

Published in final edited form as:

Nat Genet. 2018 June ; 50(6): 849–856. doi:10.1038/s41588-018-0117-9.

Frequent transmission of the *Mycobacterium tuberculosis* Beijing lineage and positive selection for EsxW Beijing variant in Vietnam

Kathryn E Holt^{1,*}, Paul McAdam^{#1}, Phan Vuong Khac Thai^{#2}, Nguyen Thuy Thuong Thuong³, Dang Thi Minh Ha², Nguyen Ngoc Lan², Nguyen Huu Lan², Nguyen Thi Quynh Nhu³, Hoang Thanh Hai³, Vu Thi Ngoc Ha³, Guy Thwaites^{3,4}, David J Edwards¹, Artika P Nath^{5,6}, Kym Pham⁷, David B Ascher¹, Jeremy Farrar^{3,4}, Chiea Chuen Khor^{8,9}, Yik Ying Teo^{10,11}, Michael Inouye^{6,7,12}, Maxine Caws^{#13,14}, and Sarah J Dunstan^{#15,*}

¹Department of Biochemistry and Molecular Biology, Bio 21 Molecular Science and Biotechnology Institute, University of Melbourne, Parkville, Victoria 3010, Australia

²Pham Ngoc Thach Hospital for Tuberculosis and Lung Disease, Ho Chi Minh City, District 5, Viet Nam

³Oxford University Clinical Research Unit, Ho Chi Minh City, District 5, Viet Nam

⁴Centre for Tropical Medicine, Nuffield Department of Clinical Medicine, Oxford University, Oxford, UK

⁵Department of Microbiology and Immunology, University of Melbourne, Parkville, Victoria 3010, Australia

⁶Systems Genomics Lab, Baker Heart and Diabetes Institute, Melbourne 3004, Victoria, Australia

⁷Department of Clinical Pathology, University of Melbourne, Parkville, Victoria 3010, Australia

⁸Genome Institute of Singapore, Singapore

⁹Singapore Eye Research Institute, Singapore

¹⁰Department of Statistics and Applied Probability, National University of Singapore, Singapore

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:http://www.nature.com/authors/editorial_policies/license.html#terms

* Corresponding authors: Correspondence should be addressed to KEH (kholt@unimelb.edu.au) and SJD (sarah.dunstan@unimelb.edu.au).

Author Contributions

SJD, KEH, MC, MI, YYT, CCK, are the study principal investigators who conceived and obtained funding for the project. SJD provided overall project co-ordination; MI organized and supervised the DNA sequencing and KEH devised the overall analysis plan and wrote the first draft of the manuscript along with PM. MC and SJD established the TB cohort for this genetics study by working with PVKT, DTMH, NNL, NHL, NTQN, NTTTT, GT and JF to coordinate the collection of clinical samples and phenotypes. KP performed DNA quality checks and genome sequencing on all Vietnamese samples, while VTNH performed Sanger sequencing on selected samples. DBA performed protein structure analyses, and HTH and NTTTT performed the macrophage growth and infection experiments of EsxW variants. KEH, PM, MI, DJE, AN analyzed the data. All authors critically reviewed manuscript revisions and contributed intellectual input to the final submission.

Competing Financial Interests

The authors declare no competing financial interests.

Life Sciences Reporting Summary

Further information is available in the Life Sciences Reporting Summary

¹¹Saw Swee Hock School of Public Health, National University of Singapore, Singapore

¹²Department of Public Health and Primary Care, University of Cambridge, Strangeways Research Laboratories, Cambridge CB1 8RN, United Kingdom

¹³Department of Clinical Sciences, Liverpool School of Tropical Medicine, Pembroke Place, Liverpool, L3 5QA, United Kingdom

¹⁴Birat-Nepal Medical Trust, 257 Lazimpat, Kathmandu, Nepal

¹⁵Peter Doherty Institute for Infection and Immunity, University of Melbourne, Parkville, Victoria 3010, Australia

These authors contributed equally to this work.

Abstract

To examine transmission dynamics of *Mtb* isolated from TB patients in Ho Chi Minh City, Vietnam we sequenced whole genomes of 1,635 isolates and compared these with 3,144 isolates from elsewhere. The data reveal an underlying burden of disease caused by endemic *Mtb* Lineage 1 associated with activation of long-term latent infection, and a three-fold higher burden associated with more recently introduced Beijing lineage and Lineage 4 *Mtb* strains. We find that Beijing lineage *Mtb* is frequently transferred between Vietnam and other countries, and detect higher levels of transmission of Beijing lineage strains within this host population than endemic Lineage 1 *Mtb*. Screening for parallel evolution of Beijing lineage-associated SNPs in other *Mtb* lineages as a signal of positive selection, we identify a mutation in the ESX-5 type VII secreted protein EsxW, which could potentially contribute to the enhanced transmission of Beijing lineage *Mtb* in Vietnamese and other host populations.

Introduction

Tuberculosis (TB) is a leading cause of death from infectious disease and the global burden is now higher than at any point in history^{1,2}. Despite coordinated efforts to control TB transmission, the factors contributing to its successful spread remain poorly understood. Vietnam is identified as one of 30 high burden countries for TB and MDR-TB with an incidence of 137 TB cases per 100,000 individuals in 2015². Recent phylogenomic analyses of the causative agent *Mycobacterium tuberculosis* (*Mtb*) in other high-prevalence regions have provided insights into the complex processes underlying TB transmission^{3–5}.

Results

Genetic diversity and drug resistance

To characterize the diversity of *Mtb* circulating in Ho Chi Minh City (HCMC), we sequenced the genomes of 1,635 isolates (Supplementary Table 1) obtained from 2,091 HIV uninfected, smear positive adults (> 18 years) commencing anti-TB therapy at district TB units (DTUs) in eight districts of HCMC between December 2008 and July 2011 (see Methods). This identified 73,718 SNPs, which we used to reconstruct a maximum likelihood phylogeny (Fig. 1a) and to assign lineages⁶. The majority of isolates (n=957, 59%) belonged to lineage 2.2.1, a subgroup of the Beijing lineage (2.2). Lineage 1 (Indo-Oceanic lineage;

n=388, 23.7%) and Lineage 4 (Euro-American lineage; n=192, 11.7%) were also common. A single isolate belonged to Lineage 3 (East African-Indian lineage) and was excluded from further analysis. The distribution of lineages did not change during the 2.5-year period of study (Fig. 1b), and was in agreement with previous genotyping studies in urban areas of Vietnam (50% 2.2/Beijing lineage and ~20% Lineage 1.1/EIA in Hanoi and HCMC, 1998-2009)⁷⁻¹¹. Known antimicrobial resistance mutations were detected in all lineages but were more frequent in Beijing sublineage 2.2.1 (Table 1), consistent with earlier reports from Vietnam^{7-9,11}. In particular, Beijing sublineage 2.2.1 was enriched for mutations associated with resistance to streptomycin (OR 4.6 [95% CI 3.6-6.0], $p=1\times 10^{-15}$), isoniazid (OR 1.7 [1.3-2.1], $p=3\times 10^{-5}$), rifampicin (OR 5.4 [2.5-13.2], $p=1\times 10^{-7}$) and ethambutol (OR 5.6 [2.6-13.7], $p=1\times 10^{-7}$), using Fisher's exact test to compare to all non-2.2.1 *Mtb* isolates.

***Mtb* lineage and host demographics**

Whilst the majority of TB patients were male (74%, typical for TB studies in Vietnam and elsewhere^{8-10,12}), the Beijing sublineage was significantly associated with TB in females (OR 1.28 [95% CI 1.01-1.62], $p=0.043$ using Fisher's exact test; Table 1), consistent with prior observations in Vietnam⁸ and Nepal¹³. Beijing sublineage 2.2.1 was also significantly associated with younger people: its frequency declined with age, from 74% of cases in <20 year olds to 50% in 60 year olds ($p=0.0023$ Fisher's exact test, $p=0.0024$ linear trend test; Fig. 1c). In contrast, Lineage 1 was significantly associated with males (25% of male cases vs 19% of females, $p=0.017$) and increased with age regardless of gender, from 12% in <20 year olds to 35% in 60 year olds ($p=0.0007$ Fisher's exact test, $p=0.0014$ linear trend test; Fig. 1c). These data confirm that Beijing sublineage 2.2.1 is capable of causing active disease in a wider demographic range of the Vietnamese host population, particularly among females and younger people, than the endemic Lineage 17-10, which is associated with the more typical profile of TB susceptibility that is skewed towards males [M:F prevalence ratio for smear-positive TB was recently estimated at 2.5 (95% CI 2.07-3.04), based on 40 surveys in 22 countries¹²] and older people^{2,14}.

***Mtb* lineages display distinct local transmission dynamics**

We hypothesised that Beijing lineage or sublineage 2.2.1 was more transmissible than Lineage 1, and/or more capable of causing active disease in infected hosts, in the local Vietnamese Kinh host population. To investigate this, we used the whole genome phylogeny to compare diversity metrics for each lineage (Fig. 2). Terminal branch lengths, which represent an upper bound of evolutionary time since transmission for each *Mtb* case, were significantly shorter for Beijing sublineage 2.2.1 *Mtb* isolates (median 8 SNPs) than for non-Beijing lineage isolates (Lineage 1: median 53 SNPs, $p<1\times 10^{-15}$ using Kolmogorov-Smirnov test; Lineage 2.1: 30 SNPs, $p<1\times 10^{-6}$; Lineage 4: 17 SNPs, $p<1\times 10^{-9}$), and slightly shorter than Beijing sublineage 2.2.2 isolates (9 SNPs, $p=0.02$) (Fig. 2a). The distribution of mean node-to-tip distances for all internal nodes was skewed significantly lower within the Beijing sublineage 2.2.1 compared to the rest of the tree (median 16 SNPs compared to 62, 57, 39 and 60 SNPs for Lineages 1, 2.1, 2.2.2 and 4, respectively; $p<0.0015$ in all cases).

