diseases<sup>17-21</sup>. In Gambian infants, CMV infection induced profound CD8+ T-cell differentiation and activation which persisted up to 2 years after infection<sup>31,32</sup>. Consistent with this effect, we show that infants with a positive T-cell response to CMV peptides have a transcriptional signature associated with CMV specific CD8+ T-cells<sup>22</sup>. However, among CMV+ infants, T-cell activation markers were not differentially expressed between case and control infants.

In CMV+ infants, transcripts associated with NK cells had lower expression and NK cell frequency was lower in cases when compared to controls. A role for NK

cells in protection from TB disease has been demonstrated both in humans and animal models and it is possible that an impaired NK cell response is associated with TB disease risk among CMV+ infants<sup>33-39</sup>. CMV infection promotes the expansion of NK cells expressing the CD94/NKG2C activating receptor, and these cells are important for control of viral replication<sup>40</sup>. The CD94 and NKG2C transcripts *KLRC1* and *KLRC3* had among the greatest decreases in fold-change of expression and the highest predictive power in identifying cases among CMV+ infants. The NKG2C receptor is encoded by the *KLRC2* gene, heterozygous and homozygous deletion of which is present in a significant proportion of individuals across different populations<sup>41</sup>. *KLRC2* deletion is associated with reduced numbers of mature NK cells and increased susceptibility to HIV infection, certain autoimmune conditions and cancer<sup>41-43</sup>. We hypothesize that increased susceptibility to TB in CMV+ infants may result from loss of control of CMV replication and/or impairment of NK cell function due to *KLRC2* gene deletions in some individuals<sup>17</sup>.

Among CMV- infants, who went on to develop TB, we observed upregulation of transcripts associated with T-cell activation, including *LAG3*<sup>23</sup> which is induced during active TB in a non-human primate model<sup>44</sup>. We also found multiple transcripts and pathways that are typically altered during viral infection among CMV- infants. This may be due to underestimation of the prevalence of viral infection, as we used only IFN- $\gamma$  responsiveness to EBV and CMV CD8 epitopes as a measure of viral infection. Based on previous studies of viral prevalence in infants in Africa we would expect a CMV prevalence of 24-41%<sup>45,46</sup> and EBV prevalence of 35% by 6 months of age<sup>47,48</sup>. Detection of viral DNA would allow more accurate diagnosis of viral infection, however, this was

not possible due to the very limited samples collected from these infants. Moreover, viruses other than those measured in this study (CMV, EBV) may be contributing to risk of TB disease among these infants. Expression of a broad range of both activating and inhibitory KIR receptor transcripts was elevated in CMV- case infants. Exposure to multiple viral infections drives high diversity of KIR expression and lowers the availability of naïve NK cells to respond to future infectious challenge, resulting in susceptibility to HIV <sup>49,50</sup>. We were able to identify different classifier signatures among CMV+ and CMV- infants and were able to verify our CMV- signature in an independent cohort of infants <sup>3</sup>. In the independent cohort, infants were younger (2-3 months of age) and few infants were classified as CMV positive. CMV infection and viral replication is low at birth, peaks at 3-6 months of age and declines to plateau at 8-10 months of age <sup>51</sup>. At 4-6 months of age, infants recruited in to the MVA85A efficacy trial <sup>1</sup> were at the peak age of CMV viral replication in infancy.

Our transcriptional evidence of antiviral immune responses in infants who develop TB disease is also consistent with the observations of Zak *et al.* <sup>2</sup> who reported increased expression of Type I/II IFN associated transcripts in *M.tb*-infected adolescents who progressed to TB disease. Increased Type I/II IFN transcripts have also been observed in patients with active TB disease, when compared to *M.tb*-infected or uninfected controls <sup>52-56</sup>. Recent data from the *M.tb*-infected adolescent study has shown that an increased Type I/II IFN response precedes a shift towards an elevated monocyte to lymphocyte ratio and an increase in T cell activation, which are detected closer to the time of TB disease <sup>57</sup>. The authors suggest that the initial elevation in Type I/II IFN could be driven by viral infection and that this could then trigger the immune events

that lead to TB susceptibility<sup>57</sup>. Consistent with this we found a correlation of the adolescent CoR score with an EBV response and with elevated CD14+CD16+ inflammatory cells. We also observed significant enrichment of the 16-gene CoR signature identified by Zak *et al.* among our CMV- case infants<sup>2</sup>. However, the CoR signature was not able to accurately classify case and control infants. In our study, the infants were not infected with *M.tb*, had no symptoms of TB disease and no known exposure to TB in the household during enrolment or when the blood samples were obtained. In addition, since incident TB disease was diagnosed months to years after blood sample collection it is unlikely that the elevated type I/II IFN associated transcriptomic signatures we observed are a result of sub-clinical TB disease. Although there was no evidence of BCG vaccine-induced disease in these infants we cannot rule out that persistent, sub-clinical, replication of BCG vaccine could be driving a type I/II IFN response in some infants.

