

China's tuberculosis epidemic stems from historical expansion of four strains of *Mycobacterium tuberculosis*

Qingyun Liu^{1,2}, Aijing Ma^{3,14}, Lanhai Wei^{4,14}, Yu Pang⁵, Beibei Wu⁶, Tao Luo⁷, Yang Zhou³, Hong-Xiang Zheng⁴, Qi Jiang^{1,2}, Mingyu Gan^{1,2}, Tianyu Zuo¹, Mei Liu¹, Chongguang Yang^{1,8}, Li Jin⁴, Iñaki Comas⁹, Sébastien Gagneux^{10,11}, Yanlin Zhao^{3*}, Caitlin S. Pepperell^{12,13*} and Qian Gao^{1,2*}

A small number of high-burden countries account for the majority of tuberculosis cases worldwide. Detailed data are lacking from these regions. To explore the evolutionary history of *Mycobacterium tuberculosis* in China—the country with the third highest tuberculosis burden—we analysed a countrywide collection of 4,578 isolates. Little genetic diversity was detected, with 99.4% of the bacterial population belonging to lineage 2 and three sublineages of lineage 4. The deeply rooted phylogenetic positions and geographic restriction of these four genotypes indicate that their populations expanded *in situ* following a small number of introductions to China. Coalescent analyses suggest that these bacterial subpopulations emerged in China around 1,000 years ago, and expanded in parallel from the twelfth century onwards, and that the whole population peaked in the late eighteenth century. More recently, sublineage L2.3, which is indigenous to China and exhibited relatively high transmissibility and extensive global dissemination, came to dominate the population dynamics of *M. tuberculosis* in China. Our results indicate that historical expansion of four *M. tuberculosis* strains shaped the current tuberculosis epidemic in China, and highlight the long-term genetic continuity of the indigenous *M. tuberculosis* population.

Mycobacterium tuberculosis complex (MTBC), which causes tuberculosis, has circulated among human populations for thousands of years¹. With more than 10 million cases and 1.7 million deaths each year, tuberculosis remains the leading cause of death due to an infectious disease. The burden of tuberculosis is unevenly distributed, with 30 tuberculosis high-burden countries accounting for 87% of all tuberculosis cases in the world². Tuberculosis is a typical disease of poverty and all the high-burden countries are in the developing world^{3,4}.

China ranks as the country with the third highest tuberculosis burden in the world, with about one million incident cases each year². Numerous literary sources describe a disease resembling tuberculosis in ancient China, suggesting that tuberculosis has affected Chinese populations for thousands of years^{5–8}. The oldest literary description suggestive of tuberculosis is from ~5,700 years ago, predating the first dynasty Xia in China⁶. MTBC DNA was detected in human skeletons from Xinjiang province dating back ~2,000 years before the present⁵. Based on these observations, the current epidemic of tuberculosis in China may have very deep historical roots. However, unlike Western Europe, which was devastated by the so-called white plague of tuberculosis during the

eighteenth and nineteenth centuries^{9,10}, the historical record does not contain any similar descriptions of severe epidemics of tuberculosis in China^{11,12}. Therefore, it remains unclear when epidemic forms of tuberculosis first arose in East Asia, and what course these epidemics may have followed throughout Chinese history^{11,13}. Starting in the second half of the past century, China underwent major social changes, including strong population growth, massive internal migrations of rural workers to urban areas, increases in household crowding, and later the incursion of the human immunodeficiency virus pandemic^{14,15}. It is not clear what role these factors have played in enabling epidemic forms of tuberculosis to flourish in the more recent past.

The global spread of MTBC strains is mainly driven by human activities. Previous research suggests it was globally disseminated as a result of human movements driven by exploration, migration, trade and conquest^{16–19}. Although more recent human movements linked to globalization might disturb phylogeographic patterns, MTBC lineages still vary in their distribution between countries and continents^{20,21}. These geographic patterns can shed light on historical phenomena that contributed to the spread of tuberculosis^{22,23}. Whole-genome sequencing (WGS) studies of MTBC populations in

¹Key Laboratory of Medical Molecular Virology, Ministry of Education and Health, School of Basic Medical Sciences, Shanghai Public Health Clinical Center, Fudan University, Shanghai, China. ²Shenzhen Center for Chronic Disease Control, Shenzhen, China. ³National Center for Tuberculosis Control and Prevention, Chinese Center for Disease Control and Prevention, Beijing, China. ⁴State Key Laboratory of Genetic Engineering and Ministry of Education Key Laboratory of Contemporary Anthropology, School of Life Sciences and Institutes of Biomedical Sciences, Fudan University, Shanghai, China. ⁵National Tuberculosis Clinical Laboratory, Beijing Key Laboratory for Drug Resistance Tuberculosis Research, Beijing Tuberculosis and Thoracic Tumor Research Institute, Beijing Chest Hospital, Capital Medical University, Beijing, China. ⁶The Institute of TB Control, Zhejiang Provincial Center for Disease Control and Prevention, Hangzhou, China. ⁷West China School of Basic Medical Sciences and Forensic Medicines, Sichuan University, Chengdu, China. ⁸Department of Epidemiology of Microbial Diseases, School of Public Health, Yale University, New Haven, CT, USA. ⁹Institute of Biomedicine of Valencia, CSIC and CIBER in Epidemiology and Public Health, Valencia, Spain. ¹⁰Swiss Tropical and Public Health Institute, Basel, Switzerland. ¹¹University of Basel, Basel, Switzerland. ¹²Department of Medicine, Division of Infectious Diseases, University of Wisconsin-Madison, Madison, WI, USA. ¹³Department of Medical Microbiology and Immunology, University of Wisconsin-Madison, Madison, WI, USA. ¹⁴These authors contributed equally: Aijing Ma, Lanhai Wei.
*e-mail: zhaoyanlin@chinabt.org; cspepper@medicine.wisc.edu; qgao99@yahoo.com

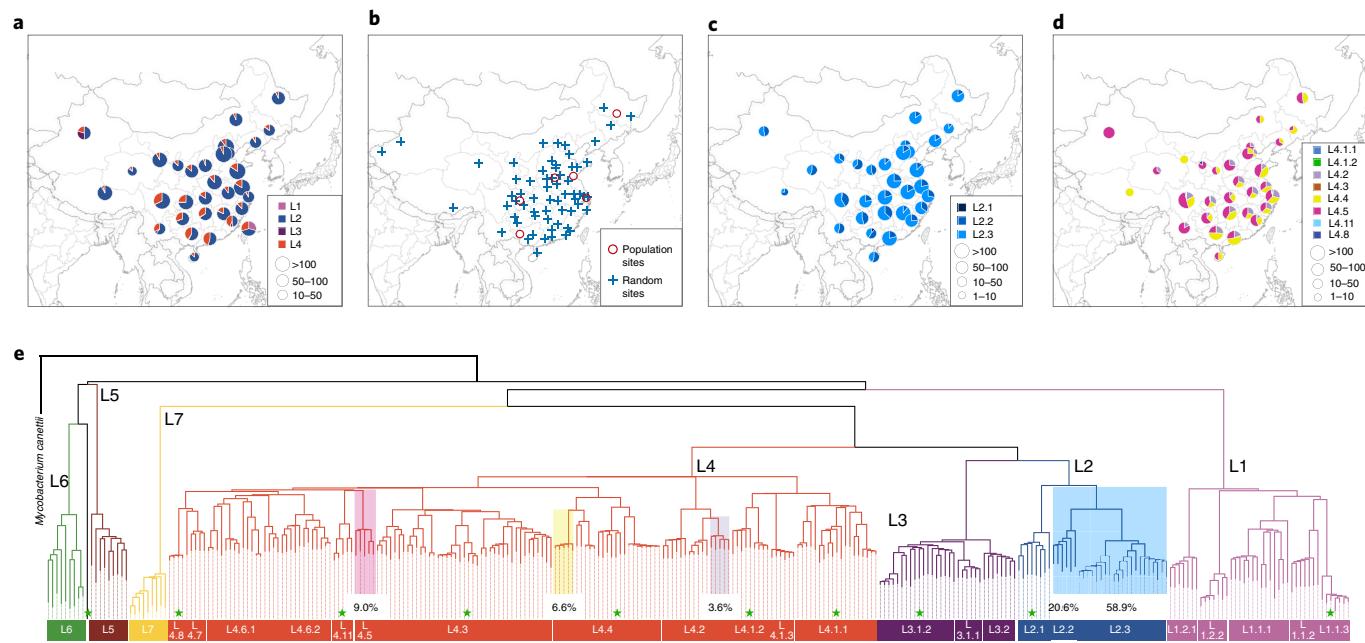


Fig. 1 | Genotyping results of countrywide collected MTBC strains in China. **a**, Prevalence of different MTBC lineages in 32 provinces based on spoligotyping data from 16,221 isolates collected throughout China. **b**, The 76 county sites from which MTBC isolates were sampled for this study. ‘Population sites’ are counties where MTBC isolates were collected through exhaustive sampling from 2009–2010. ‘Random sites’ are counties where MTBC isolates were randomly sampled in 2007. **c,d**, SNP typing results of L2 (**c**) and L4 strains (**d**) show the relative proportion of each sublineage in each province. **e**, Phylogeny of 301 MTBC isolates reflecting the diversity found worldwide. Branches are coloured according to the convention described in ref. ³¹. Sublineages found commonly in China are highlighted, with a notation of their prevalence. Sublineages rarely encountered in China are marked with green stars. The remaining unmarked sublineages were not identified in China.

