

620 Supplementary Figures

621

622 **Supplementary Figure 1: Details of the tuberculosis risk signature.** (A) Network
623 representation of the tuberculosis risk signature. Nodes (circles) represent splice junctions from
624 genes that comprise the signature. Node colors indicate to which genes the splice junctions
625 belong. Lines indicate splice junction discriminant pairs that make up the overall signature.
626 Node size is proportional to the number of discriminant pairs that contain the splice junction. All
627 splice junctions were expressed at higher levels in progressors than controls. (B) Network
628 representation of the tuberculosis risk scores during the transition from latent to active
629 tuberculosis for a given progressor. Pair-wise discriminants that vote “progressor” are indicated
630 as red lines and those that vote “control” are indicated as green lines. The overall tuberculosis
631 risk score is the proportion of all pair-wise discriminants that vote progressor. The risk scores,
632 expressed as percentages, are shown for progressor samples at the indicated time points before
633 disease diagnosis, or from a sample from a single matched healthy control at study day 0.

634

635 **Supplementary Figure 2: Benchmarking the predictive performance of the SVM-based**
636 **tuberculosis risk signature against a risk signature derived from use of Random Forests.** To
637 benchmark prediction performance of the risk signature against an alternative data mining
638 method, a Random Forest model consisting of 100,000 trees (100k RF) was generated from ACS
639 Training set RNA-Seq data. (A-D) Performance of the Random Forest model was compared to
640 the PSVM-based tuberculosis risk signature in cross-validation of the ACS Training set,
641 stratified by time preceding tuberculosis disease onset: (A) 1-180 days before tuberculosis, (B)
642 181-360 days before tuberculosis, (C) 361-540 days before tuberculosis, and (D) 541-720 days
643 before tuberculosis. (E) Despite the superiority of 100k RF in cross-validation, the prediction
644 accuracy on the ACS test set was indistinguishable for the two methods.

645

646 **Supplementary Figure 3: Performance of the tuberculosis risk signature for diagnosing**
647 **active adult and childhood tuberculosis disease using published microarray datasets.** (A)
648 Analytical strategy for evaluating the tuberculosis risk signature as a signature for active
649 tuberculosis disease using published microarray datasets. The RNA-Seq-based signature was re-
650 parameterized using whole blood microarray data from the UK training set of Berry, *et al.*,¹¹

651 which includes samples from tuberculosis cases and latently *M. tuberculosis* infected controls.
652 The fully locked-down microarray-based signature was then employed to make quasi-blind
653 predictions on data from other cohorts from the same and independent studies.⁹⁻¹³ (B-H) ROC
654 curves depicting prediction performance of the Illumina microarray-based risk signature on
655 cohorts from published studies are shown. (B) Discrimination of active tuberculosis disease from
656 latent *M. tuberculosis* infection and healthy controls in the South African test set and UK test set
657 from Berry, *et al.*¹¹; (C) Discrimination of active tuberculosis disease from latent *M. tuberculosis*
658 infection in the presence or absence of HIV co-infection, using data from Kaforou, *et al.*¹²; (D)
659 Discrimination of active tuberculosis disease from lung cancer, pneumonia, and sarcoidosis
660 using data from Bloom, *et al.*¹⁰; (E) Discrimination of active tuberculosis disease from other
661 pulmonary diseases in the presence or absence of HIV co-infection, using data from Kaforou, *et*
662 *al.*¹²; (F) Discrimination of culture positive and culture negative childhood tuberculosis disease
663 from latent *M. tuberculosis* infection, using data from Anderson, *et al.*¹³; (G) Discrimination of
664 culture positive childhood tuberculosis disease from other diseases in presence or absence of
665 HIV co-infection, using data from Anderson, *et al.*¹³; (H) Discrimination of active tuberculosis
666 disease from healthy individuals over the course of antimicrobial treatment (red lines), and
667 discrimination of baseline tuberculosis disease from tuberculosis disease after 2 weeks of
668 treatment (purple line), using data from Bloom, *et al.*⁹

669 **Supplementary Appendices**

670

- 671 1. Supplementary Appendix 1: Exclusion criteria for the adolescent cohort study (ACS)
- 672 2. Supplementary Appendix 2: Exclusion criteria for the GC6-74 study
- 673 3. Supplementary Appendix 3: Regulatory authorities that approved the ACS and GC6-74
- 674 4. Supplementary Appendix 4: Supplementary Methods

675

676 **Supplementary Appendix 1: Exclusion criteria for the adolescent cohort study**

677

678 Parent study:

- 679 • Pregnant or lactating women
- 680 • Any reported acute or chronic medical condition resulting in hospitalization within 6
- 681 months prior to enrollment.

682 Case-control study:

- 683 • Refer to main text

684 **Supplementary Appendix 2: Exclusion criteria for the GC6-74 study**

685

686 Parent study:

- 687 • Resident in the study area for less than 3 months prior to enrollment
- 688 • No permanent address
- 689 • Previous treatment for tuberculosis
- 690 • Previous or current anti-retroviral therapy
- 691 • Current participation in a drug or vaccine trial, or participation within 6 months of
- 692 enrollment
- 693 • Concomitant cancer
- 694 • Diabetes mellitus
- 695 • Chronic bronchitis/emphysema/asthma requiring systemic steroid therapy
- 696 • Other steroid therapy within 6 months of enrollment
- 697 • Current pregnancy, or pregnancy within 3 months of enrollment

698

699 Case-control study:

- 700 • Refer to main text

701 **Supplementary Appendix 3: Regulatory authorities that approved the ACS and GC6-74**
702 **protocols**

703

704 ACS:

705

- 706 • University of Cape Town Research Ethics Committee, Cape Town, South Africa

707

708

709 GC6-74 South African and Gambian sites:

- 710 • Scientific and Ethics Committee of the Faculty of Health Sciences of Stellenbosch
711 University, Cape Town, South Africa