All infants in this study received BCG at birth and the T cell response to BCG peaks at 2-3 months of age<sup>58</sup>. Viral infections during the development of the BCG-specific immune response may impair the development of protective immunity, as has been suggested by studies in Malawi<sup>59</sup>. In the Gambia, exposure to HIV *in utero* and being born in the wet season, associated with increased respiratory and diarrheal disease, have been shown to impact the BCG antigen specific T-cell response following vaccination<sup>31,60</sup>. We observed an inverse correlation between CD4+ and CD8+ T-cell activation and antigen specific T-cell responses following immunization with MVA85A, suggesting that T-cell activation may be associated with decreased vaccine boosting in infants. This confirms previous observations where increased immune activation

associated with lower immune responses to MVA85A in both infants and adults<sup>61,62</sup> It could also explain why immune responses to MVA85A were lower when administered within a week of vaccination with DTwP-Hib and hepatitis B in the Expanded Programme on Immunization (EPI) schedule<sup>63</sup>.

Childhood TB is difficult to diagnose and treat and improved strategies are needed to control TB in children<sup>5</sup>. We found that CMV infection was associated with increased risk of developing TB disease in infants and that distinct immune pathways were associated with TB disease risk in CMV+ and CMV- infants. We suggest that viral infection can increase the risk of progression to TB disease and may compromise the immune response to TB vaccines given in infancy. Strategies which include the use of vaccines or antivirals to reduce chronic viral infection in infancy could enhance TB vaccine efficacy, reduce TB risk and help to reduce the global burden of childhood TB.

## METHODS

### Case-control design

BCG-vaccinated infants who were enrolled in an efficacy trial of the candidate TB vaccine MVA85A were included in this study, ClinicalTrials.gov number NCT00953927<sup>1</sup> (Figure 1). This trial was approved by the University of Cape Town Faculty of Health Sciences Human Research Ethics Committee, Oxford University Tropical Research Ethics Committee, and the Medicines Control Council of South Africa. All infants received BCG within 7 days of birth. HIV and *M.tb* uninfected infants without known TB exposure were randomized at 16-24 weeks of age to receive a single intradermal dose of MVA85A or placebo (Candin™, a candida skin test antigen)<sup>1</sup>.

Transcriptomic analysis was performed using PBMC from infants who were included in our previously described case-control study<sup>8</sup>. Briefly, infants who met the primary case definition for TB disease were included as cases and for each case, three infants were randomly selected from a pool of controls (Figure 1). Infants were included in the control pool if they did not demonstrate *M.tb* infection as defined by a positive QuantiFERON TB Gold In-tube test (Cellestis, Australia); had not received TB treatment and had not received isoniazid preventive therapy during study follow-up. Matching was based on gender, ethnic group, CDC weight-for-age percentile ( $\pm 10$  points), and time on study ( $\pm 9$  months).

Infant PBMC samples collected from two time points Day 0 (D0) and Day 28 (D28) were combined for testing of both cellular parameters and transcriptional signatures associated with risk of TB disease. Only infants for whom a sample

was available at both D0 and D28 were included in the analysis resulting in combined analysis of a maximum of 98 samples from 49 cases and 258 samples from 129 controls (actual numbers in analysis vary per availability of assay data and are listed in Table 1, Figure 1).

### **Cell culture**

PBMC were retrieved from liquid nitrogen storage, thawed and rested for 2 hours in media containing DNase to aid the removal of debris from dead and dying cells. After 2 hours cells were counted and immediately transferred to cell culture plates for stimulation for RNA extraction, ELISpot assays or measurement of cell-surface markers assessed by flow cytometry. Viability of thawed PBMC was assessed using flow cytometry with Live/Dead Violet stain (Invitrogen). Phytohemagglutinin (PHA) was included as a positive control for cell viability on ELISpot plates.