Greenland and Nunavik have demonstrated a single recent origin for the regional population of pathogens^{24,25}, while studies from Ethiopia and Vietnam point to more ancient and complex origins for the endemic population of MTBC^{26,27}. In this study, we aimed to investigate historical migration events and bacterial population history underlying the current tuberculosis epidemic in China. To this end, we integrated analyses of genotyping and WGS data from 4,578 MTBC isolates collected throughout China. Our findings demonstrate that, although the extant population of MTBC in China is large, almost all strains currently circulating in the country descend from only four ancestors introduced into China around the turn of the second millennium.

Results

A large MTBC population with low genetic diversity. We first collected spoligotyping data from MTBC isolates sampled throughout China to obtain an initial picture of the MTBC population structure in the country. Spoligotyping records for a total of 16,621 isolates were obtained from 26 studies covering all 32 provinces in the country, and 15,217 of them could be successfully assigned to known lineages (Supplementary Table 1). Among these, 12,302 (80.8%) were classified as MTBC lineage 2 (L2), 2,570 (16.9%) were assigned to lineage 4 (L4), 227 (1.5%) were assigned to lineage 1 (L1) and 118 (0.8%) were assigned to lineage 3 (L3) (Fig. 1a). These data are in line with the previous observations that the Beijing family strains (belonging to L2) are most prevalent in China^{28,29}. Our results indicate that L2 and L4 are distributed throughout the country, whereas L1 is most prevalent in Taiwan (southeast) and L3 is most prevalent in Xinjiang (northwest) (Fig. 1a). The overwhelming majority of tuberculosis cases in China were caused by L2 and L4 strains.

To characterize the genetic diversity within L2 and L4, we further genotyped a countrywide collection of 4,578 MTBC isolates using previously validated phylogenetic single nucleotide polymorphisms

(SNPs)^{21,30} (Fig. 1b and Supplementary Table 2). Of the 4,578 MTBC isolates, 99% were from L2 (79.8%) and L4 (19.6%). Within L2, sublineage L2.3 (‘modern’ Beijing) was the most prevalent throughout the country and accounted for 73.9% of all L2 isolates (Fig. 1c); L2.2 (‘ancient’ Beijing) accounted for 25.9% and was widely distributed, with a greater concentration in the south relative to the north; and L2.1 (proto-Beijing) showed a low frequency (0.2%) and was restricted to provinces in far Southwest China (Fig. 1c). A total of 8 L4 sublineages were found, but 96.7% of all L4 strains in China belonged to only 3 of these sublineages (L4.2, L4.4 and L4.5), which were widely distributed throughout the country. The other sublineages of L4 were identified sporadically across the country (Fig. 1d and Supplementary Table 2). Taken together, these results indicate that the tuberculosis epidemic in China largely traces its origin to a handful of MTBC sublineages (Fig. 1e).

To put the MTBC strains from China into the wider context of the global population of the MTBC, we selected 306 representative isolates out of 4,578 genotyped strains for WGS (Supplementary Table 3). In addition, we analysed 15,591 previously published MTBC genomes to represent the global diversity of the MTBC (Supplementary Table 4). These data illustrate the homogeneous composition of the MTBC population in China, in contrast with most other countries harbouring a greater diversity of MTBC (Fig. 2a). A formal comparison of diversity, pairwise SNP genetic distances, nucleotide diversity (π) and rarefaction analyses consistently demonstrated that the genetic diversity of MTBC in China as a whole was lower than that of regional samples from individual cities in other countries (Fig. 2b–d). It is striking that the world’s third-largest MTBC population exhibits so little genetic diversity.

Phylogeographically restricted patterns indicate single origins. Published phylogeographic studies have inferred an African origin for the MTBC^{16,18,31}, suggesting that MTBC strains in other

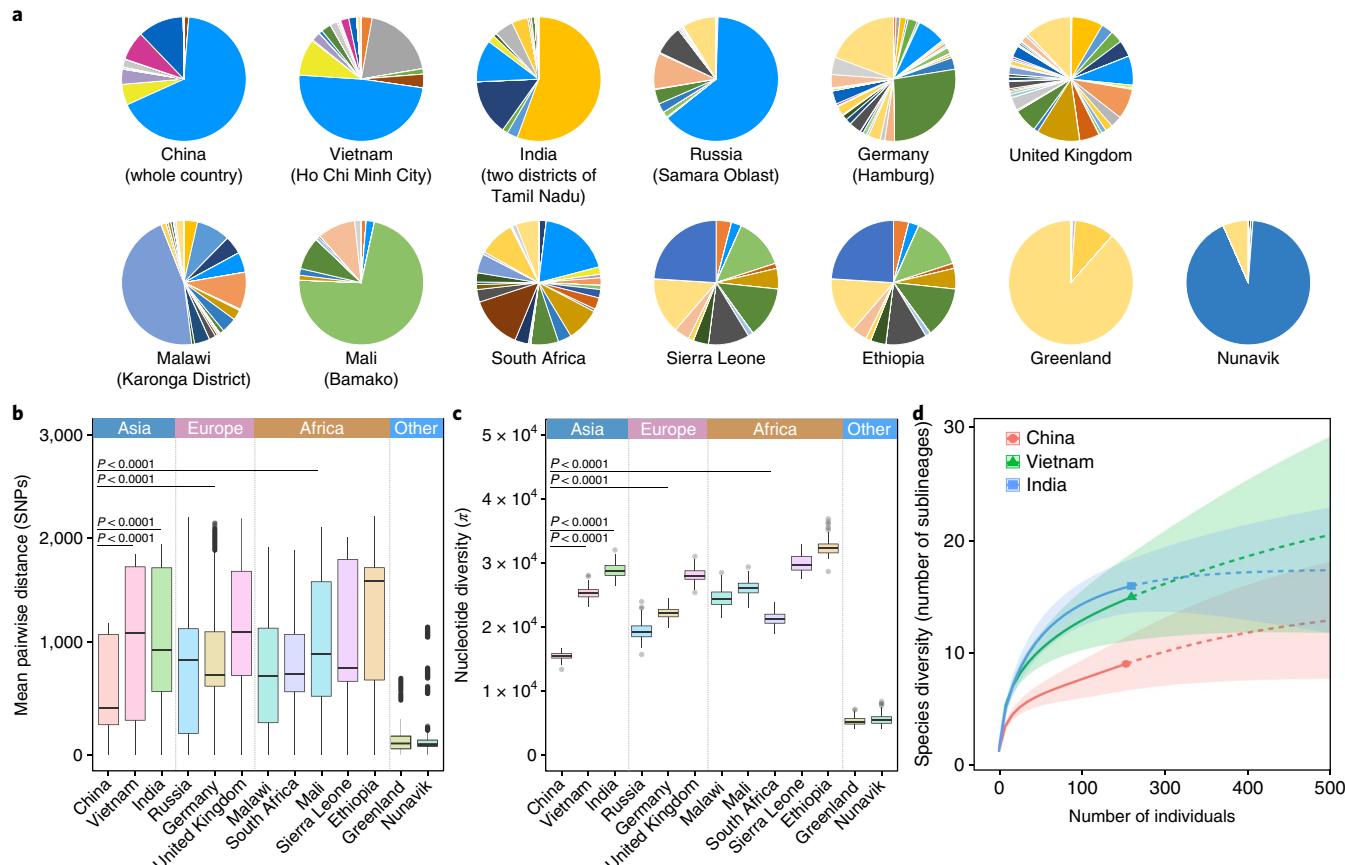


Fig. 2 | Low genetic diversity in China's MTBC population. **a**, Pie charts showing the relative prevalences of MTBC sublineages in different countries and regions, where each sublineage is assigned a colour according to a recent defining scheme. **b,c**, Mean pairwise SNP distance between MTBC strains (b) and nucleotide diversity (c) in the MTBC population from each country or region. Error bars represent 95% CIs, box edges refer to the 95% percentiles and the solid horizontal lines represent the median values. **d**, Rarefaction analysis predicted the sublineage diversity of the MTBC population in China, India and Vietnam. Two hundred isolates were randomly sampled from each of the three countries; solid lines show the captured sublineages while dashed lines show the predicted changes as the sample size is increased. Shading represents 95% CIs.