- 712 • Joint Medical Research Council and Gambian Government ethics review committee,
713 Banjul, The Gambia

714 **Supplementary Appendix 4: Supplementary Methods**

715

716 A. Detailed definition of cases and controls in the ACS for identifying and validating
717 signatures of tuberculosis risk

718 B. Sequencing of whole blood transcriptomes

719 C. Derivation of the tuberculosis risk signature from the ACS training RNA-Seq dataset

720 D. Adaptation of the tuberculosis risk signature from RNA-Seq to qRT-PCR

721 E. Adaptation of the tuberculosis risk signature to the Illumina microarray platform

722

723 **A. Detailed definition of cases and controls in the ACS and in GC6-74 for identifying**
724 **and validating signatures of tuberculosis risk**

725

726 The ACS determined the prevalence and incidence of *M. tuberculosis* infection and disease
727 among adolescents from the Cape Town region of South Africa.^{30,37} Overall, 53% of ACS
728 participants had latent *M. tuberculosis* infection at enrollment. For the ACS signature of risk
729 study, adolescents with latent *M. tuberculosis* infection at enrollment were eligible; *M.*
730 *tuberculosis* infection was diagnosed by a positive QuantiFERON TB GOLD In-Tube Assay
731 (QFT, Cellestis; >0.35 IU/mL) and/or a positive tuberculin skin test (TST, 0.1mL dose of
732 Purified Protein Derivative RT-23, 2-TU, Statens Serum Institute; >10mm). According to South
733 African policy, these QFT and/or TST positive adolescents were not given therapy to prevent
734 tuberculosis disease.²¹

735

736 Adolescents who developed active tuberculosis disease during 2 years of follow-up were
737 included as “progressors” (cases). Participants that were either exposed to tuberculosis patients,
738 or had symptoms suggestive of tuberculosis, were evaluated clinically and by sputum smear,
739 culture and chest roentgenography. Tuberculosis was defined as intrathoracic disease, with either
740 two sputum smears positive for acid-fast bacilli or one positive sputum culture confirmed as *M.*
741 *tuberculosis* complex (mycobacterial growth indicator tube, BD BioSciences). Participants who
742 developed tuberculosis within 6 months of enrolment were excluded on the basis that they may
743 represent individuals with active but as yet asymptomatic tuberculosis disease.

744

745 Five ACS participants who were not infected with *M. tuberculosis* at enrollment but who
746 converted to a positive QFT and/or TST, and ultimately developed tuberculosis disease at least 6
747 months post QFT/TST conversion, were also included as progressors. As a subset of ACS
748 participants from the parent study, who had a negative QFT at baseline, were followed for
749 incident tuberculosis for up to 3 years after the last QFT (5 years in total) through biannual study
750 visits and passive surveillance of health facility records, the follow-up of these participants was
751 longer than the 2 years applying to most participants³⁸.

752
753 All ACS patients with tuberculosis disease were offered a HIV test; HIV infected patients were
754 excluded from the case controls study. HIV testing of healthy study participants was not
755 permitted by the human research ethics committee of the University of Cape Town; this
756 committee also did not allow post-hoc, anonymous HIV testing. Regardless, the HIV incidence
757 rate in adolescents diagnosed with active tuberculosis was <2% (1 out of 61 who were offered
758 and accepted testing), and since HIV is a risk factor for tuberculosis, we expect the HIV
759 prevalence among healthy adolescents (from whom controls were identified) to be minute³⁰.

760
761 For each ACS progressor, two matched controls were identified. Controls were selected from
762 ACS participants that remained healthy during follow-up, and were matched to progressors by
763 age at enrolment, gender, ethnicity, school of attendance, and presence or absence of prior
764 episodes of tuberculosis disease.

765
766 Among GC6-74 participants, progressors were defined as having intrathoracic tuberculosis by
767 one of three categories: First, two positive sputum cultures (MGIT); second, one positive sputum
768 culture and/or a positive sputum smear, and clinical signs and symptoms compatible with
769 tuberculosis and/or a chest roentgenogram compatible with active pulmonary tuberculosis; third,
770 two positive sputum smears with clinical signs and symptoms compatible with tuberculosis or a
771 chest roentgenogram compatible with active pulmonary tuberculosis. Progressors were excluded
772 if they developed disease within 90 days of enrollment, for reasons mentioned above. Controls
773 were matched to progressors based on age category (<18, 18-25, 26-35, ≥36 years of age),
774 gender and year of enrollment.

775

776 *Time to tuberculosis*: For the ACS study, two *time to tuberculosis* values were calculated for
777 each progressor. First, original values were assigned early after sample collection and were
778 employed throughout signature construction. Second, per protocol values were assigned during
779 manuscript preparation when it was revealed that some original *time to diagnosis* assignments
780 had been wrong. All prediction results are reported in terms of the per protocol values.

781

782 **B. Sequencing of whole blood transcriptomes in the ACS**

783

784 ***RNA-Seq***: Whole blood was collected in PAXgene Blood RNA Tubes (PreAnalytiX), frozen,
785 and RNA was later extracted using of PAXgene Blood RNA Kits. Globin transcripts were
786 depleted and RNA was sequenced by Expression Analysis Inc. (Durham, NC) using a 30M read,
787 50bp paired-end sequencing strategy (60M reads/sample).

788

789 ***Alignment***: Read pairs were preprocessed to adjust base calls with phred scores <5 to N and to
790 remove read pairs for which either end has fewer than 30 unambiguous (non-N) base calls. This
791 method indirectly removes pairs containing adaptor sequences. The median depth of RNA
792 sequencing after post-processing was 31 million read pairs. Read pairs were aligned to the
793 human genome using the gsnap program²², allowing for novel splice junction detection. The
794 mean percentage of reads mapped was 90% and the mean GC content was 50.6%.