### **RNA processing**

Rested PBMC ( $1 \times 10^6$ ) were resuspended in single wells of 200 $\mu$ l RPMI supplemented with 10% FBS and L-glutamine. Cells were stimulated for 12 hours at 37°C in single wells containing live BCG SSI from pooled vaccine vials ( $\sim 2 \times 10^5$  CFU/ml). After 12 hours of stimulation cells were pelleted and lysed in RNA lysis buffer (RLT, Qiagen). RNA was extracted using the RNeasy Mini Kit (Qiagen) per the manufacturer's instructions with the following modification; in the first step an equal volume of 80% ethanol was added to cells lysed in RLT buffer, mixed and total volume transferred to an RNeasy column. RNA was quantified by Nanodrop and stored at -80°C until use. Extracted RNA was

amplified and labeled with biotin using the Illumina Total Prep kit (Ambion) per manufacturer's instructions. Amplified RNA was assessed by nanodrop and Bioanalyzer for quantity and quality prior to hybridization. Hybridization to Illumina HT-12 arrays was performed per manufacturer's instructions. Arrays were scanned using an Illumina iScan machine and data extracted using Genome Studio software.

### ***Ex vivo* IFN- $\gamma$ ELISpot assay**

The *ex vivo* IFN- $\gamma$  ELISpot assay was performed as previously described using a human IFN- $\gamma$  ELISpot kit (capture mAb -D1K) (Mabtech) <sup>8</sup>. Briefly, duplicate wells containing  $3 \times 10^5$  PBMC were stimulated for 18 hours with antigen, PHA or media alone. Antigens included a single pool of Ag85A peptides (2  $\mu\text{g/ml/peptide}$ ) (Peptide Protein Research); BCG ( $2 \times 10^5$  CFU/ml (Statens Serum Institute)); purified protein derivative (PPD) from *M. tuberculosis* (20  $\mu\text{g/ml}$ ) (Statens Serum Institute); peptide pools containing known CD8<sup>+</sup> T-cell epitopes from EBV (15 peptides), and CMV (5 peptides), (2  $\mu\text{g/ml/peptide}$ , ANASPEC). Results are reported as spot-forming cells (SFC) per million PBMC, calculated by subtracting the mean of the unstimulated wells from the mean of antigen wells and correcting for the number of PBMC. A response was considered positive if the mean number of spots in the antigen well was at least twice the mean of the unstimulated wells and at least 5 spots greater.

### **Cell surface flow cytometry**

As previously described <sup>8</sup> PBMC were washed and stained with 5 $\mu\text{l}$  Live/Dead Violet (Invitrogen) followed by surface staining with the following titrated



antibodies: 0.5 $\mu$ l CD3-AF700 (clone UCHT1, Ebioscience), 2 $\mu$ l CD4-APC (clone RPA-T4, Biolegend), 2 $\mu$ l CD8-Efluor605 (clone RPA-T8, Ebioscience), 2 $\mu$ l CD14-PE/Cy7 (clone HCD14, Biolegend), 2 $\mu$ l CD16-AF488 (clone 3G8, Biolegend), 1 $\mu$ l CD19-PE/Cy5 (clone HIB19, Biolegend), 2 $\mu$ l CD25-APC/Cy7 (clone BC96, Biolegend), 2 $\mu$ l CD127-NC650 (clone eBioRDR5, Ebioscience) and 15 $\mu$ l HLA-DR-PE (clone L243 Biolegend). Fluorescence minus one (FMO) controls were used to set gates for CD25, CD127 and HLA-DR. Samples were acquired on a BD LSR II flow cytometer. Results are presented as percentages of cells after excluding dead cells and doublets. CD4<sup>+</sup> and CD8<sup>+</sup> T-cells were identified as CD3<sup>+</sup> cells, while CD14<sup>+/-</sup> and CD16<sup>+/-</sup> cells were identified as CD3<sup>-</sup> and CD19<sup>-</sup> populations. CD25<sup>+</sup> CD27<sup>-</sup> populations were gated on the CD4<sup>+</sup> cells. The network representation of cell populations positively correlated among all infants was done using the igraph package in R. To identify closely related clusters (communities) within the network, the 'cluster\_optimal' function was used implementing an algorithm described in Brandes *et. al.*<sup>4</sup>

### **Transcriptional analysis**

Raw, probe level summary values exported from Illumina GenomeStudio 2011 of Illumina HumanHT 12 V4 microarrays were imported into R using beadarray<sup>64</sup>. Probes were background corrected using negative control probes followed by quantile normalization using the neqc command<sup>65</sup>. The analysis was restricted to probes with a detection p-value of <0.01 in at least 10% of the samples and probes matching to the transcript definition of the following databases (in descending importance) with at most 2 mismatches, no insertions

and a minimum mapping length of 40 bases: GENCODE version 23, RefSeq (refMrna.fa) and GenBank (mrna.fa) downloaded in August 2015 from <http://hgdownload.cse.ucsc.edu/goldenPath/hg38/bigZips/>.