continents were introduced via human migration. The diversity found within MTBC sublineages in China could have been generated before or after the introduction of the ancestral strains. To distinguish between these possibilities, we determined the phylogenetic positions of strains from China in global MTBC phylogenies. In the L2 phylogeny, although the most recent common ancestor of L2 has a mixed probability of Chinese or Southeast Asian origin, the ancestors of L2.2 and L2.3 have a predicted origin in China, with posterior probabilities of 99.6 and 98.6%, respectively (Fig. 3a and Supplementary Figs. 1 and 2). This observation was robust to resampling analyses, in which we randomly reduced the number of isolates from China (Supplementary Fig. 3). These results indicate that L2 diversified locally following the migration event that established its most recent common ancestor (L2.2) in China. Moreover, the globally extant L2 appears to trace to more recent migration events out of China and Southeast Asia.

In contrast, the tree topology of L4 revealed a deep separation between distinct L4 sublineages found in China, with other global clades interspersed between the Chinese clades (Fig. 3b). This indicates that the indigenous sublineages of L4 diverged before their introduction to China. Chinese strains of L4.2 and L4.4 formed sister clades to strains from the same sublineage found in other regions, while the L4.5 branch is almost entirely composed of strains from China, except for two early-diverged strains that were isolated in Russia. Average pairwise genetic distances between strains from

China versus other regions were 407 and 526 SNPs for L4.2 and L4.4, respectively, while the corresponding distances between strains from within China were 267 and 298 SNPs, respectively ($P < 0.0001$, permutation test) (Supplementary Fig. 4). The ancestors of the Chinese clades in L4.2, L4.4 and L4.5 consistently had a most likely origin in China, with posterior probabilities of 99.4, 96.1 and 99.6%, respectively (Supplementary Fig. 5). These results suggest that L4 strains in China diversified locally, with minimal global exchanges following their original introduction. We identified a novel sublineage provisionally termed L4.11 that appears to be specific to China (Fig. 3b), and principal component analysis showed a clear division from its sister clades L4.5 and L4.6 (Supplementary Fig. 6). These results suggest that the indigenous L4 sublineages arose from separate parallel migrations followed by local diversification.

Historical origins and expansions of the indigenous MTBC population. We estimate that L2 emerged around AD 223, L2.2 in AD 806 and L2.3 in AD 1520 under the MTBC-6 substitution rate model. These results are very similar to previous estimates using the same model^{1,16}. The indigenous L4 sublineages originated over a similar timescale (Table 1 and Fig. 3b). It is remarkable that new sublineages formed *in situ* and important external introduction events occurred during a short window (AD 1150–1268) (Supplementary Fig. 7). These numerous contemporaneous emergences indicate that epidemic expansion of tuberculosis was occurring during this period,

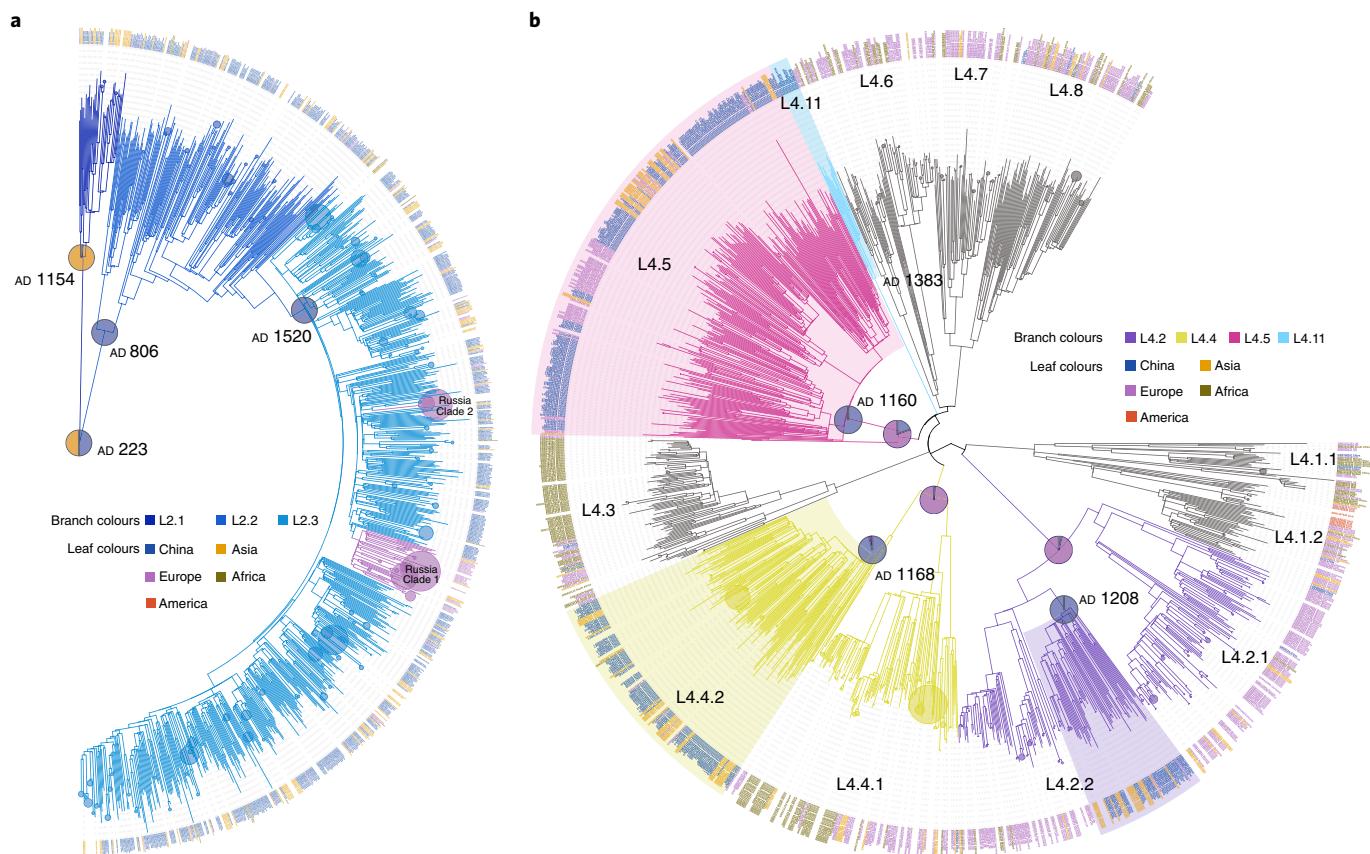


Fig. 3 | Single origins of the four indigenous genotypes. **a,b,** The phylogenetic trees of L2 (**a**) and L4 (**b**) were reconstructed with 1,242 and 1,569 isolates, respectively. To reduce the complexity in both trees, terminal branches with branch lengths <0.008 (indicating that clusters diversified very recently; for example, the two Russian clades in L2.3) were automatically collapsed into circles. The circle sizes of these collapsed branches are proportional to the number of leaves that were collapsed. The estimated origin times of indigenous genotypes are shown at the relevant nodes, and their inferred geographic states are shown as pies with the colours indicating the isolates' country of origin. 'Asia' refers to Asian countries and regions excluding China. The three dominant clades of L4 in China are highlighted.