795

796 ***Splice junction filtering***: Gene expression abundance was measured at the level of splice
797 junction counts, which quantifies the relative frequency of specific mRNA splicing events in
798 expressed genes. This facilitates adaptation of the signatures from RNA-Seq to qRT-PCR
799 because, in practice, PCR primer sets are designed to target splice junctions as a means to
800 prevent amplification of contaminating genomic DNA. Splice junctions that were detected by at
801 least five reads in at least ten samples were retained for signature analysis, leaving 141,140
802 splice junctions. In total, 355 samples from the ACS cohort were analyzed by RNA-Seq
803 (Supplementary Table 6): 264 samples from the ACS training set and 91 samples from the ACS
804 test set. RNA-Seq analysis of the ACS training and test sets was performed independently, over a
805 year apart.

806

807 **Normalization:** Splice junction counts for each sample were first pre-normalized for library size
808 by adding “1” to the raw counts, dividing the counts in a given sample by the sum of all counts
809 in that sample, and then taking the logarithm (base 2). “Reference junctions” for use as internal
810 controls in all subsequent analyses were then identified from the 20 splice junctions with the
811 smallest coefficient of variance computed across all samples from the pre-normalized table. The
812 final normalized log₂-based splice junction table was finally constructed by subtracting the mean
813 of the reference junction counts for each sample. Reference junctions were identified by using
814 the 264 samples that comprise the full ACS training set RNA-Seq sample set (Supplementary
815 Table 6), which included a small number of samples that were collected after the initiation of
816 treatment (Supplementary Table 2). The set of reference junctions is provided in Supplementary
817 Table 11. Expression levels of the reference junctions and junctions comprising the tuberculosis
818 risk signature are provided in Supplementary Table 10. The Microsoft Excel-based worksheet
819 for computing tuberculosis risk scores (Supplementary Table 11) includes reference junction
820 normalization.

821

822 **C. Derivation of the tuberculosis risk signature from the ACS training RNA-Seq** 823 **dataset**

824

825 Mining the RNA-Seq data from the ACS training set to generate the tuberculosis risk signature
826 required (1) selecting the appropriate **mathematical framework** in which the signature would be
827 formulated, and (2) exploiting the **longitudinal structure** of the progression cohort to extract
828 and train the signature.

829

830 **Mathematical framework – SVM:** The mathematical framework for the signatures is a
831 generalization of the k-top-scoring pairs (k-TSP) methodology, was developed for discovery of
832 cancer biomarkers from microarray datasets.²³ Signatures derived using the k-TSP approach are
833 collections of gene-pair discriminators that can vote “progressor” (1) or “control” (0) (for
834 example). For a given sample, the classification “score” is the average of all of the “0” or “1”
835 votes computed for the whole collection of discriminators for that sample. In this manner, k-TSP
836 combines many “weak” discriminators to improve the reliability of the predictions. The pair-

837 wise discriminators underlying k-TSP are very simple, involving only a pair of genes for which
838 gene1 > gene2 in progressors and the reverse is true in controls (for example).

839
840 The k-TSP framework is desirable in the present study for three reasons. First, it has the potential
841 to identify combinations of genes that better predict progression than either gene individually, a
842 characteristic common to bivariate approaches.²⁴ Second, being based on an ensemble of models,
843 rather than a single model, the methodology is tolerant to failed measurements. For example, if a
844 particular primer fails for a particular sample, the overall score can still be computed from the
845 unaffected pairs. In this regard, k-TSP is similar to Random Forests.³⁹ Third, the underlying
846 models, involving only two genes, are parsimonious and are therefore unlikely to suffer from
847 overfitting.⁴⁰

848
849 In the present study, we replaced the simple rank-based gene pair models in k-TSP with linear
850 SVM gene pair discriminant models, and call the approach “PSVM” (pair-wise support vector
851 machine ensembles). This generalization allows for greater flexibility in the selection of gene
852 expression patterns that predict tuberculosis progression. While the k-TSP approach requires the
853 relative ranking of the genes to change between the two conditions (effectively favoring gene
854 pairs that are differentially expressed in opposite directions) any pair of genes that provides non-
855 redundant information for predicting tuberculosis can be combined in a linear SVM discriminant.
856 This was important for tuberculosis progression, where genes with the largest magnitude
857 expression differences between progressors and controls tend to be expressed higher in
858 progressors. By merging the k-TSP approach with SVMs, PSVM is similar to the k-TSP
859 modification proposed by Shi et al., (2011).²⁵ The difference between the method of Shi et al.
860 (2011)²⁵ and PSVM is that the former replaces the ensemble-based structure with a single SVM
861 model, while PSVM retains the ensemble structure and replaces the rank-based pairs with SVMs
862 internally.

863
864 ***Extracting and training the tuberculosis risk signature from the ACS training set:***
865 The PSVM approach was applied to the ACS training set in a manner designed to optimally
866 extract a predictive risk signature from the data. At a high level, the strategy involves three
867 steps (1) changing the cohort time scale to reflect the time before diagnosis with tuberculosis

868 instead of time since enrolment (since the enrolment time is arbitrary and not related to clinical
869 outcome); (2) deriving a set of candidate predictor genes by comparing gene expression in
870 controls with gene expression in progressors at the time points most proximal to diagnosis; and
871 (3) filtering the candidate predictor genes to retain only those that, in pairwise combinations,
872 robustly discriminate progressors from controls at time points that are more distal to tuberculosis
873 diagnosis. Although this derivation of the tuberculosis risk signature employs multiple time
874 points in the gene selection, applying the signature to assess tuberculosis disease risk only
875 requires expression measurements for the final selected set of genes at single time point and is
876 done in a blinded manner, without any requirement for sample metadata or participant
877 information.

878
879 Although it may appear counter-intuitive that samples collected proximally to diagnosis are used
880 to generate candidate genes for the signature of risk, initial cross-validation analyses
881 demonstrated that this approach yields a more predictive signature than approaches that rely
882 solely on samples that are distal to tuberculosis to generate the signature. One interpretation is
883 that while a strong coherent signal is necessary to discover a predictive signature, once
884 developed, the signature can correctly predict tuberculosis on samples where the signal is
885 weaker. This is because less signal is required to recognize an established gene expression
886 pattern than to discover a new one.