A linear model was fitted using limma<sup>66</sup> to determine differential expression adjusted for vaccine, day, stimulus, gender, age, ethnicity and batch effects. Array quality weights were incorporated<sup>67</sup> to account for between array quality differences. To account for between patient correlations, the duplicateCorrelation command from the limma package was used. Nominal p-values were corrected for multiple hypotheses testing using the Benjamini-Hochberg procedure<sup>68</sup>. Due to the heterogeneity of the samples, a lenient cut off at an FDR of 20% was chosen to identify genes as significantly differentially expressed. In total 221, 101 and 16 probes mapping to 203, 95 and 16 genes were significantly differentially expressed between Case and Control infants within CMV-, CMV+ and the combined group respectively. For the comparison of EBV+ versus EBV-, CMV+ versus CMV- and CMV strongly positive (ELISpot >100 SFC/million) versus CMV- in total 334, 14 and 103 probes mapping to 296, 14 and 97 genes were significantly differentially expressed respectively. Gene set enrichment analysis was carried out using the cerno test from the tmod package in R. Modules for the enrichment analysis were taken from Li et al.<sup>69</sup>. Datasets used for comparative analysis to CMV infected CD8 T-cells were obtained from Gene Expression Omnibus by downloading GSE12589 and GSE24151.

## Classification

For classification into cases and control samples, we used normalized microarray intensities adjusted for scan date, sample collection time point (Day 0 or Day 28) and stimulation (unstimulated or BCG stimulated). Features for training were selected by using the significant probes at 20% FDR and selecting only the probe with the highest average expression per gene giving the final transcriptional signature size of 95, 203 and 16 genes for CMV+, CMV- and all infants, respectively. Each infant was represented either by sets of two samples (unstimulated and BCG stimulated either at Day 0 or Day 28) or by sets of four samples (unstimulated or BCG stimulated at Day 0 and Day 28). To avoid overfitting, we implemented a modified nested cross validation scheme such that only complete sample sets per infant were assigned to either test or training set at each splitting iteration during the cross-validation process.

Model training was performed using a neural network model as implemented in the nnet package through the caret interface in R<sup>70</sup>. In the outer loop, samples were split 50 times with replacement into training (about 70%) and test sets (about 30%) to evaluate model performance and feature importance. For model parameter tuning in the inner loop, each training set was split into training and validation sets using a leave-one-infant-out cross validation scheme and the AUROC was recorded as performance metric over a grid of 4 size and 4 decay parameter combinations. All accuracies stated in the manuscript are provided as balanced accuracies<sup>71</sup> to account for the imbalance of the case and control infants within our cohort.

To validate our classification results, an independent cohort of 10-week-old South African infants vaccinated with BCG at birth was obtained as raw

expression data from Gene Expression Omnibus by downloading GSE20716. Study batch was removed with the ComBat command in R using parametric adjustment, risk of TB as null model and the validation cohort as reference in order to avoid any bias on the validation cohort, For prediction of risk of TB, 112, 50 and 5 probes overlapped between Illumina HumanHT 12 V4 and Illumina HumanRef-8 V2 for within CMV-, CMV+ and the combined group respectively. CMV status prediction was performed using 55 overlapping probes which were differentially expressed between CMV negative and CMV strongly positive (ELISpot >100 SFC/million) infants at an FDR of 20%. Only infants with at least two positive samples were labelled as suspected CMV+.

### **Data availability**

Raw and normalized expression data have been deposited at Gene Expression Omnibus under the accession number GSE98550

### **Acknowledgements**

We thank study participants and their families, the community of Cape Winelands East district, and South African Tuberculosis Vaccine Initiative (SATVI) personnel. This work was funded by Aeras and The Wellcome Trust with support from the European Commission within the 7th framework program (FP7) NEWTBVAC (Grant No. HEALTH-F3-2009-241745) and by the European Commission within Horizon2020 TBVAC2020 (Grant No. H2020 PHC-643381). HM is a Wellcome Trust Senior Clinical Research Fellow.