Table 1 | Time to most recent common ancestor and population growth rates of the indigenous sublineages

Sublineage	Time to most recent common ancestor	Start of growth (AD) ^a	Fast growth interval (AD) ^a	Growth rate by interval (%) ^a		
					MTBC-6 (AD)	95% HPD (AD)
L2.3	1520	1311-1726	1504	1504-1816	2.320	
L2.2	806	250-1272	858	1105-1750	0.669	
L4.2 ^b	1208	823-1528	1365	1365-1560	1.189	
L4.4 ^b	1268	946-1576	1285	1400-1610	1.556	
L4.5 ^b	1160	787-1510	1152	1240-1560	1.610	

^aValues in these columns were calculated using the MTBC-6 model results. ^bHere, L4.2, L4.4 and L4.5 only refer to the Chinese-specific clades in the relative sublineages.

L2.2 and the three L4 sublineages appear to have followed parallel demographic trajectories: each started to expand around the twelfth century, then experienced two or three waves of major expansions before plateauing and then going through a precipitous decline starting in the 1950s (Fig. 4a).

To explore a potential correlation between MTBC population growth and human demography, we compared the human population growth curve³² with the MTBC N_e curve (Fig. 4b). We observed similar trends between human and MTBC population growth: the first sharp increase in the MTBC N_e followed the major human population expansion during the Song dynasty (AD 960–1279); the second substantial increase in the MTBC N_e was parallel to the population boom during the Qing dynasty (AD 1616–1912); and there was a period (thirteenth to sixteenth century) when both human and MTBC populations tended to be stationary. More recently (second half of the twentieth century), we observe a dramatic decline in the MTBC N_e that coincides with the availability of anti-tuberculosis drugs and improved control of the tuberculosis epidemic: China has halved the prevalence of tuberculosis in the past 20 years and maintained a steady decrease in the tuberculosis incidence each year³³. The effects of drug therapy on MTBC diversity may extend beyond a decrease in case counts, in that positive selection imposed by antibiotics could promote displacement of diverse lineages via clonal expansion of drug-resistant strains³⁴. Twentieth-century environments may also have introduced greater individual-to-individual variability in tuberculosis transmission as high-density, crowded locations allow extended chains of

probably as a result of environmental conditions that favoured the pathogen. None of the currently prevalent sublineages was introduced after AD 1383, indicating that the current tuberculosis epidemic in China has largely been shaped by introduction events in the early second millennium. We identified stepwise growths in the effective population size (N_e) in all of the indigenous sublineages between the twelfth and eighteenth centuries (Fig. 4a). Intriguingly,

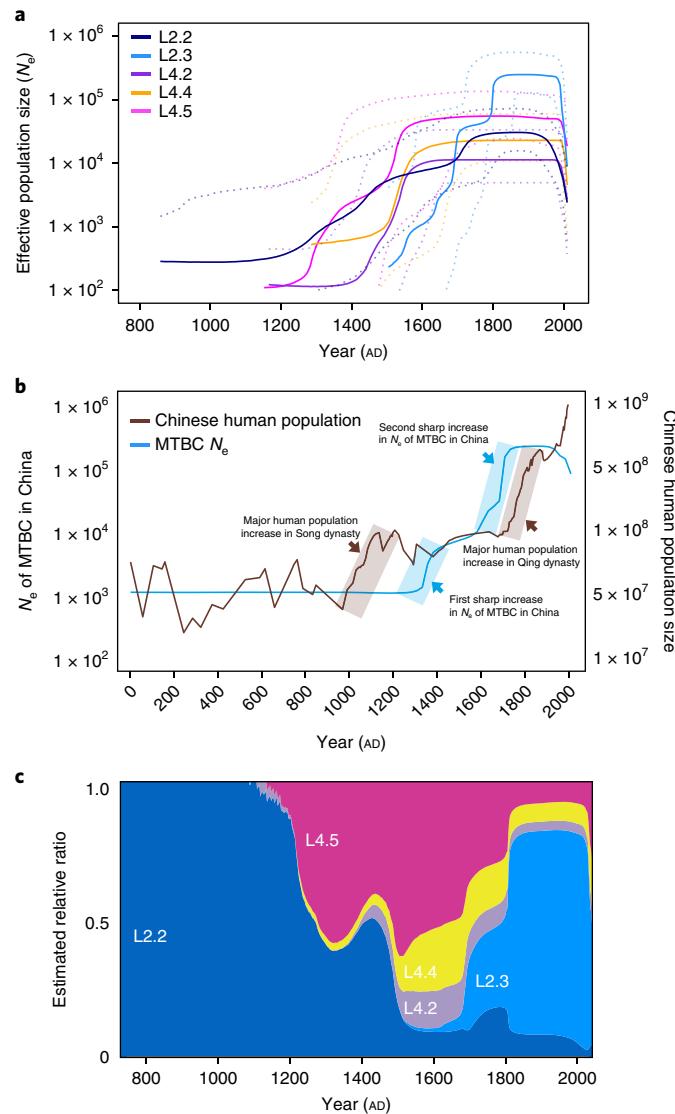


Fig. 4 | Historical expansions of indigenous MTBC genotypes. **a**, Estimated effective population size changes of the major MTBC sublineages in China. L2.3 was separated from L2.2 in the BSP analysis, and the dashed lines represent the 95% HPD. **b**, Comparison of the Chinese human population growth curve and MTBC N_e curve (all indigenous genotypes). **c**, Inferred past population dynamics of each sublineage in China, estimated from the effective population growth.

transmission to develop; this transmission variability is also expected to be reflected in reduced pathogen N_e ³⁵. These concordances between human historical phenomena and bacterial demography suggest that the current epidemic of tuberculosis in China was enabled by the historical growth of human populations, probably including nonlinear effects such as crowding and urbanization. These results also highlight the continuity of the MTBC population in China over the past 1,000 years.

L2.3 dominates recent population dynamics. It is notable that L2.3 emerged relatively recently; that is, ~450 years after the first expansion of L2.2 and ~300 years later than the three indigenous L4 sublineages. However, L2.3 is likely to have swept rapidly through the population (Fig. 4c). Our estimates of the population growth rate per year (based on changes in the median N_e) imply that L2.3 may have expanded 1.4~3.5 times faster than the other sublineages (Table 1). It is also interesting that the growth of the other sublineages slowed or ceased as L2.3 started to expand

(Fig. 4a,c). We further compared 15-locus mycobacterial interspersed repetitive-unit (MIRU)-variable-number tandem-repeat (VNTR) cluster rates (an indicator of recent transmission) of sublineages in 6 geographically distinct populations (population sites in Fig. 1b). Our data showed that isolates belonging to L2.3 were more likely to be clustered compared with the other sublineages (odds ratio = 3.7, 95% confidence interval (CI): 2.9–4.8), consistent with higher rates of transmission (Table 2). This finding remained statistically significant when we repeated the comparison using different sets of VNTR loci to define clusters (Supplementary Table 5). These analyses suggest that relatively high transmission rates have contributed to L2.3's dominance of modern MTBC populations. We also observed bias in the frequencies of L2.2 and L2.3 outside China (Fig. 5a). Among the L2 strains sampled in other non-Southeast Asian countries, 94.3% were caused by L2.3 while L2.2 accounted for only 5.7% (chi-squared test, $P < 0.0001$). This result suggests that the emergence of L2 strains across the globe has been driven primarily by L2.3.

Table 2 | VNTR cluster rates of the indigenous sublineages in six distinct populations

Sublineage	Individual county site ^a						All sites	Odds ratio	Pvalue
	Guangxi	Heilongjiang	Henan	Sichuan	Shandong	Shanghai			
L2.3	41% (24/58)	66% (95/145)	57% (70/123)	23% (15/65)	61% (76/125)	58% (132/226)	56% (412/742)	1	-
L2.2	31% (15/49)	42% (5/12)	40% (17/42)	16% (7/45)	25% (6/24)	36% (29/80)	31% (79/252)	0.366 (0.267–0.500)	<0.0001
L4.2	13% (2/16)	0% (0/0)	22% (2/9)	24% (4/17)	0% (0/5)	36% (8/22)	23% (16/69)	0.242 (0.127–0.440)	<0.0001
L4.4	20% (4/20)	60% (6/10)	0% (0/3)	0% (0/23)	27% (4/15)	0% (0/17)	16% (14/88)	0.152 (0.078–0.277)	<0.0001
L4.5	40% (6/15)	13% (2/15)	0% (0/8)	23% (10/44)	26% (5/19)	8% (2/25)	20% (25/126)	0.198 (0.120–0.319)	<0.0001

^aThe six county sites refer to the population sites from which the MTBC isolates were collected using the population-based approach from 2009 to 2010.