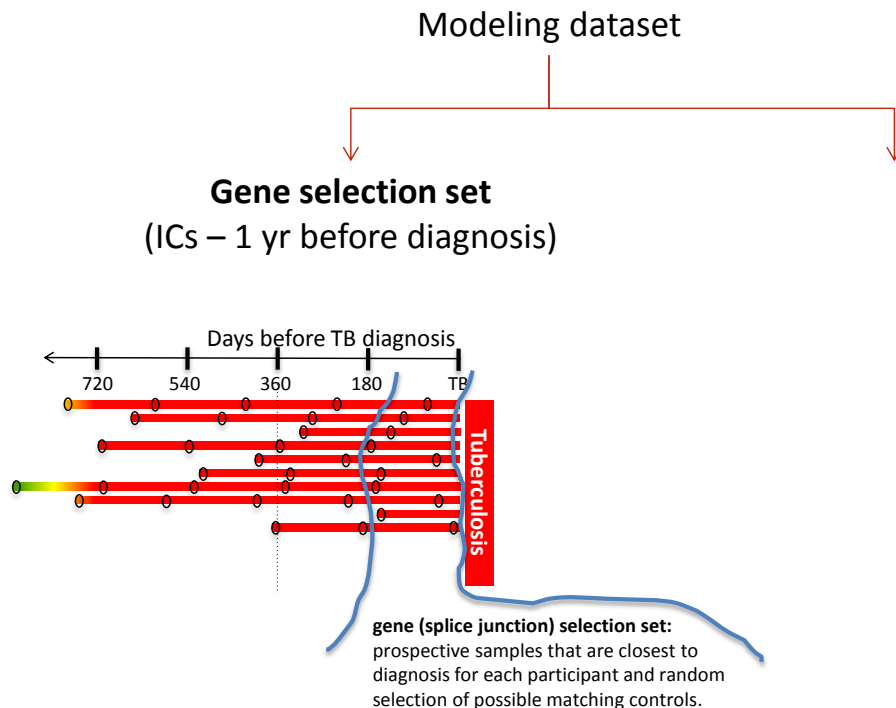
887
888 Having described the signature development process at a high level above, the following is a
889 description of the process in detail. Development of the signature starts from the normalized
890 splice junction count table of gene expression for the ACS training set. This same procedure was
891 employed for both development of the final tuberculosis risk signature and for cross-validation
892 of the signature generation method. In the case of the former, the signature generation procedure
893 is applied to the entire ACS training set. For the latter, the procedure is applied to the ACS
894 training set after having excluded 1/5 of progressors and matched controls. For simplicity, the
895 dataset that is used for generating the signature will be referred to as the “Modeling dataset” for
896 either case.

897

898 **(1) Re-aligning the Modeling dataset time scale in terms of diagnosis with tuberculosis:** The
899 dataset was first synchronized according to the tuberculosis diagnosis time point instead of the
900 study enrollment time point (Figure 2B). This allowed identification of which progressor
901 samples (and matched control samples) that had been collected most proximally and most
902 distally to diagnosis.

903
904 **(2) Defining the gene (splice junction) selection set:** This first partition of the data defines
905 which samples are be used to select the genes (splice junctions) that are evaluated as candidates
906 predictors in the signature. The partition was made in the following manner:

- 907 • Only one sample from any progressor or control was included in the junction selection
908 set.
- 909 • For subjects with at least one sample less than a year before diagnosis (original estimates,
910 Supplementary Table 2) or at the time of diagnosis (IC = incident case samples,
911 Supplementary Table 2), the sample closest to or at diagnosis was put into the junction
912 selection set.
- 913 • Matching control samples for each progressor were randomly selected from possible
914 matches (same blood draw time and demographic bin as the progressor).