## REFERENCES

- 1 Tameris, M. D. *et al.* Safety and efficacy of MVA85A, a new tuberculosis vaccine, in infants previously vaccinated with BCG: a randomised, placebo-controlled phase 2b trial. *Lancet* **381**, 1021-1028, doi:10.1016/S0140-6736(13)60177-4 (2013).
- 2 Zak, D. E. *et al.* A blood RNA signature for tuberculosis disease risk: a prospective cohort study. *Lancet* **387**, 2312-2322, doi:10.1016/S0140-6736(15)01316-1 (2016).
- 3 Fletcher, H. A. *et al.* Human newborn bacille Calmette-Guerin vaccination and risk of tuberculosis disease: a case-control study. *BMC Med* **14**, 76, doi:10.1186/s12916-016-0617-3 (2016).
- 4 Brandes, U. *et al.* On Modularity Clustering. *IEEE Trans Knowl Data Eng* **20**, 172-188, doi:10.1109/TKDE.2007.190689 (2008).
- 5 Organisation, W. H. Global Tuberculosis Report 2016. (2016).
- 6 Lamb, G. S. & Starke, J. R. Tuberculosis in Infants and Children. *Microbiol Spectr* **5**, doi:10.1128/microbiolspec.TNMI7-0037-2016 (2017).
- 7 Hawkrige, A. *et al.* Efficacy of percutaneous versus intradermal BCG in the prevention of tuberculosis in South African infants: randomised trial. *BMJ* **337**, a2052, doi:10.1136/bmj.a2052 (2008).
- 8 Fletcher, H. A. *et al.* T-cell activation is an immune correlate of risk in BCG vaccinated infants. *Nat Commun* **7**, 11290, doi:10.1038/ncomms11290 (2016).
- 9 Klenerman, P. & Hill, A. T cells and viral persistence: lessons from diverse infections. *Nat Immunol* **6**, 873-879, doi:10.1038/ni1241 (2005).
- 10 Appay, V. & Sauce, D. Immune activation and inflammation in HIV-1 infection: causes and consequences. *J Pathol* **214**, 231-241, doi:10.1002/path.2276 (2008).
- 11 Wittkop, L. *et al.* Effect of cytomegalovirus-induced immune response, self antigen-induced immune response, and microbial translocation on chronic immune activation in successfully treated HIV type 1-infected patients: the ANRS CO3 Aquitaine Cohort. *The Journal of infectious diseases* **207**, 622-627, doi:10.1093/infdis/jis732 (2013).
- 12 Evans, T. G. *et al.* Expansion of the CD57 subset of CD8 T cells in HIV-1 infection is related to CMV serostatus. *Aids* **13**, 1139-1141 (1999).
- 13 Gronborg, H. L., Jespersen, S., Honge, B. L., Jensen-Fangel, S. & Wejse, C. Review of cytomegalovirus coinfection in HIV-infected individuals in Africa. *Rev Med Virol* **27**, doi:10.1002/rmv.1907 (2017).
- 14 Gianella, S. *et al.* Replication of Human Herpesviruses Is Associated with Higher HIV DNA Levels during Antiretroviral Therapy Started at Early Phases of HIV Infection. *J Virol* **90**, 3944-3952, doi:10.1128/JVI.02638-15 (2016).
- 15 Klatt, N. R., Chomont, N., Douek, D. C. & Deeks, S. G. Immune activation and HIV persistence: implications for curative approaches to HIV infection. *Immunol Rev* **254**, 326-342, doi:10.1111/imr.12065 (2013).
- 16 Adland, E., Klenerman, P., Goulder, P. & Matthews, P. C. Ongoing burden of disease and mortality from HIV/CMV coinfection in Africa in the