To better understand the differences in transmission dynamics, we explored the distribution of potential transmission clusters using a range of maximum pairwise SNP distance

thresholds to define a cluster¹⁵ (Fig. 2b-c). Using the smallest cut-off of five SNPs (transmission age of <5 years) there were n=109 clusters, of which 76 (70%) belonged to Beijing sublineage 2.2.1; these had a mean size of 2.4 strains per cluster, compared to mean 2.1 for other clusters (Fig. 2b-c). Using cut-offs of 10 and 20 SNPs there were n=164 and n=220 clusters, respectively; of these, 118 (72%) and 156 (71%), respectively, were Beijing 2.2.1 and these showed significantly larger cluster sizes (means, 2.5 and 3.2 strains per cluster) than those of other lineages (means, 2.2 and 2.4). Notably, the proportion of cases that belonged to transmission clusters was significantly lower amongst Lineage 1 cases (7.7% at the 10 SNP threshold) compared to Lineage 4 cases (20.3%, $p=2\times 10^{-5}$), which in turn was significantly lower than among Beijing lineage cases (31.5%, $p=3\times 10^{-3}$; see Fig. 2b). Therefore sublineage 2.2.1 infections were more likely to result from recent detectable transmission within the local population, despite our study having low power to detect recent transmission due to sampling only a small fraction of all incident cases in HCMC (~30% of those in the study districts, <10% across HCMC, and restricted to new cases of smear-positive, culture-positive HIV-negative TB only; see Methods). Household data was not available; however pairs of TB cases whose infecting *Mtb* strains were separated by 10 or 20 SNPs were significantly more likely to be diagnosed in the same DTU than more distantly related pairs of the same lineage (27% or 21%, vs 15% amongst strains separated by >20 SNPs; $p<1\times 10^{-12}$, see Supplementary Figure 1). This phenomenon was observed for both Beijing 2.2.1 and non-Beijing clusters, but the effect was significantly less for the Beijing lineage (of strain pairs separated by <20 SNPs, 21% of Beijing pairs vs 29% of non-Beijing pairs were diagnosed in the same DTUs; $p=0.036$, see Supplementary Figure 1), suggesting they may be more readily transmitted across greater geographical distances within HCMC than other *Mtb*. This could potentially be associated with the higher frequency of Beijing lineage TB amongst younger adults (Fig. 1c), who may be more likely to travel regularly between districts, promoting onward transmission; however the age distribution amongst transmission clusters did not differ by lineage. Taken together, the phylogenomic data reveal significant differences in transmission dynamics between the various *Mtb* lineages circulating in HCMC, and suggest that newly diagnosed cases of Lineage 1 *Mtb* in this population often result from activation of longer-term latent infections with many private SNPs and no evidence of recent transmission, whilst new cases of Beijing sublineage 2.2.1 *Mtb* often result from more recent transmission and shorter time to develop active disease.

Geographical relationships of Vietnamese *Mtb* strains

It has been suggested that the Beijing lineage is slowly displacing the resident Lineage 1 strains in Vietnam, following the introduction of the Beijing strain into urban areas and subsequent spread to rural areas where Lineage 1 still dominates^{8,10}. Our data are consistent with this, showing a higher frequency of Beijing lineage (65%) amongst *Mtb* infections in HCMC in 2008-2011 compared to the frequencies reported in the city a decade earlier⁷ (53% in 1998), or in rural areas outside the city⁸ (32-37% in 2003-2005). We therefore hypothesized that Beijing 2.2.1 isolates from HCMC may represent a locally established epidemic subclade of the Beijing lineage, similar to that previously described in Russia³. To investigate this we combined our HCMC *Mtb* genome data with 3,146 publicly available *Mtb* whole genome sequences from Russia³, Malawi^{4,5}, Argentina¹⁶, and

China¹⁷; and globally dispersed Lineage 1 and 2 genomes^{18,19,20} (Supplementary Table 2); then inferred phylogenies for each lineage (Fig. 3). HCMC Lineage 1 strains were quite distinct from those in other locations (Fig. 3a), with little evidence of transfer between Vietnam and other regions. The vast majority (n=319, 82%) belonged to a localised subclade (1.1.1.16) that included only seven previously sequenced strains, all of which were from Vietnam. A further 46 (12%) of the HCMC strains belonged to a related sister clade (1.1.1) which also included strains from neighbouring Cambodia (n=1) and Thailand (n=2); the remainder (n=23, 6%) belonged to subclade 1.2.1, in which the HCMC strains were intermingled with others from the Philippines (n=9) and China (n=1). These data suggest that Lineage 1 associated TB in HCMC results mainly from a local endemic *Mtb* population. In contrast, Beijing 2.2.1 isolates from HCMC formed several distinct clusters that each shared a recent common ancestor with isolates from outside Vietnam (Fig. 3b). Notably, isolates from Russia, Malawi, China and numerous other countries were interspersed throughout the HCMC Beijing 2.2.1 population (Fig 3b), suggesting multiple, frequent transfers between host populations in HCMC and other geographic regions. HCMC Lineage 4 isolates were drawn from eight of the ten recognised sublineages⁶ (including those identified as specialist, generalist and intermediate in their geographic range²¹) and were interspersed with isolates from other geographical locations, consistent with multiple imports into HCMC (Fig. 3c). In further support of these observations, stochastic mapping of locations onto the phylogenies predicted dozens of strain transfer events between Vietnam and other locations for Lineages 2 and 4, but not for Lineage 1 (Fig. 3d), strongly supporting that *Mtb* sublineage 1.1.1.1 is endemic in Vietnam.

Beijing lineage-defining SNPs under positive selection

The population structure (Figs 1-2) provides evidence that Beijing lineage strains are more transmissible within this HIV-negative HCMC population than are other *Mtb* lineages. Genomic evidence for enhanced transmission of the Beijing lineage has been documented in Russia (associated with antimicrobial resistance)³ and Malawi (independent of antimicrobial resistance)⁴. While antimicrobial resistance was common amongst HCMC Beijing lineage isolates, the majority of transmission clusters (defined by 10 SNPs) comprised groups of isolates that did not share any known resistance mutations that could account for their transmission success (Supplementary Figure 2). This is consistent with previous reports that the Beijing lineage is highly transmissible and more likely to progress to active disease in various host populations and is also more virulent and less pro-inflammatory in various cellular assays, independent of antimicrobial resistance^{22–25}. We therefore aimed to interrogate the *Mtb* genome data to identify mutations that may contribute to the success of the Beijing lineage (2.2). Evolutionary convergence has previously been used as a signal of positive selection to identify mutations associated with antimicrobial resistance in *Mtb*^{26,27}. We reasoned that advantageous polymorphisms contributing to the enhanced transmissibility of Lineage 2.2 should be fixed in this lineage, and should also be under positive selection that is detectable as convergent or parallel evolution at the same variant sites in other lineages. We identified a total of 424 homoplastic nonsynonymous SNPs (nsSNPs) across the HCMC phylogeny. The most frequent of these occurred in genes in which convergent evolution has previously been associated with antimicrobial resistance including *gidB*, *embB*, *gyrA*, *rpoB*, and *inhA26*, which together accounted for 12.4% of all homoplastic

nsSNPs. The distribution of common homoplastic nsSNPs in these genes is shown in Supplementary Figure 3; in particular, rifampicin resistance-associated mutations in *rpoB* and ethambutol-associated mutations in *embB* arose independently many times in Beijing sublineage 2.2.1 (n=33/35 *rpoB*-450 mutations; n=16/17 *embB*-306 mutations; n=8/10 *embB*-406 mutations). The homoplastic nsSNPs included three that arose on the branch defining Lineage 2.2 and also elsewhere in the HCMC tree (Table 2, Supplementary Figure 4). One was a mutation in *esxW* (Rv3620c) codon 2 (EsxW-Thr2Ala), which arose on nine other branches (six times in Lineage 4, three times in Lineage 1; see Supplementary Figure 4) and showed evidence of onward transmission on 4/9 occasions. Comparison to the global tree detected the same *esxW* mutation on a further ten Lineage 4 branches in Malawi and Russia, with onward transmission detected on 6/10 occasions. The other two mutations were in *Rv3081* (conserved hypothetical protein) and *gidB* (mutations in which are often associated with streptomycin resistance) and arose less frequently (Table 2). In contrast, homoplastic nsSNPs on the branches defining Lineages 1 or 4 were each detected on only 1-2 other branches in the HCMC tree and no additional branches of the global tree (Supplementary Table 3). No homoplastic SNPs were associated with sublineage 2.2.1, and although synonymous or intergenic SNPs can have functional consequences, we found no such homoplasies associated with Beijing or other lineages.

EsxW mutation

EsxW is included in multiple *Mtb* vaccines currently under development (including H65, ID83, ID93)^{28–31} due to its demonstrated immunogenicity in mice, safety and immunogenicity in non-human primates³¹, demonstrated T-cell targeting in humans³², and epitopes predicted to bind a wide range of human HLA-DRB1 alleles^{28,33}. Hence we considered whether the EsxW-Thr2Ala mutation could affect epitope binding. However, residue two lies in the N-terminal loop of EsxW (see Fig. 4e), outside the experimentally demonstrated epitope region (residues 24-34)³⁴. *In silico* epitope binding analysis on the Beijing and non-Beijing EsxW protein sequences using the Immune Epitope Database Analysis Resource³⁵ predicted HLA binding to the first 9-10 residues of EsxW, but no differences in binding affinities for the wildtype and Thr2Ala mutant alleles. This is consistent with experimental data showing that immunization with ID93 (a recombinant fusion protein containing H37Rv (Lineage 4) wildtype alleles of EsxW/V) is protective against Beijing (2Ala) as well as non-Beijing (2Thr) strains^{30,31}.

Next we considered whether the EsxW-Thr2Ala mutation could impact on gene expression or protein structure and function. *esxW* is one of 23 *esx* genes in the *Mtb* genome, including 11 clustered pairs of *esx* genes whose products form heterodimers that are each secreted by one of five type VII secretion systems (T7SS; ESX-1 to ESX-5). The most-studied of these pairs is *esxB* (CFP10)/*esxA* (ESAT-6), secreted by ESX-1 and encoded in the RD1 locus which also encodes the ESX-1 system. EsxW and its heterodimerization partner EsxV are encoded by adjacent genes in the RD8 locus and secreted by ESX-5. The ESX-5 system is the most recently evolved T7SS in *Mtb* and is present only in the slow-growing *Mycobacteria*^{36,37} (including the *Mtb* complex, *M. leprae*, *M. ulcerans* and *M. marinum*). ESX-5 is unique amongst the *Mtb* T7SS in that (a) it secretes most of the PE/PPE proteins³⁸, which are also unique to slow-growing *Mycobacteria*, comprise a substantial

amount of protein coding capacity (~10%) in the *Mtb* genome and play various roles in virulence^{39,40}; and (b) it is associated with five pairs of Esx proteins, resulting from duplication and expansion of *esxM* and *esxN* (which are encoded within the locus encoding the ESX-5 T7SS machinery) to create four paralogous copies elsewhere in the *Mtb* genome⁴¹ (Fig. 4). Each pair includes a member of the QILSS family (EsxM paralog, including EsxW) and a member of the Mtb9.9 family (EsxN paralog, including EsxV), which partner to form heterodimers. The paralogous proteins differ from one another by just a few amino acids (Fig. 4b), and the reason for this apparent redundancy is not clear. One pair, EsxJ/EsxI, has been shown to play a role in substrate selection for ESX-5 secretion⁴², and it is hypothesized that the other paralogs including EsxW/V play a similar role⁴¹. Notably, two of the loci encoding QILSS paralogs (EsxP/O, RD5; and EsxW/V, RD8) are missing from the *M. bovis* BCG vaccine strain.