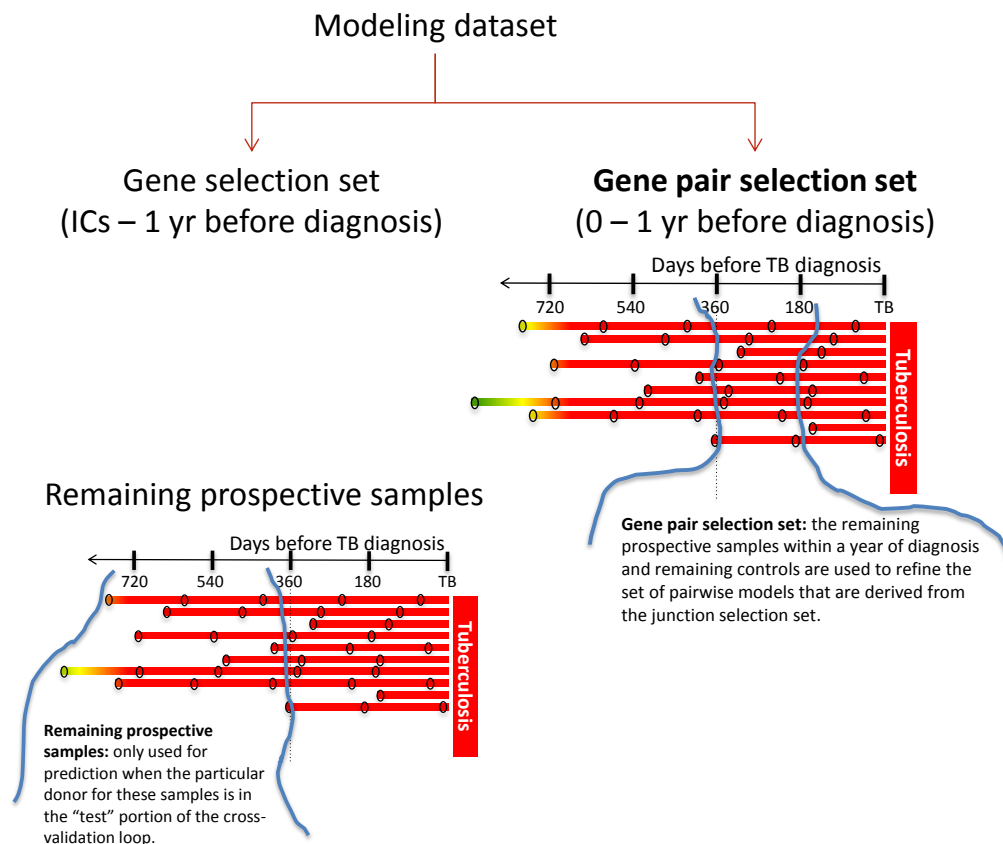


916

917 **(3) Defining the gene pair selection set:**

918 This step defines which samples are used to filter gene pairs for accurate discrimination of
919 progressors from controls at time points that are more distal to diagnosis.

- 920 • All control samples and pre-diagnosis progressor samples within a year of tuberculosis
921 diagnosis (by original estimates, Supplementary Table 2) that remained after defining the
922 junction selection set comprised the pair selection set.
- 923 • Pre-diagnosis progressor samples that were collected over a year before tuberculosis
924 diagnosis were evaluated in cross-validation but were not used for gene or gene pair
925 selection.



926

927 **(4) Gene (splice junction) selection:** Splice junctions that significantly discriminate progressor
928 samples from matched control samples within the gene selection set were identified by
929 permutation test. Briefly, t-statistics comparing progressors and controls were computed for
930 each splice junction. The dataset was then randomly permuted 1,000,000 times. Splice junctions

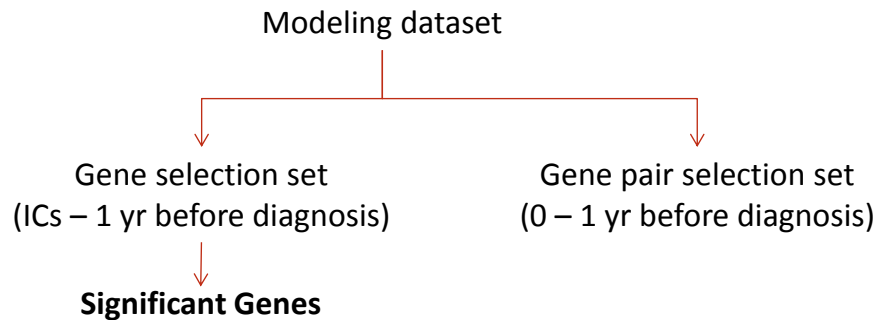
931 with t-statistics that significantly deviated from the empirical distribution (FDR<0.0001) were
932 retained.

933

934

935

936



Gene (splice junction) selection

- Compute t-statistic comparing progressors to controls
- Estimate significance by 1,000,000 permutations
- Keep Genes with FDR < 0.0001

937

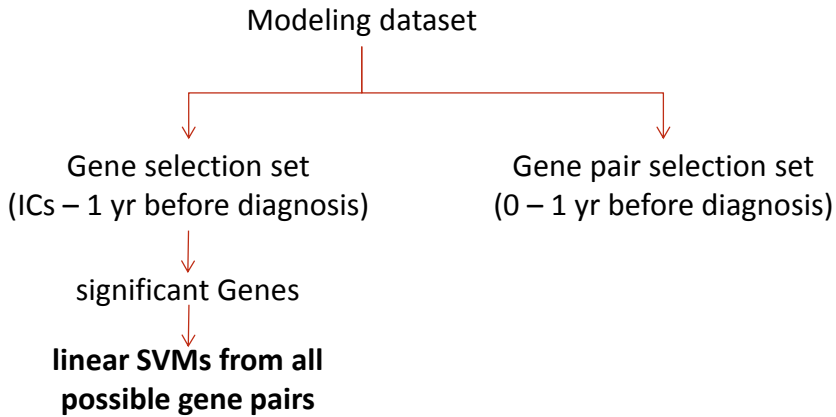
938

939

(5) Gene pair model fitting:

941 Linear models that discriminate progressors from controls for all possible pair-wise
942 combinations of the selected genes were then fit to the gene selection data set using the
943 Sequential Minimal Optimization SVM algorithm⁴¹ in C++. Fitting of each pair-wise linear
944 SVM to the entire gene selection set was performed independently.

945



Gene pair model fitting

- Construct an ensemble of Support Vector Machine (SVM) models for all possible pairs
- SVM Gene pair models are trained using **Gene selection** set data

946

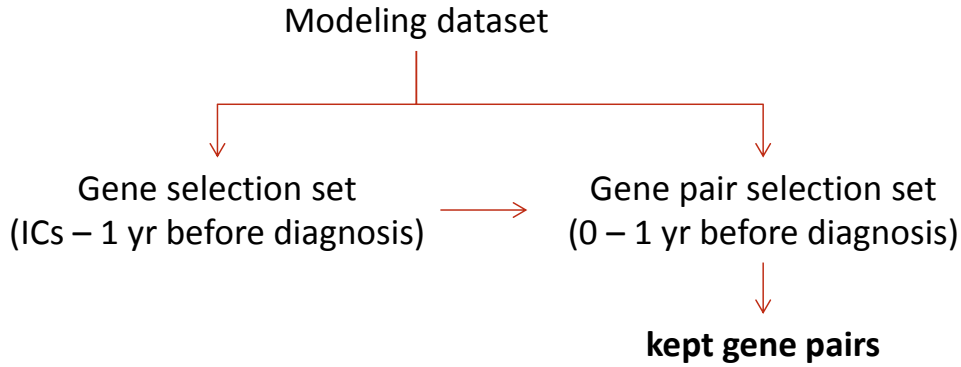
947

948

(6) Gene pair scoring:

950 Each gene pair model was used to predict on the gene pair selection set. Pairs that correctly
951 predicted 70% of control samples, 70% of progressor samples from within 6 months of diagnosis
952 (by original estimates, Supplementary Table 2), and 70% of progressor samples from between 6
953 months and 1 year before diagnosis (by original estimates, Supplementary Table 2) in the pair
954 selection set were given a score of “1”. Junction pairs that did not meet these thresholds were
955 given a score of “0”.