- antiretroviral therapy era. *Front Microbiol* **6**, 1016, doi:10.3389/fmicb.2015.01016 (2015).
- 17 Cobelens, F., Nagelkerke, N. & Fletcher, H. The convergent epidemiology of tuberculosis and human cytomegalovirus infection. *F1000Res* **7**, 280, doi:10.12688/f1000research.14184.2 (2018).
- 18 Olaleye, O. D., Omilabu, S. A. & Baba, S. S. Cytomegalovirus infection among tuberculosis patients in a chest hospital in Nigeria. *Comparative immunology, microbiology and infectious diseases* **13**, 101-106 (1990).
- 19 Sirenko, I. A. *et al.* Impact of cytomegalovirus infection on the course of tuberculosis in children and adolescents. *Problemy tuberkuleza i boleznei legkikh*, 7-9 (2003).
- 20 Stockdale, L. *et al.* Human cytomegalovirus epidemiology and relationship to tuberculosis and cardiovascular disease risk factors in a rural Ugandan cohort. *PloS one* **13**, e0192086 (2018).
- 21 Stockdale, L. *et al.* HIV, HCMV and mycobacterial antibody levels: a cross-sectional study in a rural Ugandan cohort. *Tropical medicine & international health : TM & IH* **24**, 247-257, doi:10.1111/tmi.13188 (2019).
- 22 Hertoghs, K. M. *et al.* Molecular profiling of cytomegalovirus-induced human CD8+ T cell differentiation. *J Clin Invest* **120**, 4077-4090, doi:10.1172/JCI42758 (2010).
- 23 Triebel, F. *et al.* LAG-3, a novel lymphocyte activation gene closely related to CD4. *J Exp Med* **171**, 1393-1405 (1990).
- 24 Djaoud, Z. *et al.* Amplified NKG2C+ NK cells in cytomegalovirus (CMV) infection preferentially express killer cell Ig-like receptor 2DL: functional impact in controlling CMV-infected dendritic cells. *J Immunol* **191**, 2708-2716, doi:10.4049/jimmunol.1301138 (2013).
- 25 Nielsen, C. M., White, M. J., Goodier, M. R. & Riley, E. M. Functional Significance of CD57 Expression on Human NK Cells and Relevance to Disease. *Front Immunol* **4**, 422, doi:10.3389/fimmu.2013.00422 (2013).
- 26 van Stijn, A. *et al.* Human cytomegalovirus infection induces a rapid and sustained change in the expression of NK cell receptors on CD8+ T cells. *J Immunol* **180**, 4550-4560 (2008).
- 27 Graham, S. M. *et al.* A prospective study of endothelial activation biomarkers, including plasma angiopoietin-1 and angiopoietin-2, in Kenyan women initiating antiretroviral therapy. *BMC Infect Dis* **13**, 263, doi:10.1186/1471-2334-13-263 (2013).
- 28 Legat, A., Speiser, D. E., Pircher, H., Zehn, D. & Fuertes Marraco, S. A. Inhibitory Receptor Expression Depends More Dominantly on Differentiation and Activation than "Exhaustion" of Human CD8 T Cells. *Front Immunol* **4**, 455, doi:10.3389/fimmu.2013.00455 (2013).
- 29 Netea, M. G. *et al.* Interleukin-32 induces the differentiation of monocytes into macrophage-like cells. *Proc Natl Acad Sci U S A* **105**, 3515-3520, doi:10.1073/pnas.0712381105 (2008).
- 30 Kagina, B. M. *et al.* Specific T cell frequency and cytokine expression profile do not correlate with protection against tuberculosis after bacillus Calmette-Guerin vaccination of newborns. *Am J Respir Crit Care Med* **182**, 1073-1079, doi:10.1164/rccm.201003-0334OC (2010).