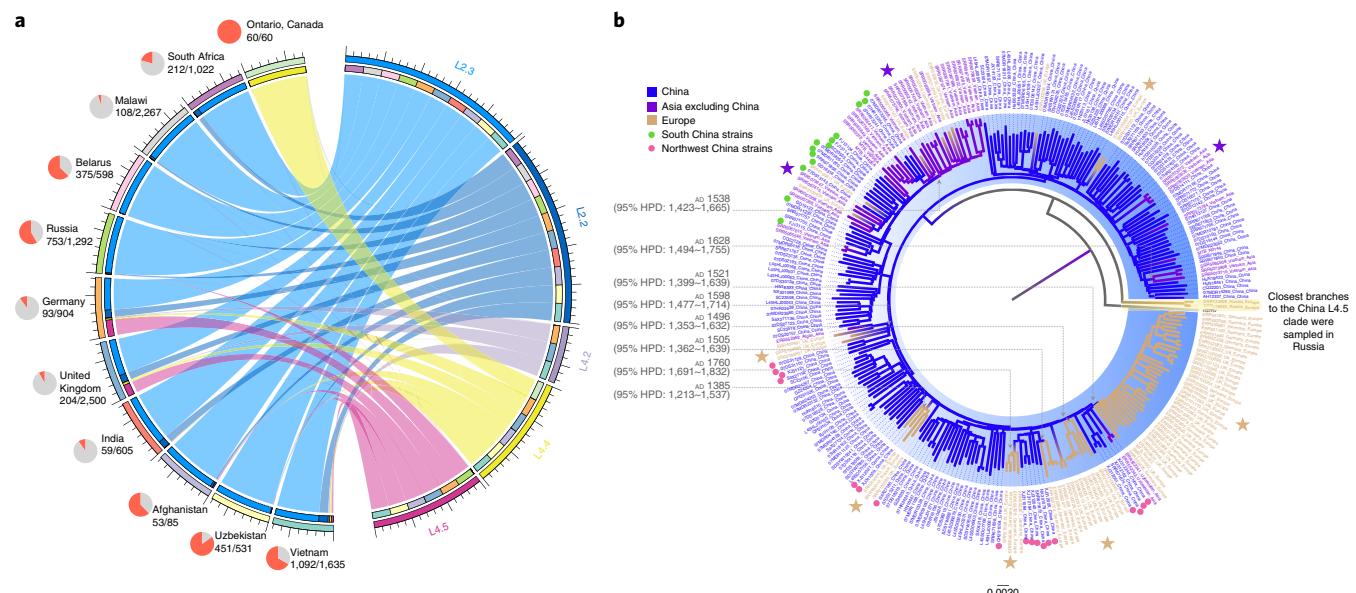


Fig. 5 | Global dispersal of Chinese indigenous genotypes. **a**, Circle plot with ribbons depicting the dispersal flows that led to the global emergence of Chinese indigenous sublineages. All the flows refer to 'one-way' outflows and indicate direct or indirect exportation events, with the ribbon width at each end proportional to the number of strains sampled in each country. The pie charts next to the country names show the proportion (red sector) of strains in the relative dataset that were found to be descendants of Chinese indigenous genotypes. **b**, Global dispersal of Chinese L4.5 strains. Different leaf colours indicate the diverse geographic origins of those isolates. The Chinese L4.5 clade is highlighted in blue, and the branches are coloured according to their geographic attributions. The major country transition events are marked with stars. The transition time of each event was estimated using the MTBC-6 model. The European-specific clades are nested within the Chinese L4.5 clade, with the closest branches sampled from Northwest China. The closest branches to the strains sampled from Vietnam were mostly collected in South China.

Internal and global dispersal routes of indigenous genotypes. To explore the possible dispersal routes of the indigenous MTBC genotypes in China, we created contour maps showing the relative concentrations of each sublineage across the country (Fig. 6). We identified hotspots for L2.2 in the southwest and northwest, suggesting that this sublineage expanded out of two distinct geographic foci (Fig. 6a). L2.3 was concentrated in the northeast, with a single hotspot around Beijing and a gradual decrease in prevalence as a function of distance from the city (Fig. 6b). Beijing has been the capital city of China since the Yuan dynasty (AD 1271–1368). Since then, the human population in Beijing has doubled and continues

to grow rapidly³⁶. As the political and trade centre in China, Beijing was the central hub of population flow and interchange of resources between the thirteenth and twentieth centuries³⁷. L2.3's emergence in the context of Beijing's expanding and mobile human population may have contributed to its success.

All three indigenous sublineages of L4 are concentrated in Southern China (Fig. 6c–e), suggesting a role of migrations to this region in the original establishment of these sublineages. In Central China, L4.2 and L4.4 are concentrated in Sichuan province, which is consistent with the local history in that the inhabitants of Sichuan are mostly immigrants from Guangxi, Guangdong and other

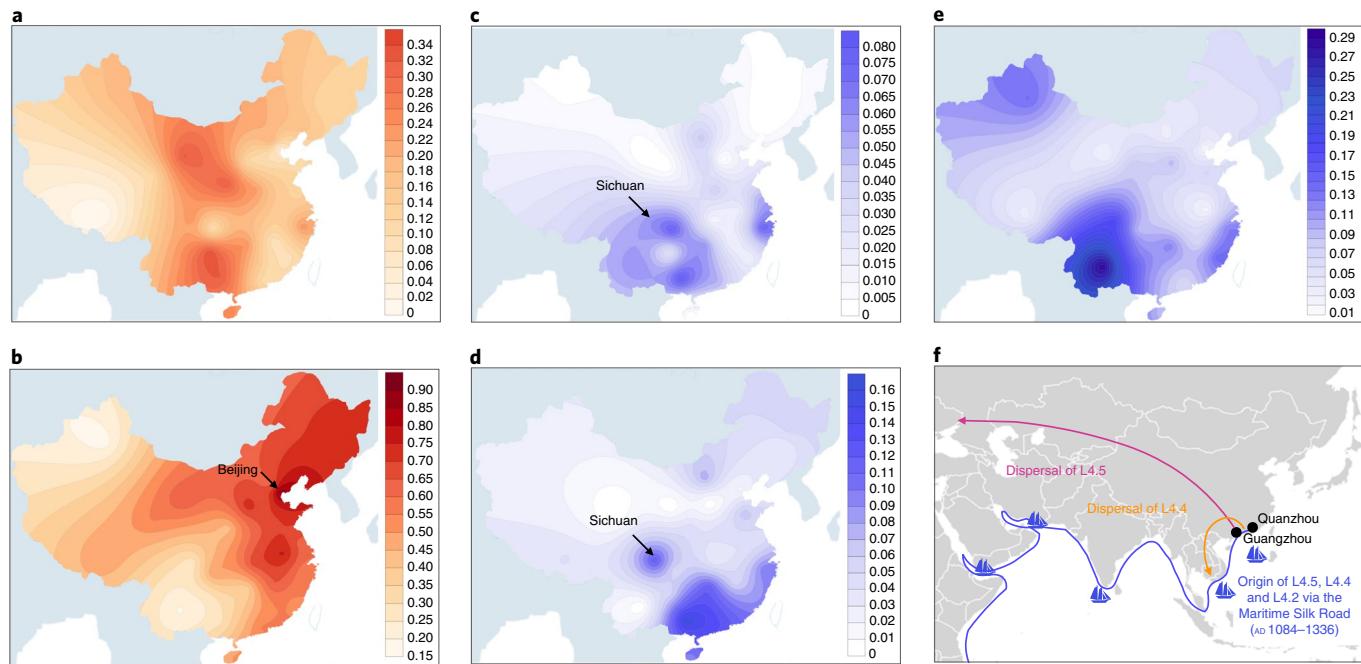


Fig. 6 | Contour maps showing the countrywide prevalence of indigenous sublineages. a–e, Colour ranges showing the percentage prevalence of L2.2 (a), L2.3 (b), L4.2 (c), L4.4 (d) and L4.5 (e) based on the SNP typing data. **f,** A scenario of Maritime Silk Road origins for L4 sublineages in China. The ship symbols mark the major ports on historical sea trade routes. The directions of dispersal of L4.5 and L4.4 are shown.

southern provinces where L4.2 and L4.4 are most prevalent (almost one million migrants moved to Sichuan following the great massacre in the late Ming dynasty that almost emptied that province)³⁸. Taken together, these results highlight the importance of human historical phenomena in shaping the genetic diversity of MTBC within China.