956



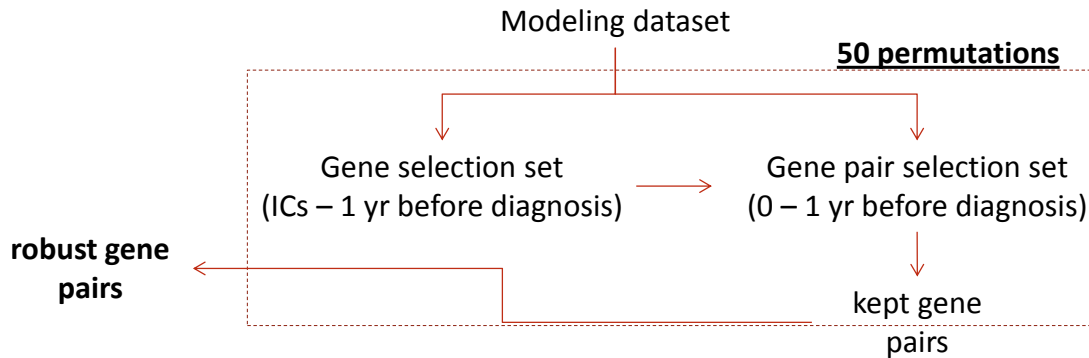
Gene pair scoring

- Test the ability of gene pair models to predict on the gene pair selection set
- Evaluate 3 groups: controls, 0-0.5 yr progressors, 0.5-1 yr progressors
- Retain only those gene pairs that achieve 70% accuracy in all groups

957
958
959

(7) Gene pair selection:

961 After completing one iteration of steps 1-6, the process was repeated starting from a different
962 randomized matching of controls to progressors, yielding slightly different gene selection and
963 gene pair selection sets (steps 1-2), and a different set of scores for the retained gene pairs (step
964 6). The process was repeated 50 times, with junction pairs receiving scores of 0-50. Gene pairs
965 with scores ≥ 45 were retained for the final ensemble. In this manner, the ensemble is comprised
966 of the most robustly predictive gene pairs.



Gene pair selection

- Repeat the Gene selection/Gene pair selection split 50X (randomized progressor/control matching)
- Retain Gene pairs that are selected in at least 45/50 permutations

967
968

(8) Final model parameterization:

969 The gene pairs that were retained after step 6 comprise the PSVM ensemble. Final linear
970 discriminants for each of these pairs were parameterized using the entire set of progressor
971 samples from the Modeling set that were collected less than a year before diagnosis (as indicated
972 by original estimates, Supplementary Table 2) or at the time of diagnosis (IC = incident case
973 samples, Supplementary Table 2) and matched control samples using the Sequential Minimal
974 Optimization SVM algorithm⁴¹ in C++ . As in step 4, fitting of each pair-wise linear SVM to the
975 Modeling set was performed independently.

976

977 **(9) Estimating the prediction accuracy of the TB risk signature by cross-validation:** Prior to
978 making predictions on the ACS test set or adapting the TB risk signature to qRT-PCR, the
979 predictive accuracy of the approach was estimated by performing 100 iterations of cross-
980 validation involving random 4:1 splits of ACS training set participants into training and
981 prediction sets. In each iteration, the entire procedure was repeated starting from the entire ACS
982 training junction expression matrix, ensuring unbiased estimates of prediction accuracy.

983

984 **(10) Benchmarking the signature performance against Random Forest.** Prior to testing the
985 PSVM signature on the GC6-74 dataset, prediction performance of the signature was
986 benchmarked against the performance of a Random Forest³⁹ signature consisting of 100,000
987 trees. We compared the prediction performance of the 100k RF signature to the PSVM signature
988 on two levels: cross-validation of the ACS training set RNA-Seq data (100 iterations of 5-fold
989 splits), and prediction of the ACS test set RNA-Seq data. For both analyses, RF models were
990 constructed in two-step procedure: (1) junction selection was performed by permutation analysis
991 as for the PSVM signature, except that all progressor samples within one year of diagnosis and
992 averaged matched controls were used to compute the t-statistics. (2) RF models were generated
993 from the selected junctions using the R package *randomForest*⁴², using default settings, except
994 specifying 100,000 trees. Applying this approach to the entire ACS training set generated the
995 final 100k RF model, which consisted of 1,712 junctions and 631 genes.

996

997 **D. Adaptation of the tuberculosis risk signature from RNA-Seq to qRT-PCR for**
998 **validation in the ACS and GC6-74**

999

1000 The signature was adapted to the qRT-PCR platform by matching each splice junction to Applied
1001 BioSystems TaqMan® gene expression assays in the following manner. Exact matches between
1002 splice junctions and TaqMan assays were identified first and selected. When exact matches were
1003 not commercially available, an attempt was made to construct custom Taqman assays spanning
1004 the desired junction. If the custom-design process failed, the commercial assay was selected that
1005 matched the splice junction of the same gene that was the most strongly correlated to the splice
1006 junction of interest. The resulting matches between RNA-Seq designed splice junctions and
1007 commercial Taqman assays are provided in Supplementary Tables 17-21. Parameterization of the
1008 qRT-PCR-based versions of the signature was performed using qRT-PCR expression data for the
1009 ACS Training set and retaining the ensemble structures. Normalization of the cycle threshold
1010 data was performed by comparing expression of the signature genes to the set of reference genes.
1011 The qRT-PCR-based risk signature was finally generated by re-training the pairwise SVM
1012 models to the normalized Ct data using the network structures obtained from RNA-Seq.