- 31 Miles, D. J. *et al.* CD4(+) T cell responses to cytomegalovirus in early life: a prospective birth cohort study. *The Journal of infectious diseases* **197**, 658-662, doi:10.1086/527418 (2008).
- 32 Miles, D. J. *et al.* Cytomegalovirus infection in Gambian infants leads to profound CD8 T-cell differentiation. *J Virol* **81**, 5766-5776, doi:10.1128/JVI.00052-07 (2007).
- 33 Bai, X. *et al.* IL-32 is a host protective cytokine against Mycobacterium tuberculosis in differentiated THP-1 human macrophages. *J Immunol* **184**, 3830-3840, doi:10.4049/jimmunol.0901913 (2010).
- 34 Bai, X. *et al.* Human IL-32 expression protects mice against a hypervirulent strain of Mycobacterium tuberculosis. *Proc Natl Acad Sci U S A* **112**, 5111-5116, doi:10.1073/pnas.1424302112 (2015).
- 35 Bozzano, F. *et al.* Functionally relevant decreases in activatory receptor expression on NK cells are associated with pulmonary tuberculosis in vivo and persist after successful treatment. *Int Immunol* **21**, 779-791, doi:10.1093/intimm/dxp046 (2009).
- 36 Esin, S. & Batoni, G. Natural killer cells: a coherent model for their functional role in Mycobacterium tuberculosis infection. *J Innate Immun* **7**, 11-24, doi:10.1159/000363321 (2015).
- 37 Lu, C. C. *et al.* NK cells kill mycobacteria directly by releasing perforin and granulysin. *J Leukoc Biol* **96**, 1119-1129, doi:10.1189/jlb.4A0713-363RR (2014).
- 38 Montoya, D. *et al.* IL-32 is a molecular marker of a host defense network in human tuberculosis. *Sci Transl Med* **6**, 250ra114, doi:10.1126/scitranslmed.3009546 (2014).
- 39 Portevin, D., Via, L. E., Eum, S. & Young, D. Natural killer cells are recruited during pulmonary tuberculosis and their ex vivo responses to mycobacteria vary between healthy human donors in association with KIR haplotype. *Cell Microbiol* **14**, 1734-1744, doi:10.1111/j.1462-5822.2012.01834.x (2012).
- 40 López-Botet, M. *et al.* Human Cytomegalovirus Infection Is Associated with Increased Proportions of NK Cells That Express the CD94/NKG2C Receptor in Aviremic HIV-1-Positive Patients. *The Journal of Infectious Diseases* **194**, 38-41, doi:10.1086/504719 (2006).
- 41 Goncalves, A. *et al.* Differential frequency of NKG2C/KLRC2 deletion in distinct African populations and susceptibility to Trachoma: a new method for imputation of KLRC2 genotypes from SNP genotyping data. *Human genetics* **135**, 939-951 (2016).
- 42 Thomas, R. *et al.* NKG2C deletion is a risk factor of HIV infection. *AIDS research and human retroviruses* **28**, 844-851 (2012).
- 43 Goodier, M. R. *et al.* Rapid NK cell differentiation in a population with near-universal human cytomegalovirus infection is attenuated by NKG2C deletions. *Blood* **124**, 2213-2222 (2014).
- 44 Phillips, B. L. *et al.* LAG3 expression in active Mycobacterium tuberculosis infections. *Am J Pathol* **185**, 820-833, doi:10.1016/j.ajpath.2014.11.003 (2015).
- 45 Anigilaje, E. A., Dabit, J. O., Nweke, N. O. & Agbedeh, A. A. Prevalence and risk factors of cytomegalovirus infection among HIV-infected and HIV-

- exposed uninfected infants in Nigeria. *J Infect Dev Ctries* **9**, 977-987, doi:10.3855/jidc.6131 (2015).
- 46 Hsiao, N. Y., Zampoli, M., Morrow, B., Zar, H. J. & Hardie, D. Cytomegalovirus viraemia in HIV exposed and infected infants: prevalence and clinical utility for diagnosing CMV pneumonia. *J Clin Virol* **58**, 74-78, doi:10.1016/j.jcv.2013.05.002 (2013).
- 47 Piriou, E. *et al.* Early age at time of primary Epstein-Barr virus infection results in poorly controlled viral infection in infants from Western Kenya: clues to the etiology of endemic Burkitt lymphoma. *J Infect Dis* **205**, 906-913, doi:10.1093/infdis/jir872 (2012).
- 48 Jayasooriya, S. *et al.* Early virological and immunological events in asymptomatic Epstein-Barr virus infection in African children. *PLoS Pathog* **11**, e1004746, doi:10.1371/journal.ppat.1004746 (2015).
- 49 Della Chiesa, M. *et al.* Human cytomegalovirus infection promotes rapid maturation of NK cells expressing activating killer Ig-like receptor in patients transplanted with NKG2C<sup>-/-</sup> umbilical cord blood. *J Immunol* **192**, 1471-1479, doi:10.4049/jimmunol.1302053 (2014).
- 50 Strauss-Albee, D. M. *et al.* Human NK cell repertoire diversity reflects immune experience and correlates with viral susceptibility. *Sci Transl Med* **7**, 297ra115, doi:10.1126/scitranslmed.aac5722 (2015).
- 51 Slyker, J. A. *et al.* Acute cytomegalovirus infection in Kenyan HIV-infected infants. *Aids* **23**, 2173-2181, doi:10.1097/QAD.0b013e32833016e8 (2009).
- 52 Berry, M. P. *et al.* An interferon-inducible neutrophil-driven blood transcriptional signature in human tuberculosis. *Nature* **466**, 973-977, doi:10.1038/nature09247 (2010).
- 53 Cliff, J. M. *et al.* Distinct phases of blood gene expression pattern through tuberculosis treatment reflect modulation of the humoral immune response. *The Journal of infectious diseases* **207**, 18-29, doi:10.1093/infdis/jis499 (2013).
- 54 Maertzdorf, J. *et al.* Human gene expression profiles of susceptibility and resistance in tuberculosis. *Genes and immunity* **12**, 15-22, doi:10.1038/gene.2010.51 (2011).
- 55 Anderson, S. T. *et al.* Diagnosis of childhood tuberculosis and host RNA expression in Africa. *The New England journal of medicine* **370**, 1712-1723, doi:10.1056/NEJMoa1303657 (2014).
- 56 Kaforou, M. *et al.* Detection of tuberculosis in HIV-infected and -uninfected African adults using whole blood RNA expression signatures: a case-control study. *PLoS medicine* **10**, e1001538, doi:10.1371/journal.pmed.1001538 (2013).
- 57 Scriba, T. J. *et al.* Sequential inflammatory processes define human progression from M. tuberculosis infection to tuberculosis disease. *PLoS Pathog* **13**, e1006687, doi:10.1371/journal.ppat.1006687 (2017).
- 58 Soares, A. P. *et al.* Longitudinal changes in CD4(+) T-cell memory responses induced by BCG vaccination of newborns. *The Journal of infectious diseases* **207**, 1084-1094, doi:10.1093/infdis/jis941 (2013).
- 59 Ben-Smith, A. *et al.* Differences between naive and memory T cell phenotype in Malawian and UK adolescents: a role for Cytomegalovirus? *BMC Infect Dis* **8**, 139, doi:10.1186/1471-2334-8-139 (2008).