We also observed historical and recent exports of Chinese L4 strains. L4.5 strains were sampled in the United Kingdom and Germany at relatively high frequencies (Fig. 5a), and those strains formed European-specific clades nested within the Chinese L4.5 clade (Fig. 5b). Their deeply rooted positions suggest that they were exported in ancient times. Intriguingly, we noticed a consistent pattern in that the early-diverging branches closest to European L4.5 clades were sampled exclusively in Northwest China (especially in Xinjiang province). This pattern could result from the dispersal of L4.5 to Europe through the ancient overland Silk Road, on which Xinjiang acted as one of the crossroads of Eurasia³⁹. A more recent example is a strain that caused an outbreak in Ontario, Canada in the 1990s, which turned out to be a descendant of the Chinese L4.4 clade (Supplementary Fig. 8). The origin of the relative ancestor was estimated to be around AD 1973 and the closest branches to that strain were sampled in Guangdong province (Southern China), suggesting that this strain hitchhiked the Chinese migrant waves from Southern China to North America in the second half of the past century^{40–42}.

A scenario of Maritime Silk Road origin. We inferred a European origin for the ancestors of L4.2, L4.4 and L4.5 (Supplementary Fig. 5), consistent with the designation of L4 as a ‘Euro-American’ lineage⁴³. In our reconstructions of the individual ancestral states for L4.2, L4.4 and L4.5, we found the three China-specific clades within those sublineages to have the highest probability of an origin in Southern China (Supplementary Figs. 10 and 11). The Chinese clades of L4 could have been introduced from Europe to South China via two major routes. The overland route traverses the Middle East, South Asia and Southeast Asia⁴⁴. Alternatively, MTBC

could have migrated by sea; for example, along the well-known Maritime Silk Road^{45,46}. In the former scenario, we would expect the L4 sublineages to be prevalent in the countries along the migration route to China, and the corresponding L4 sublineages there should form early-diverged branches to China’s L4 strains. This is the case for L3, which is hypothesized to have been dispersed via the overland Silk Road⁴⁶ and is concentrated in Northwest China (Fig. 1a). In contrast, the three indigenous L4 sublineages (L4.2, L4.4 and L4.5) were rarely detected in a large sample from countries including India, Afghanistan and Uzbekistan (Supplementary Table 6). Importantly, isolates from these countries were nested among samples from China, suggesting that L4.2, L4.4 and L4.5 diversified in China before spreading to other countries in Southern and Central Asia (Supplementary Fig. 5). The emergence of the three L4 sublineages (AD 1160–1268) coincides with the period of peak activity for the Maritime Silk Road in China, during the Song dynasty (AD 960–1279). Trading in this period was intense and extensive between multiple ports in Europe and South China, such as Guangzhou and Quanzhou^{46,47} (Fig. 6f). The Southern Chinese origin and timing of the origin of the indigenous L4 lineages are consistent with an emergence in the context of the Silk Road era maritime trade between Europe and China (Fig. 6f).

Discussion

Through reconstruction of the evolutionary history of MTBC strains circulating in China, we demonstrate that the current tuberculosis epidemic stems almost entirely from historical human migration events that established L2 and three sublineages of L4 starting around 1,000 years before the present. Our data show that these strains underwent massive expansion over the past 1,000 years. This is remarkable in that there is no historical record of a severe tuberculosis outbreak in China^{1,12}, unlike the situation in Europe where the epidemic of eighteenth to nineteenth centuries constitutes a clearly defined and documented epidemiological phenomenon with marked cultural and demographic impacts¹⁰. A potential explanation for this disparity is differences in the timing and pace of

industrialization in China and Europe. The devastating tuberculosis epidemic in Europe is believed to have been triggered by transitions in social conditions such as overcrowding and malnutrition linked to the industrial revolution¹³. The environment in China differed markedly in this period, as China did not participate in the industrial revolution⁴⁸ and the historical urbanization rate rarely exceeded 10% of the total population⁴⁹. Historical tuberculosis expansions in China could reflect cryptic but frequent transmission, as reflected by the numerous descriptions of ‘Lao-bing’ (tuberculosis) in historical medical texts¹². Our results indicate there was a period of four centuries during which the four indigenous Chinese sublineages expanded simultaneously, suggesting that growth of the pathogen population was driven by common ecological transitions such as growth of the host population and an increase in urbanization (the urbanization rate increased from 5% in the Tang dynasty (AD 618–907) to 10~13% in the Song dynasty (AD 960–1279))⁵⁰.

Why so little genetic diversity? A simple explanation for the limited number of introductions of MTBC to China is its distance from the African continent, where MTBC appears to have first emerged^{16,51}. However, there is evidence to suggest that MTBC dispersed readily between Southeast Asia and the African continent¹⁶, which argues against a simple model in which migration events scale with distance. Historically isolationist policies may be of greater relevance: for example, during the Ming and Qing dynasties (AD 1368–1912), feudal rulers adopted a policy of seclusion that hampered exchanges with the outside world^{51,52}. Interestingly, we did not identify any indigenous sublineages that were introduced during that time interval. Alternatively, past genetic diversity may have been higher than it is at present if, for example, multiple lineages were displaced by a sweep of L2.3, or if they failed to survive the bacterial population reduction that followed the widespread implementation of anti-tuberculosis therapy. The newly discovered sublineage L4.11 could support this notion, as its wide geographical distribution at low frequency suggests that the population has been through expansion and contraction.

Propagation of the tuberculosis epidemic in China. Tuberculosis in China continues to be characterized by substantial ongoing transmission, with a reported recent transmission rate of ~30%⁵³. Our analyses suggest that the current tuberculosis situation in China represents the waning era of an epidemic that expanded over the past millennium. As discussed above, historically isolationist policies may have limited the incursion of foreign MTBC in China, thereby contributing to the minimal genetic diversity observed here. Our results also suggest that historical patterns of mixing were uneven, and enabled a ‘winner takes all’ dynamic in which a small subset of the MTBC strains introduced to China came to dominate the population. The contrast between the indigenous genotypes and L3 or L1 is illustrative: L3 and L1 remain at low frequencies and are geographically restricted to the northwest (Xinjiang province) and southeast (Taiwan), respectively. We hypothesize that the disparate fates of these lineages and sublineages reflect historical variation in mobility and/or growth among subpopulations of hosts.

Our analyses here suggest that L2.3 spread out of Beijing, and we posit that the dominance of this sublineage reflects the centralization of the city within China over the past 700 years. Specific historical phenomena that are likely to have contributed to L2.3’s dominance include growth and urbanization of the Beijing population and the city’s role as a hub for migration within China³⁶. Phenomena that are extrinsic to the pathogen, such as crowding and host malnutrition, can be powerful drivers of tuberculosis transmission and can lead to the dominance of specific MTBC strains that land under favourable conditions^{17,54}. Our results also suggest that L2.3 strains spread more rapidly than strains of the other extant sublineages. Hence, the rapid expansion of the L2.3 population may also have

been facilitated by relatively high rates of disease progression following infection and/or larger numbers of secondary cases per source case^{55,56}. Recent studies have shown increased virulence of ‘modern’ Beijing (L2.3) strains in mouse infection models and induction of lower levels of proinflammatory cytokines than ‘ancient’ Beijing (L2.2) strains^{57–59}. More recently, a study showed that mutations of *ppe38* in L2.3 strains could completely block the secretion of two large subsets of ESX-5 substrates and lead to a hypervirulent phenotype⁶⁰. Enhanced virulence could plausibly lead to more rapid disease progression following infection, and consequently high rates of spread. We posit that the success of the L2.3 sublineage results both from factors that are intrinsic to the bacteria (for example, increased virulence leading to rapid onset of transmissible disease) as well as extrinsic conditions, such as an expanding and mobile host population. It is possible that these phenomena interacted to make L2.3 a dominant sublineage. For example, the large size and high density of the human population may have better sustained ongoing transmission of rapidly progressive forms of tuberculosis than low-density communities of hosts.

Genetic continuity of the MTBC population. A recent study of MTBC genomes from eighteenth-century Hungary found them to be nested within a phylogeny of contemporary strains, which points to continuity of MTBC lineages over the past two centuries in Europe⁶¹. We inferred that MTBC strains currently circulating in China trace to introductions that occurred around 1,000 years before the present. These findings demonstrate that the MTBC lineages that become established in favourable environments can persist for centuries. We hypothesize that the capacity of MTBC to establish long-term latent infections contributes to this continuity, which contrasts with bacterial pathogens such as *Salmonella Typhi*, *Vibrio cholerae* and *Yersinia pestis*, for which massive lineage replacements have been observed^{62–65} (see additional discussion in the Supplementary Information).