The upstream sequences of *esxW* and its homologs differ substantially (Fig. 4a). This suggests their expression is subject to different regulatory controls, which could provide a reason for their expansion in the *Mtb* genome despite the lack of differentiation at the protein level⁴³ (Fig. 4b); indeed, it has been shown that the different QILSS family members are all expressed at different levels by *Mtb* strain H37Rv during growth in broth, sputum and macrophages⁴⁴. This uniqueness also allows unambiguous read mapping and confident SNP calling at the N-terminal region of EsxW; however we also used PCR and capillary sequencing to confirm the EsxW-2Ala allele in all non-Lineage 2 strains in which the SNP was identified from Illumina reads (see Methods). We sought to investigate the expression of QILSS family proteins, and the potential impact of the EsxW-Thr2Ala mutation on gene expression, during growth in macrophages. For these experiments we selected four non-Beijing *Mtb* isolates harbouring the EsxW-2Ala allele, and the closest genetic relative of each with the wildtype allele EsxW-2Thr, and used RNAseq to measure genome-wide *Mtb* expression levels (see Methods). The results showed that the ESX-5 system was highly expressed in all isolates following 24h growth in macrophages (mean 0.25% of total *Mtb* RNAseq reads); however the QILSS paralogs were expressed at different levels relative to ESX-5 (Fig. 4c). Compared to *esxM*, which is encoded within the ESX-5 locus itself and was expressed at the highest levels in all isolates, *esxK* and *esxP* were expressed at intermediate levels (mean 42% of *esxM* level) and *esxJ* and *esxW* were expressed at low levels (mean 2.2% of *esxM* level). There were no significant differences between the EsxW-2Ala mutant vs wildtype in terms of *Mtb* growth in macrophages (measured at 4, 7 or 11 days post-infection), or in *esxW* (Fig. 4d), ESX-5 or global gene expression measured at 24h post-infection. It has been reported that *esxW* was significantly expressed in the lungs of *Mtb* aerosol-infected mice⁴⁵; hence future experiments with wildtype and mutant EsxW in this animal system could potentially help to unravel its functional effects.

Given the proposed role for the EsxW/V heterodimer in substrate selection for ESX-5 secretion⁴¹, we investigated whether the EsxW-Thr2Ala mutation could affect protein structure and function. In *M. canetti*, the reference *Mtb* genome H37Rv, and the majority of non-Beijing *Mtb* isolates, EsxW carries the polar threonine (codon ACC) at residue 2, while the other QILSS proteins in the *Mtb* complex and other slow growing Mycobacteria carry the hydrophobic alanine (GCC) at this position (Fig. 4b). In the *Mtb* Beijing lineage, EsxW

residue 2 is converted to the more typical residue alanine (GCC), making it identical at the protein level to EsxJ, which has been shown to be involved in substrate selection for ESX-5 secretion⁴². We analysed the quantitative effects of the Thr2Ala mutation on the stability and affinity of the EsxW protomer and the EsxW/V heterodimer, using computational modeling and direct biophysical experiments (see Methods and Fig. 4e). This indicated that the 2Ala mutation was likely to lead to a mild increase in the affinity and stability of the heterodimer complex (mean ΔG 0.28±0.06 Kcal/mol), and this was supported by biophysical measurements that showed the mutant bound slightly stronger to EsxV than the wildtype (K_D^{wt} =0.6 μ M; K_D^{T2A} =0.4 μ M; p =0.05 using two-tailed t-test; Fig. 4f). We hypothesise that, if EsxW/V does indeed play a role in substrate selection for ESX-5 secretion, then increased stability of the heterodimer and/or the hydrophobicity of EsxW-2Ala could potentially affect the efficiency of secretion of certain PE/PPE proteins. This could thereby have downstream impacts on one or more of the known functions of ESX-5 including inflammasome activation, IL-1 β secretion or escape from macrophages⁴⁶, any of which could potentially work to promote transmission between hosts.

Discussion

The shorter terminal branch lengths and node-to-tip distances for local Beijing lineage *Mtb* (Fig 2) could be explained by (i) a slower mutation rate in the Beijing lineage, resulting in slower accumulation of SNPs over time; (ii) sampling bias, whereby new cases of active TB arising in the study population were more likely to be detected and included in the study if they were caused by Beijing lineage strains; or (iii) strain-specific transmission dynamics in the study population, whereby the average time taken to progress to active disease is shorter for Beijing lineage than for other *Mtb* strains. The mutation rate for *Mtb* has been estimated at ~0.5 SNPs per year for Lineage 4, and ~2x faster for Beijing strains⁴⁷ (Lineage 1 has not been specifically measured but is assumed to be similar to Lineage 4, which is considered typical). Hence if all *Mtb* lineages were subject to the same transmission dynamics within the study area, we would expect to see longer terminal branch lengths for Beijing lineage isolates, whereas here we observe the opposite (significantly shorter branch lengths). Patients were recruited into the study following positive diagnosis at eight central DTUs in HCMC (map in Supplementary Figure 5). The identity of the infecting *Mtb* lineages was not known at the time of recruitment, and we are unaware of any factors that would bias the inclusion of Beijing lineage cases over others presenting to the clinics for treatment; hence we conclude the frequency of infections with Beijing strains, and their shorter branch lengths, reflect differences in the transmission dynamics of these strains within the study population.

The lack of transfer of *Mtb* sublineage 1.1.1.1 between Vietnam and other geographical locations (Fig 3) may be associated with adaptation to the local host population with which it has co-evolved for centuries, similar to the ‘host-specialist’ clades recently identified within Lineage 4²¹. In contrast, whilst the direction of transfer of Lineage 2.2.1 strains between HCMC and other geographical regions cannot be determined from our data, the frequency of transfer events and the scale of diversity amongst the HCMC strains (Fig 3) makes it unlikely that the rise of Beijing sublineage 2.2.1 in HCMC represents clonal spread of a locally established subclade. Regardless of direction, the frequency of transfer between

Vietnam and diverse geographically dispersed populations supports previous contentions that the Beijing lineage is a host-generalist, capable of moving between ethnically diverse host populations²⁴.

While the mechanism remains to be elucidated, our results provide evidence that the Beijing lineage carries a variant of *esxW* that is under positive selection in natural *Mtb* populations. This is consistent with the idea that the protein is important for host interactions, potentially through substrate selection for ESX-5 secretion under certain conditions. Immunizing against *EsxW* has already been shown to be protective against infection^{29–32}, and positive selection for this protein suggests vaccines including it are likely to remain effective in the long term.

Taken together, our data show that the burden of TB in HCMC comprises (i) an underlying burden of disease caused by endemic Lineage 1 *Mtb* strains (24% of all TB cases), which disproportionately affects men and older people and is associated more with activation of long-term latent infection than short-term transmission clusters; and (ii) an additional disease burden caused by more recently introduced Lineage 2 and 4 *Mtb* strains (76% of all TB cases). In particular, Beijing sublineage 2.2.1 was associated with a wider demographic host range, infecting women and young people significantly more frequently than other lineages, and was associated with shorter time to active disease and frequent onward local transmission. One third of all Beijing strains were involved in transmission clusters (< 10 SNPs), and these were associated with wider geographic dispersal within HCMC. Notably, 75% of TB cases associated with transmission clusters involved the Beijing lineage, accounting for 20% of all cases included in the genomic study.

Importantly, these data show that not all *Mtb* are equal: genetically diverse strains display distinct transmission dynamics even within a single localized host population, suggesting that more detailed understanding of lineage-specific variation in *Mtb* could be informative to tailor local TB control in HCMC and other settings. For example, TB contact tracing is commonly used in low-incidence high-income countries, but results vary in high-incidence low- and middle-income countries and there is a need to prioritize resources towards cases where contact tracing is most likely to yield results, which includes considering the likelihood of transmission and progression to active disease within the time frame of a contact tracing program (1-2 years)^{48,49}. In HCMC it may thus be advantageous to direct contact tracing resources towards Beijing lineage cases, as they pose the greatest risk of onward transmission resulting in new active TB cases.

Online Methods

Bacterial isolates used in this study

Between December 2008 and July 2011, 2,091 individuals of the Vietnamese Kinh ethnic group attending the outpatient department of Pham Ngoc Thach Hospital or from 8 District Tuberculosis Units (District 1, 4, 5, 6, 8, Tan Binh, Binh Thanh and Phu Nhuan) in HCMC were recruited into a clinical study investigating predictors of failure and relapse in isoniazid resistant TB 50. The 8 TB units were chosen for inclusion from amongst the 24 servicing HCMC as they are centrally located and close to the TB reference laboratory, which was

essential for the logistics of sample collection and processing (see map in Supplementary Figure 5). Inclusion criteria were: (1) 18 years or older, (2) negative HIV test, (3) provision of written informed consent, (4) smear positive pulmonary TB. Exclusion criteria were: (1) under 18 years of age, (2) HIV infected, (3) unable or unwilling to provide consent, (5) pregnancy, (6) prior history of TB antibiotic therapy, (7) will receive TB-DOTS (directly observed treatment, short-course) outside the study centres. Over the 2.5 year study period there were N=5036 new smear positive cases of TB (HIV positive and negative) at the district tuberculosis units, of which N=2091 were eligible for recruitment. Of these, N=1822 *Mtb* strains were isolated from the study participants. The annual incidence of pulmonary TB in HCMC is ~82,000, of which ~11,000 are in HCMC; therefore our sample represents ~6.6% of all cases in HCMC during the 2.5-year study period, and >30% of cases in the study districts.

Ethics

The study protocol was approved by the Institutional Research Board of Pham Ngoc Thach Hospital (the supervisory institution of the District TB Units in southern Vietnam), Ho Chi Minh City Health Services and the Oxford University Tropical Research Ethics Committee, UK (Oxtrec 030-07). Written informed consent was obtained from all patients.

DNA extraction and sequencing

Mtb isolates were subcultured on Lowenstein Jensen media and DNA extracted at the Oxford University Clinical Research Unit in HCMC using the cetyl trimethylammonium bromide (CTAB) extraction protocol as described previously 51. DNA was successfully obtained from N=1,728 isolates and shipped to the University of Melbourne for whole genome sequencing. DNA extracts were purified using AxyPrep™ Mag PCR Normalizer Protocol prior to library preparation. A total of N=1,655 DNA samples passed QC, were included for sequencing and subjected to library preparation using the Nextera XT protocol. Libraries were quantified using Quant-iT PicoGreen (dsDNA kit, Invitrogen), then normalised and pooled to 4 nM concentration. DNA underwent 150 bp paired end sequencing (Rapid mode v2) on the Illumina HiSeq 2500 platform (Illumina, San Diego). Sequence data was excluded for N=19 *Mtb* isolates that yielded less than the pre-established criteria of 10x mean read depth across the *Mtb* genome, as SNPs can not be reliably called below this depth. Sequence data was successfully generated for N=1,635 *Mtb* isolates from HCMC (representing 90% of those isolated from eligible patients in the cohort) with median three million reads per sample, providing median 99.2% coverage and 86x depth for each *Mtb* genome (Supplementary Table 1). To confirm the *esxW* codon 2 SNP, we performed PCR and capillary sequencing of the region flanking the SNP in all non-Beijing lineage strains in which the SNP was detected from Illumina data. The primers used are listed in Supplementary Table 4.