- 60 Miles, D. J. *et al.* Human immunodeficiency virus (HIV) infection during pregnancy induces CD4 T-cell differentiation and modulates responses to Bacille Calmette-Guerin (BCG) vaccine in HIV-uninfected infants. *Immunology* **129**, 446-454, doi:10.1111/j.1365-2567.2009.03186.x (2010).
- 61 Matsumiya, M. *et al.* Inflammatory and myeloid-associated gene expression before and one day after infant vaccination with MVA85A correlates with induction of a T cell response. *BMC Infect Dis* **14**, 314, doi:10.1186/1471-2334-14-314 (2014).
- 62 Tanner, R. *et al.* Serum indoleamine 2,3-dioxygenase activity is associated with reduced immunogenicity following vaccination with MVA85A. *BMC Infect Dis* **14**, 660, doi:10.1186/s12879-014-0660-7 (2014).
- 63 Ota, M. O. *et al.* Immunogenicity of the tuberculosis vaccine MVA85A is reduced by coadministration with EPI vaccines in a randomized controlled trial in Gambian infants. *Sci Transl Med* **3**, 88ra56, doi:10.1126/scitranslmed.3002461 (2011).
- 64 Dunning, M. J. *et al.* The importance of platform annotation in interpreting microarray data. *Lancet Oncol* **11**, 717, doi:10.1016/S1470-2045(10)70115-7 (2010).
- 65 Shi, W., Oshlack, A. & Smyth, G. K. Optimizing the noise versus bias trade-off for Illumina whole genome expression BeadChips. *Nucleic Acids Res* **38**, e204, doi:10.1093/nar/gkq871 (2010).
- 66 Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* **43**, e47, doi:10.1093/nar/gkv007 (2015).
- 67 Ritchie, M. E. *et al.* Empirical array quality weights in the analysis of microarray data. *BMC Bioinformatics* **7**, 261, doi:10.1186/1471-2105-7-261 (2006).
- 68 Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate - a Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B-Methodological* **57**, 289-300 (1995).
- 69 Li, S. *et al.* Molecular signatures of antibody responses derived from a systems biology study of five human vaccines. *Nat Immunol* **15**, 195-204, doi:10.1038/ni.2789 (2014).
- 70 Kuhn, M. caret Package. *J. Stat. Softw.* **28**, 1-26 (2008).
- 71 Brodersen, K. H., Ong, C. S. x., Stephan, K. E. & Buhmann, J. M. The balanced accuracy and its posterior distribution. *Pattern Recognition (ICPR), 2010 20th International Conference on*, 3121-3124 doi:DOI 10.1109/ICPR.2010.764 (2010).