In conclusion, we demonstrate that China’s tuberculosis epidemic is a historical heritage that stems from just a handful of introductions of its causative agent, *M. tuberculosis*. These introduced strains expanded in parallel, presumably in response to common ecological drivers of epidemic tuberculosis. The long-term genetic continuity and distinctive structure of local MTBC populations in this high-incidence region highlight the potential for MTBC population dynamics to guide the epidemiological surveillance of tuberculosis.

Methods

Countrywide sampling of MTBC isolates. The MTBC isolates analysed in this study consisted of two sample sets. The first was a whole-country dataset (70 random sites). In 2007, a total of 3,929 culture-positive MTBC isolates were collected from 70 counties for the purpose of surveillance for drug resistance. These counties covered 31 out of the 34 provincial regions of China (Supplementary Table 2)⁶⁶. The number of isolates collected from each province was proportional to the number of smear-positive cases reported in that province relative to the total number of cases nationwide. These isolates were recovered on Lowenstein-Jensen medium from stored MTBC samples that were previously kept in a –80 °C freezer. However, 702 isolates failed in recovery, possibly due to multiple freeze and thaw cycles. The second was a whole-population dataset (six population sites). From 1 June 2009 to 31 December 2010, a total of 1,375 MTBC isolates were collected at six county sites from six different provinces in China in a population-based molecular epidemiological study. Of these, 1,351 were available for genotyping in this study (Supplementary Table 2)⁵⁵. These 6 county sites represented different geographic settings spreading south to north and west to east, and covered a total population of about 5.8 million inhabitants. All the county site locations were marked on maps using Tableau version 10.4 (<https://www.tableau.com>).

DNA extraction and SNP typing. Genomic DNA of MTBC isolates was extracted using the boiling method for lineage and sublineage genotyping⁶⁷. As MTBC does not appear to engage in lateral gene transfer and homoplastic SNPs are rare^{23,68}, SNPs are ideal markers for typing MTBC sublineages⁶⁹. L2 strains were classified into L2.1 (proto-Beijing), L2.2 ('ancient' Beijing or atypical Beijing) and L2.3 ('modern' Beijing

or typical Beijing) sublineages^{30,70,71}, while L4 strains were classified according to the previously defined ten sublineages^{21,30}. We developed six real-time PCR melting curve assays for SNP typing using the well-characterized sublineage-specific SNP markers (Supplementary Table 7). The principle of this SNP typing assay is similar to the drug-resistant mutation detection assay we developed previously⁶⁷. Briefly, one dually labelled probe was designed for each sublineage-specific phylogenetic SNP. For each single-stranded nucleotide probe, one end was labelled with a fluorophore (fluorescein or rhodamine X) and the other end was labelled with a quencher. As the targeted sublineage strains differed by one nucleotide in the detecting region compared with other strains, the melting curve analysis would show different melting temperature values due to the altered annealing ability⁶⁷. The targeted sublineage was thereby differentiated from the remaining sublineages. One probe was designed to differentiate L2 strains from the remaining isolates. For L2 sublineage typing, two probes were designed to detect L2.2 and L2.3, while the remaining isolates were sequenced by WGS to detect L2.1. For L4 sublineage typing, three probes were designed to detect the common L4 sublineages (L4.5, L4.4 and L4.2) based on our pilot typing results, while the remaining L4 isolates and those showing ambiguous typing results were further sequenced by WGS. A total of 32 isolates that could not be assigned to any known sublineage or that showed ambiguous typing results underwent WGS. All the real-time PCR melting curve analysis assays were performed on a Bio-Rad CFX96 platform.

Public data collection. *Sporolotyping data collection.* As spoligotyping results are generally concordant with lineage classifications based on WGS data⁷², we first collected spoligotyping data from MTBC isolates sampled throughout China to obtain an initial picture of the MTBC population structure in the country. We searched for research articles that published spoligotyping data of MTBC isolates collected from China on PubMed. We identified a total of 96 articles in a preliminary search, a large proportion of which did not provide detailed spoligotyping results or reported previously published data. A total of 26 articles provided valid spoligotyping results for 16,621 MTBC isolates with either original typing records or summarized typing results (Supplementary Table 1). Each isolate was assigned to the relevant MTBC lineage according to previously identified links between spoligotypes and phylogenetic lineages⁷. A total of 15,217 isolates could be successfully assigned to known lineages. Among the remaining isolates in the sample, 1,030 (6.2%) were classified as 'orphan', 371 (2.2%) as 'Manu2' and 3 (< 0.01%) as '*Mycobacterium bovis*'. Orphan strains could be variant types from any lineage that had not been recorded in the SpolDB4.0 database⁷³, whereas Manu2 types could be DNA samples that were a mixture of L2 and L4 strains⁷⁴. These two types were excluded from our subsequent analyses. The number of MTBC isolates from each province was normalized to the relative human demographic data when calculating the total prevalence of each MTBC lineage (Supplementary Table 1).

WGS data collection. To identify WGS data from global MTBC isolates, we searched for articles with WGS data published in PubMed and downloaded the original sequencing reads from the European Nucleotide Archive (EMBL-EBI). The geographic origin and year of collection for each isolate were extracted from the relevant article. We sent an enquiry to the authors of papers that did not include this information. Sequencing data were downloaded for 15,047 MTBC isolates and we obtained geographic information for 12,596 of them (Supplementary Table 1).

WGS and SNP calling. A minimum spanning tree was constructed based on 15-locus MIRU-VNTR data for all the MTBC isolates genotyped here. We selected MTBC isolates from each of the clades to represent countrywide genetic diversity (a Perl script was written to sample isolates from each clade randomly). We purposefully did not select L2 isolates because previous studies have generated an abundance of WGS results for this lineage from China^{31,70,75}. Genomic DNA was extracted from the 306 MTBC isolates (three of L1, 23 of L3 and 280 of L4) following the cetyltrimethylammonium bromide lysozyme method as described before⁷⁰. A 300-base pair fragment length library was constructed for each DNA sample, and WGS was performed on an Illumina HiSeq 2500 system with either the single-end or paired-end strategy. We used a previously validated pipeline for the mapping of short sequencing reads to the reference genome⁷⁵. In brief, the Sickle⁷⁶ tool was used to trim WGS data. Sequencing reads with a Phred base quality above 20 and read length longer than 30 were kept for analysis. The whole-genome sequence of the *M. tuberculosis* H37Rv strain (NC_000962.2) was used as the reference template for read mapping. Sequencing reads were mapped to the reference genome using Bowtie 2 (version 2.2.9)⁷⁷. SAMtools (version 1.3.1)⁷⁸ was used for SNP calling with a mapping quality greater than 30. Fixed mutations (with a frequency of $\geq 75\%$) were identified using VarScan (version 2.3.9)⁷⁹ with at least 10 reads supporting and the strand bias filter option on. We excluded all SNPs that were located in repetitive regions of the genome (for example, PPE/PE-PGRS family genes, phage sequence, insertion or mobile genetic elements) that are difficult to characterize with short-read sequencing technologies. Small insertions or deletions identified by VarScan (version 2.3.9) were also excluded.

Phylogenetic reconstruction. Mixed-infection isolates were excluded for phylogenetic reconstruction by investigating the genotype heterozygosity of SNPs

as described previously⁸⁰. For all phylogenetic reconstruction, SNPs of MTBC isolates were combined into a single consensus and non-redundant list while those nucleotide positions with gaps in more than 5% of the taxa were excluded (possibly due to insertions or deletions, low coverage or mapping quality at those sites). The alignments of polymorphic positions from all strains were used for phylogeny reconstruction using MEGA 6.0 (ref. ⁸¹). The neighbour-joining method was used for initial inference of the phylogeny structure when the taxa numbers were large. However, for the final estimation of phylogenies, the maximum likelihood method was applied under the general time reverse model with at least 100 replicates for bootstrapping confidence levels. Phylogeny trees were visualized in FigTree (version 1.4.3) (<http://tree.bio.ed.ac.uk/software/figtree/>) or iTOL (version 3)⁸². We adapted a recently described hierarchical nomenclature to define nodes and subclades within the tree in the definition of sublineages³⁰.