Data Availability

Mtb genome data generated from 1,635 isolates (Figure 1) has been deposited in NCBI BioProject [accession ID: PRJNA355614; see URLs section]. A total of 3,144 *Mtb* genomes were included in the analysis in Figure 3, comprising data from localized studies: 1,032 from Russia 3, 1,621 from Malawi 4,5, 248 Argentina 16, and 78 from China 17; as well as

106 globally dispersed Lineage 2 genomes 18 and 59 globally dispersed Lineage 1 genomes 19,20. Illumina *Mtb* genome sequences from various previously published studies were downloaded from the European Nucleotide Archive (see URLs section, individual accession numbers are given in (Supplementary Table 2). The H37Rv reference genome sequence (see URLs section) was used for all reference-driven analyses.

SNP analysis

Sequence reads were mapped to the H37Rv reference genome using the RedDog pipeline v0.5 (see URLs section). Briefly, Bowtie2 v2.2.3 was used for read alignment with the sensitive-local algorithm and the maximum insert length set to 2000 (via the -x parameter) 52 and variant sites (i.e. SNPs) were called using SAMTools v0.1.19 53. SNPs located in previously reported repetitive regions of the genome were excluded prior to phylogenetic analysis 54,55 (Supplementary Table 5); sites for which a definitive allele call could not be made in at least 99.5% of all isolate sequences were also excluded from the set of SNPs used for phylogenetic analysis. Two SNP alignments were compiled for analysis: one comprising the 1,635 HCMC isolates (total 73,718 SNPs), and one comprising all 4,779 isolates (including the HCMC isolates and the global collections downloaded from public data; total 133,492 SNPs).

In silico lineage and antimicrobial resistance typing

Mykrobe Predictor v0.3.6 was used to analyse raw Illumina reads generated from HCMC *Mtb* isolates and (a) assign each isolate to one of the seven *Mtb* lineages, and (b) detect known resistance associated polymorphisms 56 (summarized in Table 1, individual mutation calls are provided in Supplementary Table 1). All *Mtb* isolates were further assigned to sublineages by comparing SNPs identified using RedDog with those used in the haplotyping scheme defined by Coll *et al.* 6 (lineage assignments are in Supplementary Tables 1-2).

Phylogenomic analyses

ML phylogenetic trees were inferred using RAxML v7.7.2 57 for (a) all HCMC isolates (presented in Fig. 1); and (b) each of lineages 1, 2 and 4 using combined data from the HCMC isolates and available public data (presented in Fig. 3; see isolates list in Supplementary Tables 1-2). The trees presented are those with the highest likelihood from 5 replicate runs, constructed using the GTR model of nucleotide substitution and a Gamma model of rate heterogeneity to analyse a concatenated alignment of SNP alleles. An approximate ML tree containing all data (HCMC isolates and available public data) was inferred using FastTree v2.1.8 58. Ancestral sequence reconstruction was performed for the HCMC tree and combined tree using FastML v3.1 to infer the sequence alignment at each internal node of the ML phylogeny 59. Substitution events occurring on each branch of the tree were extracted by comparing the joint reconstruction sequences for the parent and child nodes; these data were used to identify homoplasic SNPs, and to identify lineage-specific polymorphisms as well as independent occurrences of those polymorphisms outside of the lineage of interest (data in Table 2). Terminal branch lengths reported are the number of substitutions (SNPs) mapped to each terminal branch (data in Fig. 2a) and were compared to Beijing 2.2.1 sublineage using 2-sample Kolmogorov-Smirnov tests. Metrics for genetic diversity and tree topology were calculated from the phylogenies using R. Node-to-tip

distances showed similar variances within groups (standard deviations of 27-75 SNPs) and all groups were compared to Beijing 2.2.1 sublineage using 2-sample Kolmogorov-Smirnov tests. Clusters were defined as subtrees for whom the maximum patristic distance between descendant tips fell below a specified threshold (data in Fig. 2b-c). Each cluster was checked to determine whether all members of the cluster shared any of the antimicrobial resistance mutations identified by Mykrobe Predictor; clusters in which no known antimicrobial resistance mutation was conserved in all members of the cluster are reported as not explained by antimicrobial resistance (data in Supplementary Figure 2).

Phylogeography analysis

Transmission between geographical regions was assessed separately for Lineage 1, 2 and 4 trees using an implementation of stochastic mapping on phylogenies (SIMMAP) implemented in the phytools v0.5 package for R 60,61. Region of origin was treated as a discrete trait and mapped to each tree using the ARD model (which allows each region-to-region transfer rate to vary independently) with 100 replicates. The results reported (Fig. 3d) are the median values for the number of transitions to Vietnam from any other region, summarized from 100 replicate mappings for each tree.

Esx sequence analysis

Esx protein sequences were extracted from the H37Rv reference genome using Artemis, aligned using Muscle, and subjected to phylogenetic inference using PhyML v3.0 (tree in Fig. 4). DNA sequences flanking the start codon of each *esx* gene were extracted from the H37Rv reference genome using Artemis and aligned and visualised using JalView v2.6.1 (Fig. 4).

Macrophage infections

Mtb infection of macrophages (THP1 human cell line, 88081201, Sigma-Aldrich) has been previously described 62. To study *Mtb* growth in macrophages, 2.5×10^5 THP1 cells were seeded in 24-well cell culture plates and infected with *Mtb* isolates in triplicate at a multiplicity of infection (MOI) of 1. At 4h post infection extracellular bacteria were removed by washing. At 0, 4, 7 and 11 days post-infection, intracellular bacteria were harvested and plated on Middlebrook 7H10 agar plates. Colony forming unit counts resulted after 3 to 4 weeks incubation at 37°C. To measure genome-wide *Mtb* expression levels, confluent monolayers of 1.5×10^7 THP1 in 75cm² cell culture flasks were infected with *Mtb* isolates at MOI 4. At 4h post infection extracellular bacteria were removed by washing. At 24h post-infection, intracellular bacteria were stabilized and collected in guanidine thiocyanate-based lysis buffer (4 M guanidine thiocyanate, 0.5% Na N-lauryl sarcosine, 25 mM sodium citrate, and 0.1 M β-mercaptoethanol). Total RNA from was then extracted using a previously published method 63.

RNAseq analysis

RNA samples were subjected to reverse transcription and the resulting cDNA was sequenced via Illumina HiSeq (100 bp paired end) at Macrogen. The resulting reads were first mapped to the H37Rv *Mtb* reference genome to retrieve bacterial sequences, resulting in 6-19

million bacterial sequence reads per sample. A strain-specific reference genome was prepared for each of the eight *Mtb* isolates, by polishing the H37Rv genome sequence with the Illumina reads previously obtained by sequencing whole genomic DNA extracted from that isolate (i.e. the data used for phylogenomic analyses) using Pilon (v1.22) with default settings. For each isolate, the RNAseq reads were mapped to the isolate's own reference genome sequence using Bowtie2 (v2.2.3) 52. Read alignments were visualised using the BamView function in v14.0.0 64, which was also used to filter reads to include only those with mapping score ≥ 30 for both forward and reverse reads, and to extract read counts for regions of interest. Expression values for the ESX-5 locus were calculated as the total reads mapping to the region spanning from *pe19* to *eccA5* (coordinates 2029880 to 2038532 in H37Rv). Due to the very close sequence homology in the protein-coding regions of the QILSS genes, it was not possible to assess their expression by counting total reads mapping to the length of the gene. Instead, we assessed relative expression of these genes by extracting the number of reads mapping to a unique 30 bp marker region spanning the start codon of each paralog (coordinates -21 to +9 relative to the first base pair in each protein-coding sequence; see Fig. 4 for uniqueness of these regions). To facilitate comparison of QILSS gene expression across strains with different RNAseq library sizes, we normalised the read counts for each isolate by the total number of reads mapping to the ESX-5 locus from *pe19* to *eccA5* (coordinates 2029880 to 2038532 in H37Rv). The ESX-5 locus was chosen for this purpose as all QILSS proteins are secreted by ESX-5 and thus depend on its expression (note the locus contains the QILSS gene *esxM*, which was the highest expressed of all QILSS genes in all isolates). The ratios of reads mapping to each QILSS gene vs ESX-5 locus in each strain i ($x_i/ESX-5_i$) is shown in Figure 4. To investigate differences in QILSS gene expression between wildtype (i) and mutant (j) isolate pairs (Fig. 4), we calculated the difference in these ratios relative to the wildtype level, as follows: $(x_i/ESX-5_i - x_j/ESX-5_j)/(x_i/ESX-5_i)$.

Computational protein structure analysis

Structure guided approaches can provide valuable insight into the molecular mechanism of mutations and their role in diseases 20,65–74, 75. To evaluate the structural effects of the Thr2Ala mutation in EsxW, models of EsxW and EsxV were generated using Modeller v9.1976 and MacroModel (vSchrodinger 2017, New York, NY), based on the experimental structures of homologous ESAT-6-like complexes (PDB IDs: 1WA8, 2KG7, 2VS0, 3OGI, 3ZBH, 4GZR, 4IOG and 4LWS; sequence identities ranging from 20-91%). The models were then minimized using the MMF94s forcefield in Sybyl-X 2.1.1 (Certara L.P., St Louis, MO), with the final structure having more than 95% of residues in the allowed region of a Ramachandran plot. The quality of the models was confirmed with Verify3D v1 (data not shown). Model structures were examined using Pymol v1.9. The model of the EsxW-EsxV complex was built using the homologous complexes to guide protein docking. The EsxW protomer and EsxW-EsxV heterodimer model were subjected to molecular dynamics to generate a series of representative conformations using Desmond (vSchrodinger 2017, New York, NY). Sodium and chloride ions were added to reach a final concentration of 150 mM Na⁺ and the system was solvated. The Amber and Charmm36 force fields were applied to the system and the default Desmond minimization and equilibration procedure was followed. Simulations were kept at constant pressure (1 atm) and temperature (335 K). The

structural consequences of the EsxW Thr2Ala variant were analysed to account for all the potential effects of the mutations 77. The effects of the variant upon the stability of EsxW was predicted using SDM v2 78, mCSM-Stability v1 79 and DUET v1 80. The effect of the mutations upon the affinity of EsxW to bind to EsxV, and the stability of the complex, were predicted using mCSM-PPI v1 79. The predictions were analysed and averaged across the conformational landscape. These computational approaches represent the wild-type structural and chemical environment of a residue as a graph based signature in order to determine the change upon mutation in Gibb's Free Energy of stability or binding. Reported results are based on analysis of EsxW (wildtype and Thr2Ala mutant) and EsxV amino acid sequences encoded in the H37Rv (lineage 4) reference genome; the only difference between these dimers and those encoded in lineage 1 or 2 is at EsxV residues 20 and 23. Structural modelling with the lineage 1/2 EsxV sequence yielded the same results as the analyses conducted with the lineage 4 background (i.e. stronger binding affinity observed with the EsxW Thr2Ala mutant).