1013
1014 After validation of the signatures of risk on the ACS test set, qRT-PCR data for the PSVM
1015 primers and reference genes was generated from GC6-74 cohort RNA, in a blinded manner, as
1016 described above. Prior to predicting on GC6-74 RNA samples, two modifications were made to
1017 the signature. First, failure of one reference primer (GRK6) on the GC6-74 samples necessitated
1018 exclusion of this primer and re-parameterization of the signatures (using ACS training set data
1019 only). Second, post-hoc inspection of signature prediction on the ACS test set identified a subset
1020 of SVM pairs that always voted progressor or always voted control, irrespective of the sample.
1021 These pairs were pruned from the networks prior to predicting on GC6-74.

1022
1023 **E. Adaptation of tuberculosis risk signature to the Illumina microarray and predicting**
1024 **on published microarray datasets**

1025
1026 *Uniform normalization of published microarray datasets:* Microarray datasets from the
1027 published studies of interest were obtained from the Gene Expression Omnibus.⁴³ For most
1028 studies¹¹⁻¹³, the “Series Matrix File (TXT)” was downloaded. For Bloom, *et al.*⁹,
1029 “GSE40553_non-normalized_SALong.txt.gz” was downloaded. For Bloom, *et al.*¹⁰,
1030 “GSE42825_non-normalized.txt.gz”, “GSE42826_non-normalized.txt.gz” and “GSE42830_non-

1031 normalized.txt.gz” were downloaded and concatenated. Expression data matrices for each study
1032 were converted to a log₂ format as needed. For each downloaded dataset, the Illumina
1033 microarray probes for all genes present in the set of reference splice junctions (Supplementary
1034 Table 11) were identified. Each sample was then normalized by subtracting the average
1035 expression of genes detected by these probes from the log₂-based data.

1036

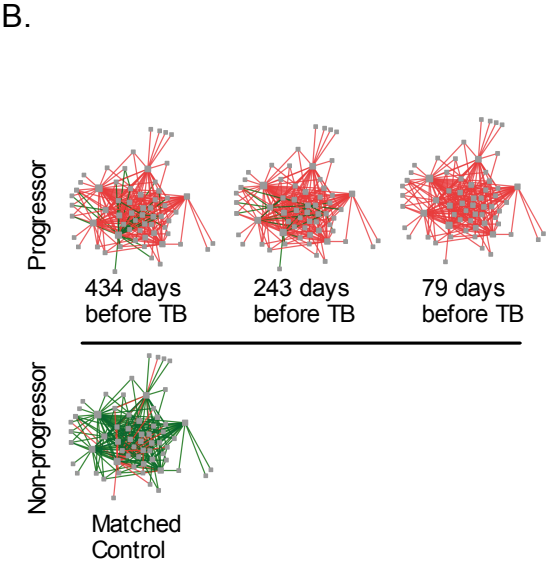
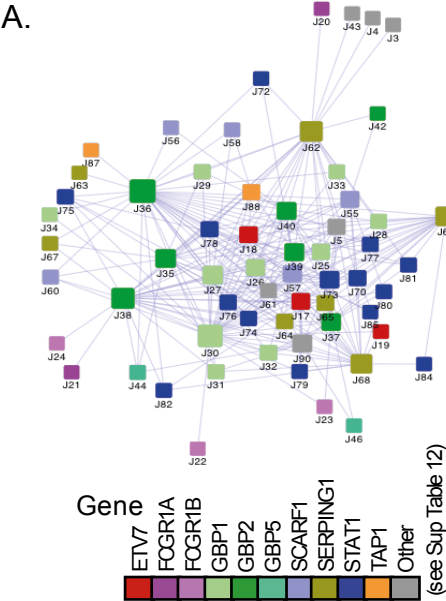
1037 ***Microarray-based parameterization of the TB risk signature:*** The tuberculosis risk signature
1038 was adapted to the Illumina microarray platform by first collapsing each splice junction rule pair
1039 into a gene pair. All Illumina probes corresponding to the two genes in the rule were identified,
1040 and all possible probe pairs using these two sets of probes were constructed. Rules were then
1041 trained (using linear SVMs) on for each probe pair using the tuberculosis disease patient and
1042 latently infected control data from the UK Training Set from Berry, *et al.*¹¹ The resulting
1043 Illumina microarray-based signature was locked down and remained unmodified for making
1044 predictions on the other microarray cohorts from the same study¹¹ and other studies.^{9,10,12,13}