Biophysical Measurements

EsxV and wildtype (2Thr) and mutant (2Ala) EsxW were cloned into a pET18 vector with a C-terminal hexahistidine tag (again using the sequences encoded in H37Rv (lineage 4) reference genome). A di-lysine motif was added prior to the hexahistidine tag of EsxW to facilitate labelling on the extended C-terminus. Protein was expressed, purified and refolded following the method of Mahmood and colleagues⁸¹. Solution microscale thermophoresis (MST) binding studies were performed to measure the binding affinity of EsxV for both the wildtype and mutant EsxW, using standard protocols on a Monolith NT.115 (Nanotemper Technologies)^{82–84}. The EsxW was chosen to be labelled due to the presence of a lysine in EsxV at the heterodimer binding interface. As there were no lysines present in EsxW, the dilysine motif was added to the C-terminal projection to facilitate labelling using a RED-NHS (Amine Reactive) Protein Labeling Kit (Nanotemper Technologies), which contains an NT-647 dye, as per the manufacturer's instructions. Labelled either wildtype or 2Ala mutant EsxW was mixed with EsxV in PBS with 0.05% Tween-20. Each replicate was undertaken using a 16-step 2-fold serial dilution series. The EsxW protein concentration was chosen such that the observed fluorescence was approximately 400 units at 70% LED power. The samples were loaded into standard capillaries and heated at 40% laser power (48 mW) for 30 s, followed by 5 s cooling. The data were normalized against the baseline obtained in the absence of any EsxV, and the maximal response obtained at the highest concentration of inhibitor. The dissociation constant K_D was obtained by plotting the normalized fluorescence (F_{norm}) against the logarithm of the concentrations of the dilution series and resulted in a sigmoidal binding curve that could be directly fitted with a nonlinear solution of the law of mass action (Figure 4). All experiments were performed with four replicates, and the normalized fluorescence thermophoresis curves were analyzed using GraphPad v6 (GraphPad, San Diego, CA, USA).

Statistical analyses

All statistical analyses were performed in R version 3.3.3 unless otherwise stated. Associations between lineages and patient demographics (age group, sex) were assessed using Fisher's exact test (two-sided in all cases) to assess $n=1,634$ HCMC *Mtb* genomes (i.e.

excluding the Lineage 3 isolate). A linear test for trend in proportions was also used to test for an increase in frequency of Lineage 1 within increasing age groups (also $n=1,634$). Tests for difference in the distributions of terminal branch lengths or node-to-tip distances between lineages (Fig 2) were conducted using the Kolmogorov-Smirnov test (two-sided in all cases, $n=1,634$ *Mtb* genomes). Differences in the proportion of pairs involved in transmission clusters between lineages, or between case pairs isolated at the same DTU within 10 SNP or 20 SNP clusters vs pairs separated by >20 SNPs (Supplementary Figure 1), were tested using a two-sided test for equal proportions ($n=1,634$ *Mtb* genomes). A two-tailed t-test was used to compare K_D values for mutant vs wildtype (Fig 4f), analysed using GraphPad Prism ($n=8$, i.e. 4 replicates each for wildtype and mutant).

URLs

NCBI BioProject accession, ID PRJNA355614, <http://www.ncbi.nlm.nih.gov/bioproject/355614>; The H37Rv reference genome sequence, accession NC_000962.3, <https://www.ncbi.nlm.nih.gov/genbank/>; RedDog pipeline v0.5, <https://github.com/katholt/RedDog>; European Nucleotide Archive, <http://www.ebi.ac.uk/ena>;

Data Availability

Mtb genome data was deposited in NCBI BioProject [ID: PRJNA355614; <http://www.ncbi.nlm.nih.gov/bioproject/355614>]; individual accession numbers for *Mtb* genomes analysed in this study are given in Supplementary Tables 1-2 (including data from previous studies).

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

We would like to thank the clinical staff who recruited patients into our study from the following District TB Units (DTU) in HCMC, Viet Nam; District 1, 4, 5, 6, 8, Tan Binh, Binh Thanh and Phu Nhuan DTUs; and also our colleagues from Pham Ngoc Thach Hospital for Tuberculosis and Lung Disease, HCMC Viet Nam. This work was supported by the National Health and Medical Research Council, Australia (Project grant #1056689 to SJD, Fellowship #1061409 to KEH, Fellowship #1061435 to MI, Fellowship #1072476 to DBA), A*STAR Biomedical Research Council, Singapore (12/1/21/24/6689 to YYT) and the Wellcome Trust UK (research training fellowship #081814/Z/06/Z to MC) and as part of their Major Overseas Program in Viet Nam (089276/Z/09/Z to JF and 106680/B/14/Z to GT).

References

1. Zumla A, et al. Eliminating tuberculosis and tuberculosis-HIV co-disease in the 21st century: key perspectives, controversies, unresolved issues, and needs. *J Infect Dis.* 2012; 205(Suppl 2):S141–6. [PubMed: 22448019]
2. World Health Organisation. WHO | Global tuberculosis report 2017. World Health Organization; 2017.
3. Casali N, et al. Evolution and transmission of drug-resistant tuberculosis in a Russian population. *Nat Genet.* 2014; 46:279–86. [PubMed: 24464101]
4. Guerra-Assuncao JA, et al. Large-scale whole genome sequencing of *M. tuberculosis* provides insights into transmission in a high prevalence area. *Elife.* 2015; 4

5. Guerra-Assuncao JA, et al. Recurrence due to Relapse or Reinfection With Mycobacterium tuberculosis: A Whole-Genome Sequencing Approach in a Large, Population-Based Cohort With a High HIV Infection Prevalence and Active Follow-up. *J Infect Dis.* 2014
6. Coll F, et al. A robust SNP barcode for typing Mycobacterium tuberculosis complex strains. *Nat Commun.* 2014; 5:4812. [PubMed: 25176035]
7. Anh DD, et al. Mycobacterium tuberculosis Beijing genotype emerging in Vietnam. *Emerg Infect Dis.* 2000; 6:302–5. [PubMed: 10827122]
8. Buu TN, et al. The Beijing genotype is associated with young age and multidrug-resistant tuberculosis in rural Vietnam. *Int J Tuberc Lung Dis.* 2009; 13:900–6. [PubMed: 19555542]
9. Maeda S, et al. Mycobacterium tuberculosis strains spreading in Hanoi, Vietnam: Beijing sublineages, genotypes, drug susceptibility patterns, and host factors. *Tuberculosis (Edinb).* 2014; 94:649–56. [PubMed: 25459163]
10. Nguyen VA, et al. High prevalence of Beijing and EAI4-VNM genotypes among M. tuberculosis isolates in northern Vietnam: sampling effect, rural and urban disparities. *PLoS One.* 2012; 7:e45553. [PubMed: 23029091]
11. Nguyen VA, et al. Mycobacterium tuberculosis lineages and anti-tuberculosis drug resistance in reference hospitals across Viet Nam. *BMC Microbiol.* 2016; 16:167. [PubMed: 27464737]
12. Horton KC, MacPherson P, Houben RM, White RG, Corbett EL. Sex Differences in Tuberculosis Burden and Notifications in Low- and Middle-Income Countries: A Systematic Review and Meta-analysis. *PLoS Med.* 2016; 13:e1002119. [PubMed: 27598345]
13. Malla B, et al. First insights into the phylogenetic diversity of Mycobacterium tuberculosis in Nepal. *PLoS One.* 2012; 7:e52297. [PubMed: 23300635]
14. Lee CH, et al. Treatment delay and fatal outcomes of pulmonary tuberculosis in advanced age: a retrospective nationwide cohort study. *BMC Infect Dis.* 2017; 17:449. [PubMed: 28646854]
15. Wlodarska M, Johnston JC, Gardy JL, Tang P. A microbiological revolution meets an ancient disease: improving the management of tuberculosis with genomics. *Clin Microbiol Rev.* 2015; 28:523–39. [PubMed: 25810419]
16. Eldholm V, et al. Four decades of transmission of a multidrug-resistant Mycobacterium tuberculosis outbreak strain. *Nat Commun.* 2015; 6:7119. [PubMed: 25960343]
17. Zhang H, et al. Genome sequencing of 161 Mycobacterium tuberculosis isolates from China identifies genes and intergenic regions associated with drug resistance. *Nat Genet.* 2013; 45:1255–60. [PubMed: 23995137]
18. Merker M, et al. Evolutionary history and global spread of the Mycobacterium tuberculosis Beijing lineage. *Nat Genet.* 2015; 47:242–9. [PubMed: 25599400]
19. Comas I, et al. Population Genomics of Mycobacterium tuberculosis in Ethiopia Contradicts the Virgin Soil Hypothesis for Human Tuberculosis in Sub-Saharan Africa. *Curr Biol.* 2015; 25:3260–6. [PubMed: 26687624]
20. Phelan J, et al. Mycobacterium tuberculosis whole genome sequencing and protein structure modelling provides insights into anti-tuberculosis drug resistance. *BMC Med.* 2016; 14:31. [PubMed: 27005572]
21. Stucki D, et al. Mycobacterium tuberculosis lineage 4 comprises globally distributed and geographically restricted sublineages. *Nat Genet.* 2016
22. Hanekom M, et al. Mycobacterium tuberculosis Beijing genotype: a template for success. *Tuberculosis (Edinb).* 2011; 91:510–23. [PubMed: 21835699]
23. Parwati I, van Crevel R, van Soolingen D. Possible underlying mechanisms for successful emergence of the Mycobacterium tuberculosis Beijing genotype strains. *Lancet Infect Dis.* 2010; 10:103–11. [PubMed: 20113979]
24. Coscolla M, Gagneux S. Consequences of genomic diversity in Mycobacterium tuberculosis. *Semin Immunol.* 2014; 26:431–44. [PubMed: 25453224]
25. van Laarhoven A, et al. Low induction of proinflammatory cytokines parallels evolutionary success of modern strains within the Mycobacterium tuberculosis Beijing genotype. *Infect Immun.* 2013; 81:3750–6. [PubMed: 23897611]
26. Farhat MR, et al. Genomic analysis identifies targets of convergent positive selection in drug-resistant Mycobacterium tuberculosis. *Nat Genet.* 2013; 45:1183–9. [PubMed: 23995135]

27. Hazbon MH, et al. Convergent evolutionary analysis identifies significant mutations in drug resistance targets of *Mycobacterium tuberculosis*. *Antimicrob Agents Chemother*. 2008; 52:3369–76. [PubMed: 18591265]
28. Knudsen NP, et al. Tuberculosis vaccine with high predicted population coverage and compatibility with modern diagnostics. *Proc Natl Acad Sci U S A*. 2014; 111:1096–101. [PubMed: 24395772]
29. Baldwin SL, et al. Intradermal immunization improves protective efficacy of a novel TB vaccine candidate. *Vaccine*. 2009; 27:3063–71. [PubMed: 19428920]
30. Baldwin SL, et al. Protection and Long-Lived Immunity Induced by the ID93/GLA-SE Vaccine Candidate against a Clinical *Mycobacterium tuberculosis* Isolate. *Clin Vaccine Immunol*. 2015; 23:137–47. [PubMed: 26656121]
31. Bertholet S, et al. A defined tuberculosis vaccine candidate boosts BCG and protects against multidrug-resistant *Mycobacterium tuberculosis*. *Sci Transl Med*. 2010; 2 53ra74.
32. Bertholet S, et al. Identification of human T cell antigens for the development of vaccines against *Mycobacterium tuberculosis*. *J Immunol*. 2008; 181:7948–57. [PubMed: 19017986]
33. Uplekar S, Heym B, Friocourt V, Rougemont J, Cole ST. Comparative genomics of Esx genes from clinical isolates of *Mycobacterium tuberculosis* provides evidence for gene conversion and epitope variation. *Infect Immun*. 2011; 79:4042–9. [PubMed: 21807910]
34. Lewinsohn DM, et al. Human *Mycobacterium tuberculosis* CD8 T Cell Antigens/Epitopes Identified by a Proteomic Peptide Library. *PLoS One*. 2013; 8:e67016. [PubMed: 23805289]
35. Kim Y, et al. Immune epitope database analysis resource. *Nucleic Acids Res*. 2012; 40:W525–30. [PubMed: 22610854]
36. Mortimer TD, Weber AM, Pepperell CS. Evolutionary Thrift: *Mycobacteria* Repurpose Plasmid Diversity during Adaptation of Type VII Secretion Systems. *Genome Biol Evol*. 2017; 9:398–413.
37. Gey Van Pittius NC, et al. The ESAT-6 gene cluster of *Mycobacterium tuberculosis* and other high G+C Gram-positive bacteria. *Genome Biol*. 2001; 2 RESEARCH0044.
38. Abdallah AM, et al. PPE and PE_PGRS proteins of *Mycobacterium marinum* are transported via the type VII secretion system ESX-5. *Mol Microbiol*. 2009; 73:329–40. [PubMed: 19602152]
39. Fishbein S, van Wyk N, Warren RM, Sampson SL. Phylogeny to function: PE/PPE protein evolution and impact on *Mycobacterium tuberculosis* pathogenicity. *Mol Microbiol*. 2015; 96:901–16. [PubMed: 25727695]
40. Brennan MJ. The Enigmatic PE/PPE Multigene Family of *Mycobacteria* and Tuberculosis Vaccination. *Infect Immun*. 2017; 85
41. Groschel MI, Sayes F, Simeone R, Majlessi L, Brosch R. ESX secretion systems: mycobacterial evolution to counter host immunity. *Nat Rev Microbiol*. 2016; 14:677–691. [PubMed: 27665717]
42. Shah S, Cannon JR, Fenselau C, Briken V. A Duplicated ESAT-6 Region of ESX-5 Is Involved in Protein Export and Virulence of *Mycobacteria*. *Infect Immun*. 2015; 83:4349–61. [PubMed: 26303392]
43. Kumar A, Chandolia A, Chaudhry U, Brahmachari V, Bose M. Comparison of mammalian cell entry operons of mycobacteria: in silico analysis and expression profiling. *FEMS Immunol Med Microbiol*. 2005; 43:185–95. [PubMed: 15681149]
44. Bukka A, Price CT, Kernodle DS, Graham JE. *Mycobacterium tuberculosis* RNA Expression Patterns in Sputum Bacteria Indicate Secreted Esx Factors Contributing to Growth are Highly Expressed in Active Disease. *Front Microbiol*. 2012; 2:266. [PubMed: 22291682]
45. Coppola M, et al. New Genome-Wide Algorithm Identifies Novel In-Vivo Expressed *Mycobacterium Tuberculosis* Antigens Inducing Human T-Cell Responses with Classical and Unconventional Cytokine Profiles. *Sci Rep*. 2016; 6:37793. [PubMed: 27892960]
46. Abdallah AM, et al. Mycobacterial secretion systems ESX-1 and ESX-5 play distinct roles in host cell death and inflammasome activation. *J Immunol*. 2011; 187:4744–53. [PubMed: 21957139]
47. Ford CB, et al. *Mycobacterium tuberculosis* mutation rate estimates from different lineages predict substantial differences in the emergence of drug-resistant tuberculosis. *Nat Genet*. 2013; 45:784–90. [PubMed: 23749189]
48. Fox GJ, Barry SE, Britton WJ, Marks GB. Contact investigation for tuberculosis: a systematic review and meta-analysis. *Eur Respir J*. 2013; 41:140–56. [PubMed: 22936710]

49. Lonnroth K, et al. Systematic screening for active tuberculosis: rationale, definitions and key considerations. *Int J Tuberc Lung Dis*. 2013; 17:289–98. [PubMed: 23407219]
50. Thai PVK, et al. Bacterial risk factors for treatment failure and relapse among patients with Isoniazid resistant tuberculosis. *BMC Infectious Diseases*. 2018; 18:112. [PubMed: 29510687]
51. Caws M, et al. The influence of host and bacterial genotype on the development of disseminated disease with *Mycobacterium tuberculosis*. *PLoS Pathog*. 2008; 4:e1000034. [PubMed: 18369480]
52. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. 2012; 9:357–9. [PubMed: 22388286]
53. Li H, et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*. 2009; 25:2078–9. [PubMed: 19505943]
54. Comas I, et al. Out-of-Africa migration and Neolithic coexpansion of *Mycobacterium tuberculosis* with modern humans. *Nat Genet*. 2013; 45:1176–82. [PubMed: 23995134]
55. Pepperell CS, et al. The role of selection in shaping diversity of natural *M. tuberculosis* populations. *PLoS Pathog*. 2013; 9:e1003543. [PubMed: 23966858]
56. Bradley P, et al. Rapid antibiotic-resistance predictions from genome sequence data for *Staphylococcus aureus* and *Mycobacterium tuberculosis*. *Nat Commun*. 2015; 6:10063. [PubMed: 26686880]
57. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*. 2014; 30:1312–3. [PubMed: 24451623]
58. Price MN, Dehal PS, Arkin AP. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS One*. 2010; 5:e9490. [PubMed: 20224823]
59. Ashkenazy H, et al. FastML: a web server for probabilistic reconstruction of ancestral sequences. *Nucleic Acids Res*. 2012; 40:W580–4. [PubMed: 22661579]
60. Bollback JP. SIMMAP: stochastic character mapping of discrete traits on phylogenies. *BMC Bioinformatics*. 2006; 7:88. [PubMed: 16504105]
61. Revell L. phytools: an R package for phylogenetic comparative biology (and other things). *Methods in Ecology and Evolution*. 2012; 3:217–223.
62. Vijay S, V DN, Hai HT, Ha VTH, Dung VTM, Dinh TD, Nhung HN, Tram TTB, Aldridge BB, Hanh NT, Thu DDA, et al. Influence of Stress and Antibiotic Resistance on Cell-Length Distribution in *Mycobacterium tuberculosis* Clinical Isolates. *Frontiers in Microbiology*. 2017; 8:1–12. [PubMed: 28197127]
63. Rohde KH, Abramovitch RB, Russell DG. *Mycobacterium tuberculosis* invasion of macrophages: linking bacterial gene expression to environmental cues. *Cell Host Microbe*. 2007; 2:352–64. [PubMed: 18005756]
64. Carver T, et al. BamView: visualizing and interpretation of next-generation sequencing read alignments. *Brief Bioinform*. 2013; 14:203–12. [PubMed: 22253280]
65. Jafri M, et al. Germline Mutations in the CDKN2B Tumor Suppressor Gene Predispose to Renal Cell Carcinoma. *Cancer Discov*. 2015; 5:723–9. [PubMed: 25873077]
66. Usher JL, et al. Analysis of HGD Gene Mutations in Patients with Alkaptonuria from the United Kingdom: Identification of Novel Mutations. *JIMD Rep*. 2015; 24:3–11. [PubMed: 25681086]
67. Jubb HC, et al. Mutations at protein-protein interfaces: Small changes over big surfaces have large impacts on human health. *Prog Biophys Mol Biol*. 2017; 128:3–13. [PubMed: 27913149]
68. Kano FS, et al. The Presence, Persistence and Functional Properties of Plasmodium vivax Duffy Binding Protein II Antibodies Are Influenced by HLA Class II Allelic Variants. *PLoS Negl Trop Dis*. 2016; 10:e0005177. [PubMed: 27959918]
69. Nemethova M, et al. Twelve novel HGD gene variants identified in 99 alkaptonuria patients: focus on 'black bone disease' in Italy. *Eur J Hum Genet*. 2016; 24:66–72. [PubMed: 25804398]
70. Silvino AC, et al. Variation in Human Cytochrome P-450 Drug-Metabolism Genes: A Gateway to the Understanding of Plasmodium vivax Relapses. *PLoS One*. 2016; 11:e0160172. [PubMed: 27467145]
71. White RR, et al. Ubiquitin-Dependent Modification of Skeletal Muscle by the Parasitic Nematode, *Trichinella spiralis*. *PLoS Pathog*. 2016; 12:e1005977. [PubMed: 27870901]