1045 **Supplementary Tables**

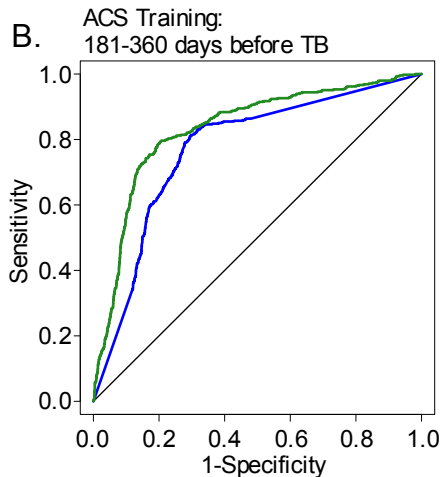
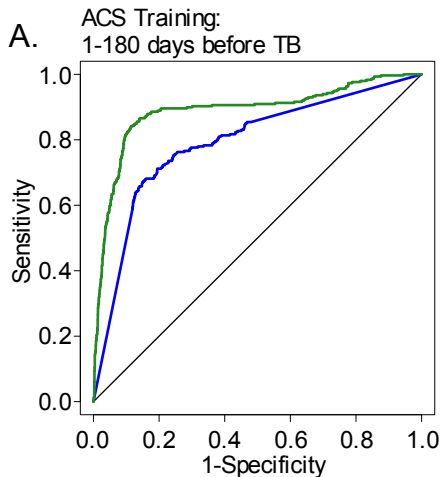
- 1046 1. (SupTab1_ACS_Progressors): Tuberculosis diagnosis in ACS progressors
- 1047 2. (SupTab2_ACS_TimeToDiagnosis): Original and per protocol (PP) sampling time points
1048 with respect to tuberculosis diagnosis for progressors from the ACS cohort.
- 1049 3. (SupTab3_GC6Progressors): Tuberculosis diagnosis in GC6-74 progressors.
- 1050 4. (SupTab4_GC6TimeToDiagnosis): Sampling time points with respect to tuberculosis
1051 diagnosis for progressors from the GC6-74 cohort.
- 1052 5. (SupTab5_SampleNumbers): Number of progressors and controls and corresponding
1053 samples for identifying and validating signatures of risk using the ACS training, ACS
1054 test, and GC6-validation cohorts.
- 1055 6. (SupTab6_RNASeqMetadata): Meta-data for RNA-Seq samples from ACS training and
1056 test cohorts.
- 1057 7. (SupTab7_ReferenceJunctions): Reference junctions used for internal normalization of
1058 the signature.
- 1059 8. (SupTab8_SignatureJunctions): Splice junctions, genes and chromosomal locations that
1060 comprise the signature
- 1061 9. (SupTab9_SignatureJunctionPairs): RNA-Seq junction pairs and SVM parameters that
1062 comprise the signature.
- 1063 10. (SupTab10_SignatureJunctionData): Raw RNA-Seq splice junction count data for all
1064 junctions in the signature.
- 1065 11. (SupTab11_TBRRiskComputationSheet): Worksheet that calculates risk scores upon input
1066 of raw splice junction count data (from SupTab10_SignatureJunctionData)
- 1067 12. (SupTab12_RefJunctionPrimers): Primers corresponding to the reference junctions that
1068 were used in the signature
- 1069 13. (SupTab13_Signatureprimers): Primers corresponding to the splice junctions contained in
1070 the signature.
- 1071 14. (SupTab14_SignaturePrimerPairs): Signature primer pairs used for blind prediction of
1072 ACS Test set and GC6-74 Validation set.
- 1073 15. (SupTab15_BenchmarkRF): ROC AUCs for the ACS training and test benchmark
1074 analysis against Random Forest

- 1075 16. (SupTab16_SignatureScores): Per-sample TB risk score and correct classification for
1076 ACS training fit (RNA-Seq and qRT-PCR), ACS test predict (RNA-Seq and qRT-PCR),
1077 GC6 predict (qRT-PCR).
- 1078 17. (SupTab17_MicroarrayAnalysis): Details and prediction statistics for the TB risk
1079 signature-based meta-analyses of published microarray datasets.
- 1080 18. (SupTab18_TBAssociations): Genes in the TB risk signature and known associations
1081 with TB and/or relevant functions.
- 1082 19. (SupTab19_SignatureGeneMatrix): Genes in the tuberculosis risk signature and
1083 diagnostic signatures (Berry *et al.*, 2010; Kaforou *et al.*, 2013; Anderson *et al.*, 2014).
- 1084 20. (SupTab20_SignatureSimilarity): Intersections between the tuberculosis risk signature
1085 and diagnostic signatures (Berry *et al.*, 2010; Kaforou *et al.*, 2013; Anderson *et al.*,
1086 2014).
- 1087 21. (SupTab21_ModuleEnrichments): Obermoser *et al.* 2013⁴⁴ module enrichments for genes
1088 that are common to TB risk and diagnostic signatures or specific to diagnostic signatures.
- 1089 22. (SupTab22_CellTypes): Heatmap and significance counts for cell-type specific
1090 expression patterns of the PSVM signature transcripts based on data from Bloom *et al.*,
1091 2013.

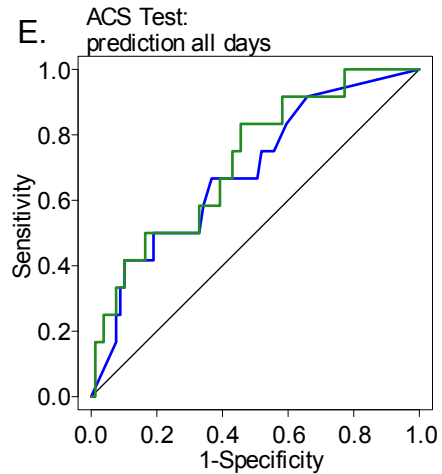
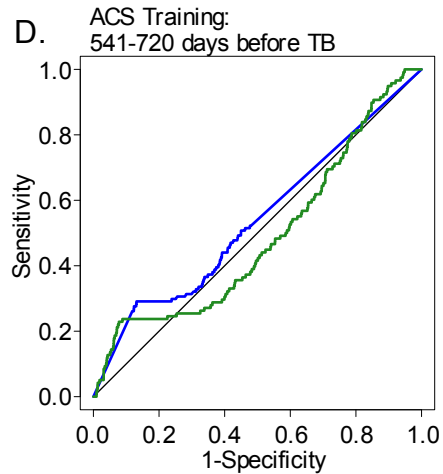
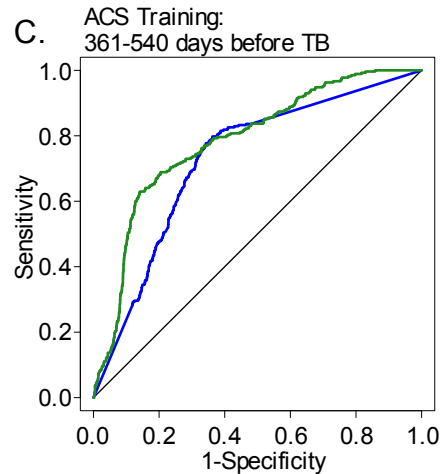
Supplementary Figure 1



Supplementary Figure 2



— PSVM
— Random Forest



Supplementary Figure 3