72. Albanaz ATS, Rodrigues CHM, Pires DEV, Ascher DB. Combating mutations in genetic disease and drug resistance: understanding molecular mechanisms to guide drug design. *Expert Opin Drug Discov.* 2017; 12:553–563. [PubMed: 28490289]
73. Casey RT, et al. SDHA related tumorigenesis: a new case series and literature review for variant interpretation and pathogenicity. *Mol Genet Genomic Med.* 2017; 5:237–250. [PubMed: 28546994]
74. Pandurangan AP, Ascher DB, Thomas SE, Blundell TL. Genomes, structural biology and drug discovery: combating the impacts of mutations in genetic disease and antibiotic resistance. *Biochem Soc Trans.* 2017; 45:303–311. [PubMed: 28408471]
75. Soardi FC, et al. Familial STAG2 germline mutation defines a new human cohesinopathy. *Npj Genomic Medicine.* 2017; 2
76. Sali A, Blundell TL. Comparative protein modelling by satisfaction of spatial restraints. *J Mol Biol.* 1993; 234:779–815. [PubMed: 8254673]
77. Pires DE, Chen J, Blundell TL, Ascher DB. In silico functional dissection of saturation mutagenesis: Interpreting the relationship between phenotypes and changes in protein stability, interactions and activity. *Sci Rep.* 2016; 6 19848.
78. Pandurangan AP, Ochoa-Montano B, Ascher DB, Blundell TL. SDM: a server for predicting effects of mutations on protein stability. *Nucleic Acids Res.* 2017
79. Pires DE, Ascher DB, Blundell TL. mCSM: predicting the effects of mutations in proteins using graph-based signatures. *Bioinformatics.* 2014; 30:335–42. [PubMed: 24281696]
80. Pires DE, Ascher DB, Blundell TL. DUET: a server for predicting effects of mutations on protein stability using an integrated computational approach. *Nucleic Acids Res.* 2014; 42:W314–9. [PubMed: 24829462]
81. Mahmood A, et al. Molecular characterization of secretory proteins Rv3619c and Rv3620c from *Mycobacterium tuberculosis* H37Rv. *FEBS J.* 2011; 278:341–53. [PubMed: 21134129]
82. Ascher DB, et al. Potent hepatitis C inhibitors bind directly to NS5A and reduce its affinity for RNA. *Sci Rep.* 2014; 4:4765. [PubMed: 24755925]
83. Chan LJ, et al. Conjugation of 10 kDa Linear PEG onto Trastuzumab Fab' Is Sufficient to Significantly Enhance Lymphatic Exposure while Preserving in Vitro Biological Activity. *Mol Pharm.* 2016; 13:1229–41. [PubMed: 26871003]
84. Chan LJ, et al. PEGylation does not significantly change the initial intravenous or subcutaneous pharmacokinetics or lymphatic exposure of trastuzumab in rats but increases plasma clearance after subcutaneous administration. *Mol Pharm.* 2015; 12:794–809. [PubMed: 25644368]

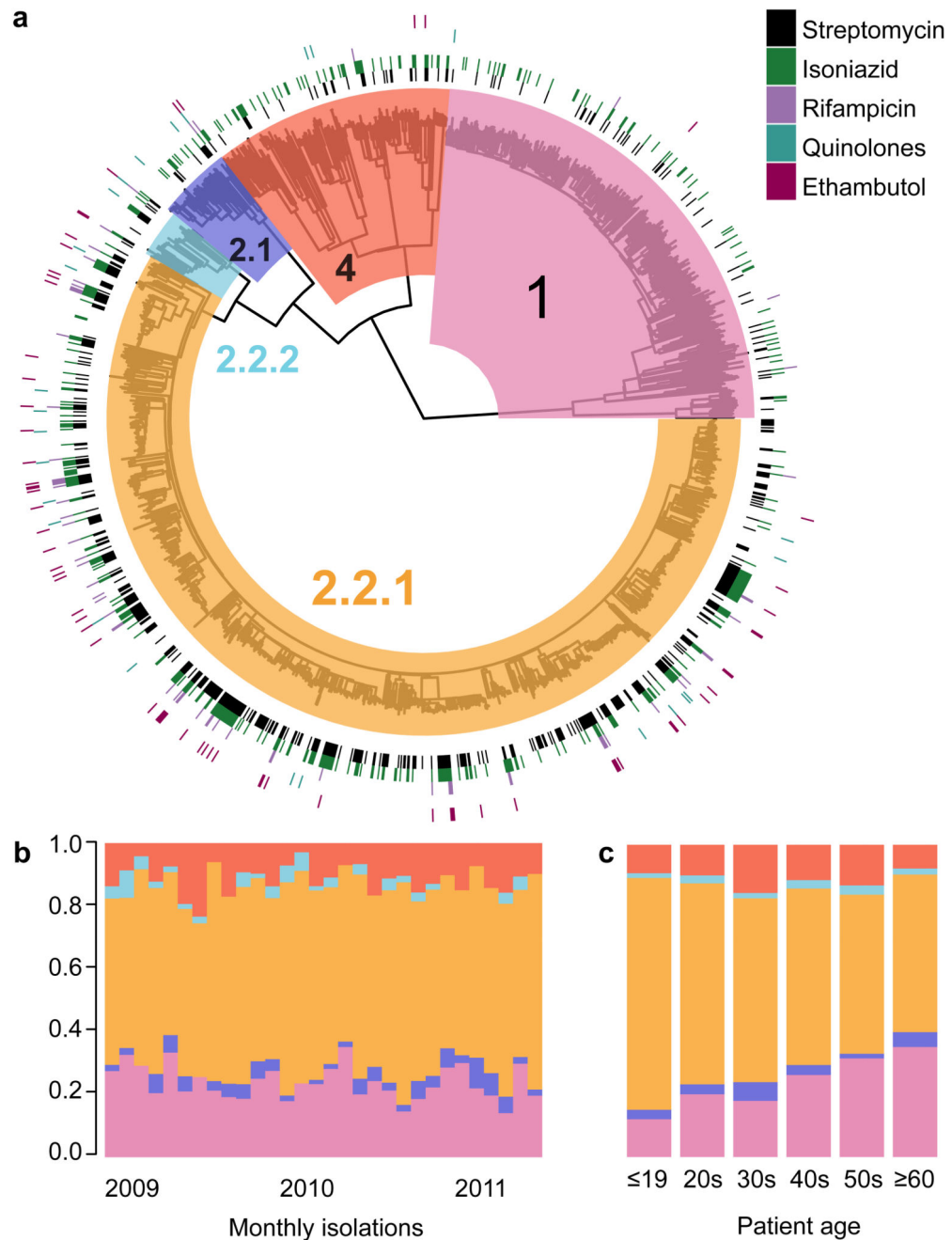


Figure 1. Circulating *M. tuberculosis* strains in HCMC are divided into multiple distinct lineages. (a) Maximum-likelihood phylogeny of 1635 *Mtb* isolates collected from TB patients in HCMC, with backgrounds shaded by lineage. Exterior rings indicate presence of known antimicrobial resistance-associated mutations (coloured by drug, according to legend in top right). (b) Frequency distribution of lineages by month. (c) Frequency distribution of lineages by patient age group.

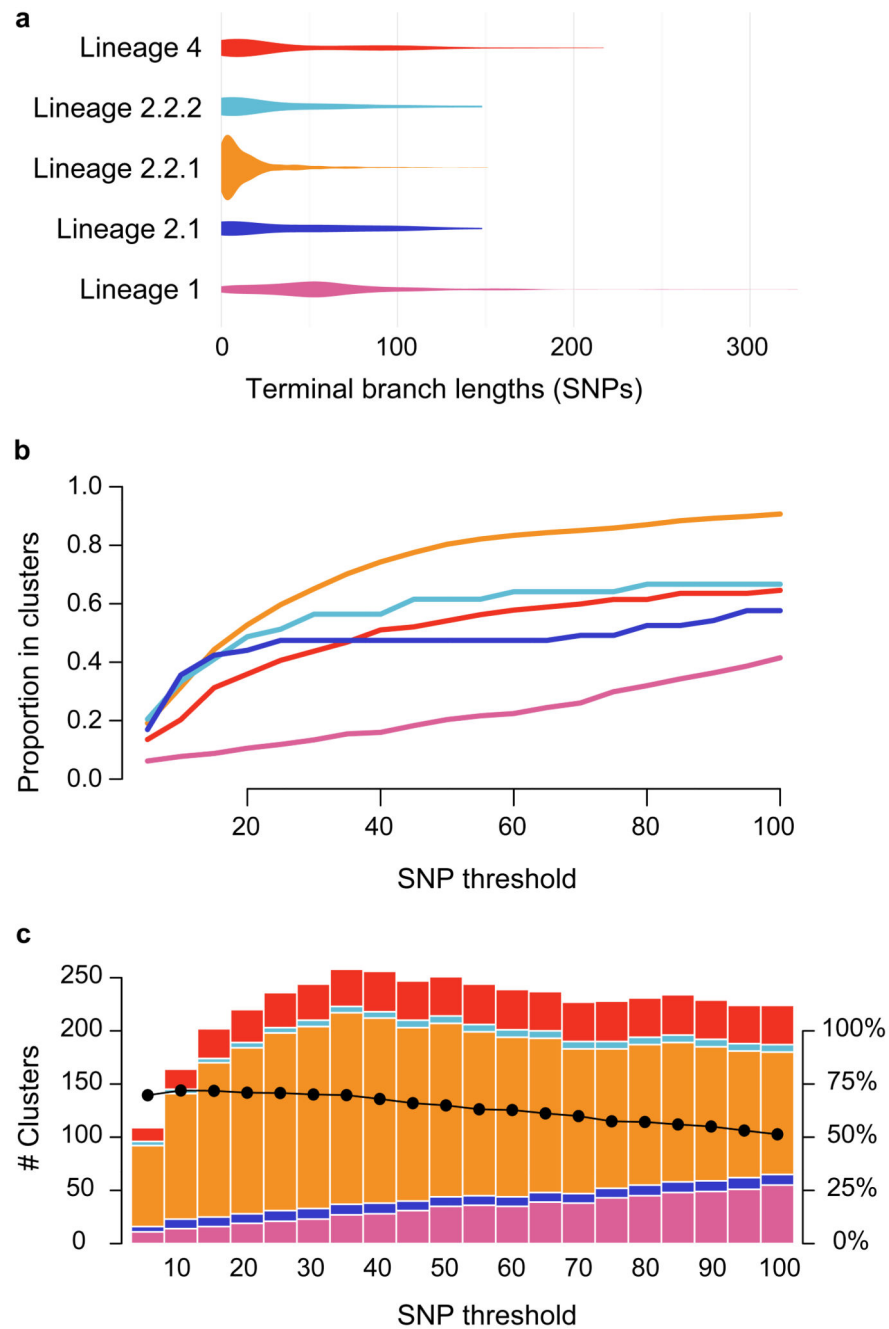


Figure 2. Properties of lineage subtrees for HCMC *M. tuberculosis* genomes.

(a) Distributions of terminal branch lengths for the 1635-strain phylogeny. **(b)** Mean subtree heights (y-axis; measured as mean node-to-tip distances for each subtree) vs subtree size (x-axis; number of descendant tips). Shaded region indicates standard error of the mean across subtrees of a given size; labels indicate lineage. **(c)** Stacked area plot showing number of clusters (y-axis) within each lineage (coloured as in panel a) identified using different maximum patristic distance thresholds to define clusters (x-axis).

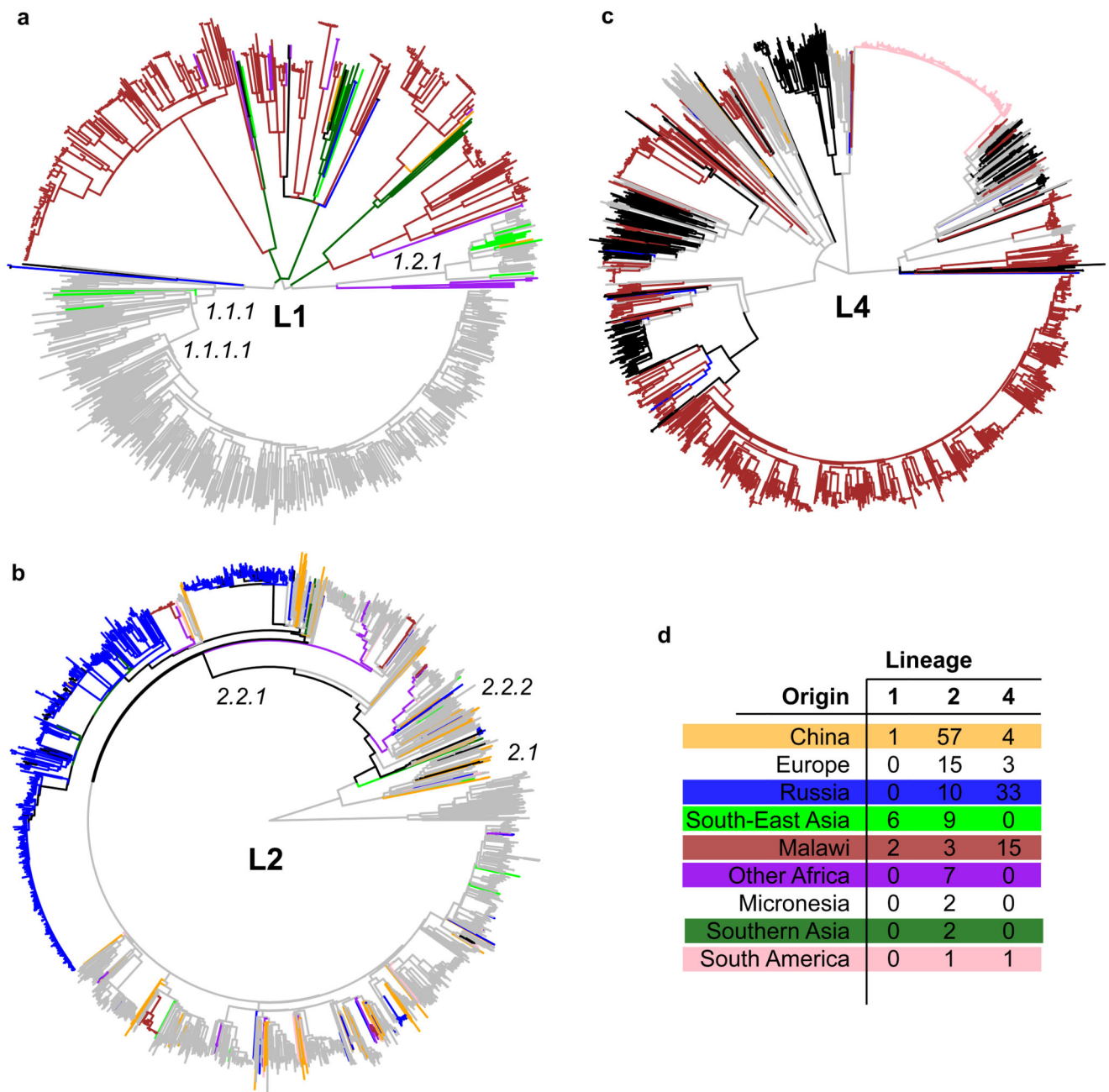


Figure 3. Phylogenies of *M. tuberculosis* showing relationships between isolates from HCMC and other locations.

HCMC isolates are coloured grey, isolates from four other localised studies are coloured as in panel (d), other locations are shown in black. (a) Lineage 1 (n=675 genomes). (b) Lineage 2 (n=1871 genomes). (c) Lineage 4 (n=2066 genomes). (d) Number of transfers between Vietnam and other locations predicted by stochastic mapping of locations onto the Lineage 2 and 4 trees.

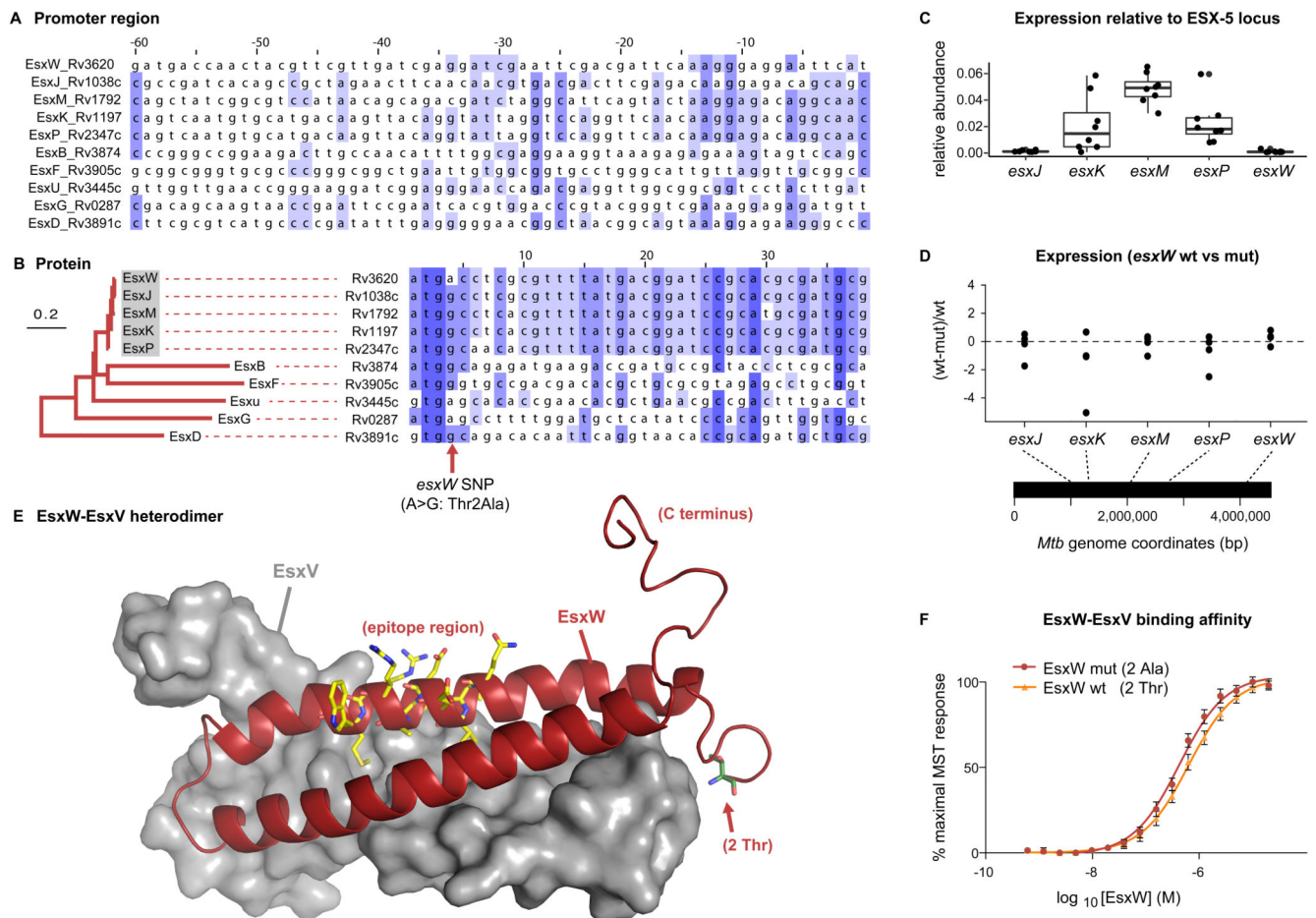


Figure 4. EsxW mutation at gene, mRNA, protein and heterodimer level.

(a) Variation in promoter region of *esxW* and other CFP10 paralogs extracted from H37Rv (Lineage 4). Sites are coloured by conservation, coordinates are relative to the start codon. (b) Variation in protein coding region. Tree shows maximum likelihood phylogeny inferred from amino acid sequences; box indicates QILSS proteins. Alignment of protein-coding DNA sequence is shown, coloured by conservation; arrow indicates site of homoplasic SNP in *esxW*, (A to G), resulting in Thr to Ala substitution at EsxW protein residue 2. (c-d) RNAseq results for *esxW* and QILSS/ESX-5 paralogs measured in 4 *Mtb* isolate pairs, each including a Lineage 1 or 4 EsxW-2Ala mutant and its genetically closest EsxW-2Thr relative, following 24h macrophage infection. mRNA levels were estimated from read counts uniquely mapping to the region -21 to +9 for each gene; normalized to total reads uniquely mapping to the locus encoding the ESX-5 machinery (and *esxM*) in each isolate. Boxes indicate interquartile range. (d) Difference between 2Ala mutant and 2Thr wildtype for each pair, relative to wildtype expression level. (e) Structural model of EsxW-EsxV heterodimer. EsxV is shown as a surface (grey) and EsxW as a ribbon (red) with key residues shown as labeled sticks. (f) Comparison of biophysical measurements of heterodimer binding affinity between wildtype and mutant EsxW. All binding curves were

determined across four replicates by microscale thermophoresis (MST), and are represented as the mean \pm standard deviation.

Table 1
Lineage characteristics for HCMC *M. tuberculosis* isolates, including known antimicrobial resistance mutations identified using Mykrobe Predictor.

	Lineage									
	1		2.1		2.2.1		2.2.2		4	
	N	%	N	%	N	%	N	%	N	%
Gender										
Female	82	21.1	9	15.3	265	27.8	10	25.6	56	29.2
Male	306	78.9	50	84.7	692	72.3	29	74.4	136	70.8
Antimicrobial										
Streptomycin	48	12.4	10	17.0	426	44.5	12	30.8	30	15.6
Isoniazid	57	14.7	12	20.3	269	28.1	9	23.1	52	27.1
Rifampicin	3	0.8	2	3.4	58	6.1	2	5.1	1	0.5
Quinolones	1	0.3	3	5.1	18	1.9	2	5.1	2	1.0
Ethambutol	1	0.3	2	3.4	60	6.3	3	7.7	2	1.0

Table 2
Homoplasic non-synonymous SNPs identified as occurring on the Beijing lineage-defining branch and also arising independently within other lineages.

The number of branches on which each SNP was identified outside the Beijing lineage-defining branch is shown. The number of such branches that have multiple descendant tips (indicating onward transmission of the SNP) is shown in no. transmitted column. HCMC refers to the 1,635 isolates from HCMC, Vietnam; Elsewhere refers to the 3,146 additional isolates from published studies^{3–5, 17–19}; trees are shown in Figure 3.

Mutation	HCMC		Elsewhere		Function
	no. branches outside Beijing lineage	no. transmitted	no. branches outside Beijing lineage	no. transmitted	
EsxW-T2A	9	4	10	6	ESX-5 secreted protein (CFP10 homolog)
Rv3081-F220L	2	1	7	3	hypothetical protein
GidB-E92D	1	1	0	0	streptomycin resistance