Population genetic analyses. *Pairwise SNP distance.* We wrote a Perl script to calculate the number of pairwise SNP distances. For those countries with a large number of MTBC isolates, we randomly subsampled to 200 isolates in the calculation of pairwise SNP distances so as to be able to compare them with other countries with samples closer to 100 isolates. The distribution and mean pairwise SNP distance for each country was plotted with ggplot2 in RStudio (version 3.4.0)⁸³. As the pairwise SNP distances from the country samples were not normally distributed, a Wilcoxon rank-sum test was used to test the differences between countries.

Nucleotide diversity. Average pairwise π values per site were calculated using the nuc.div function in the 'pegas' library⁸⁴. As sample sizes varied among countries, we generated 100 subsamples, each equal to 200 isolates, and thereby calculated the CI for those countries with samples greater than 200 MTBC isolates. The random subsampling process was completed using the rand command in Perl. We used all isolates for analysis of the data from countries with 50–200 isolates. We did not include data from countries with fewer than 50 isolates available.

Rarefaction analysis. Rarefaction is an ecological technique designed to estimate the species richness expected for a given number of individual samples, which is based on the construction of rarefaction curves. We used this method to evaluate the sublineage diversity of MTBC isolates from different countries. Rarefaction analysis was performed in RStudio (version 3.4.0) using the library 'iNEXT'⁸⁵.

Phylogeographic analysis. We used RASP⁸⁶, which implements both Bayesian and parsimony (S-DIVA) approaches, to estimate the ancestral geographic ranges of L2 and the three major L4 sublineages in China. To estimate the geographic origins of those MTBC sublineages, we divided the world map into five broad geographic areas and used them as a proxy for the most likely origin of each strain (Supplementary Fig. 1). The reason we treated China independent from Asia was because we wanted to test whether the sublineages originated and diversified locally as opposed to being imported from other countries in Asia. The maximum likelihood phylogenies of both lineages and the corresponding geographic regions of origin for those isolates were loaded as a distribution. RASP reconstruction was performed without the outgroup. For the parsimony-based analyses, a maximum of two ancestral areas per node were allowed for range reconstruction. For the Bayesian-based analyses, 5 different chains during 500,000 generations were run. For the contour plots in Fig. 6, the prevalence of each sublineage was determined by the MTBC isolates we genotyped in each province. Surfer (version 12) software (<http://www.goldensoftware.com/>) was used to contour plot each sublineage's frequency throughout the county.

Bayesian-based coalescent analysis. *Dating analysis.* We selected 96 MTBC strains from a published study to represent the global diversity of MTBC lineages³¹. These 96 strains, together with the 23 L2 strains and 41 L4 strains (10 L4.2 strains, 12 L4.4 strains, 11 L4.5 strains and 8 L4.11 strains) sampled from China, were used for phylogenetic reconstruction. A total of 24,792 concatenated genome-wide variable positions were used for phylogenetic analyses. We estimated the dates of the most common recent ancestors of L2 and L4 and their sublineages using BEAST (version 1.8.0)⁸⁷. The XML input file was modified to specify the number of invariant sites in the MTBC genomes. For the MTBC genome substitution rate, we imposed a normal distribution for the substitution rate of MTBC with a mean of 4.6×10^{-8} substitutions per genome per site per year (95% highest posterior density (HPD) interval: 3.0×10^{-8} to 6.2×10^{-8}) as in a previous study¹. We used an uncorrelated log-normal distribution for the substitution rate and a constant population size for the tree priors. We ran three chains of 5×10^7 generations sampled every 10,000 to ensure independent convergence of the chains, the first 10% of which were discarded as burn-in. Convergence was assessed using Tracer (version 1.6.0), ensuring that all relevant parameters reached an effective sample size of >100 . Phylogenetic trees were visualized using FigTree (version 1.4.3).

Bayesian Skyline Plot (BSP) analysis. A BSP analysis was applied to estimate the past effective population size dynamics of the L2 and L4 sublineages based on the substitution rate model indicated above. Ideally, one should use all the samples for BSP analysis, but this was not computationally feasible. In addition, the number

of isolates that were sequenced by WGS for each sublineage was not proportional to their epidemiological prevalence, as some sublineages (L4.2, L4.4 and L4.5) were oversampled while others (L2.2 and L2.3) were undersampled. To reduce this bias and improve the tractability of our analyses, we subsampled isolates from the original collection and constructed a WGS dataset containing 500 MTBC isolates. The isolates from different sublineages in this dataset were randomly sampled from the original sequenced genomes and the number was proportional to their prevalence in China. Sublineage-based skyline analysis was performed, and the ages of the most recent common ancestors from dating analysis were used as the tree heights. For the MTBC effective population size curve in Fig. 4b, we integrated all the above isolates for the BSP analysis. In each case, 3 chains of 5×10^7 generations were sampled every 10,000 to assure independent convergence of the chains. For the past human population size changes, the demographic data were obtained from historical investigations and records³². For the population dynamics plot in Fig. 4c, the relative prevalence of each sublineage in the past was estimated from the effective population growth curves generated by BSP analysis (each sublineage was estimated separately) and plotted using the 'streamgraph' package in RStudio (version 3.4.0).

Population growth rate estimation. The population growth rate per year was calculated using the effective population growth curves generated from BSP analysis³³. Each skyline plot consisted of 100 smoothed data points, at 5–11 yr intervals. For L2.2 and L4.5, the effective population size increase was preceded by a brief period of stationary size. The initial population size N_0 was set as the minimum population size during the period immediately preceding population growth. The population size became stationary or even decreased at later stages. Thus, we estimated the effective population growth rate for the increasing interval in our data³³. The exponential growth equation was chosen for this analysis:

$$r = \ln[N_t/N_0]/t$$

In this equation, r represents the population growth rate per year, N_0 is the initial population size and t is the duration of time since growth began.

Ethics. The analysis of anonymous MTBC isolates collected in China in this study was approved by the Tuberculosis Research Ethics Review Committee of the Chinese Center for Disease Control and Prevention and the Institutional Review Board of the Institutes of Biomedical Sciences, Fudan University.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

Sequencing reads have been submitted to the European Nucleotide Archive (EMBL-EBI) under study accession PRJEB23157. The geographic information for individual isolates is listed in Supplementary Table 3. The analysis scripts used in this study are available online at GitHub (https://github.com/StopTB/China_TB_Evolutionary_History).

Received: 24 February 2018; Accepted: 28 August 2018;

Published online: 5 November 2018

References

- Bos, K. I. et al. Pre-Columbian mycobacterial genomes reveal seals as a source of New World human tuberculosis. *Nature* **514**, 494–497 (2014).
- Global Tuberculosis Report 2017* (World Health Organization, 2017).
- Narain, J. P., Ravaglione, M. C. & Kochi, A. HIV-associated tuberculosis in developing countries: epidemiology and strategies for prevention. *Tuber. Lung Dis.* **73**, 311–321 (1992).
- Steffen, R., Rickenbach, M., Wilhelm, U., Helminger, A. & Schar, M. Health problems after travel to developing countries. *J. Infect. Dis.* **156**, 84–91 (1987).
- Fusegawa, H. et al. Outbreak of tuberculosis in a 2000-year-old Chinese population. *Kansenshogaku Zasshi* **77**, 146–149 (2003).
- Prasad, P. V. General medicine in Atharvaveda with special reference to Yaksma (consumption/tuberculosis). *Bull. Indian Inst. Hist. Med. Hyderabad* **32**, 1–14 (2002).
- Suzuki, T. & Inoue, T. Earliest evidence of spinal tuberculosis from the Neolithic Yayoi period in Japan. *Int. J. Osteoarchaeol.* **17**, 392–402 (2007).
- Li, X. et al. Archaeological and palaeopathological study on the third/second century BC grave from Turfan, China: individual health history and regional implications. *Quat. Int.* **290**, 335–343 (2013).
- Packard, R. M. *White Plague, Black Labor: Tuberculosis and the Political Economy of Health and Disease in South Africa* (Univ. California Press, Berkeley, 1989).
- Dubos, R. J. & Dubos, J. *The White Plague: Tuberculosis, Man, and Society* (Rutgers Univ. Press, New Brunswick, 1952).
- Stead, W. W. The origin and erratic global spread of tuberculosis. How the past explains the present and is the key to the future. *Clin. Chest Med.* **18**, 65–77 (1997).
- Zhang, Z. *Epidemic Chronology of Ancient China* [in Chinese] (Fujian Science and Technology Press, Fuzhou, 2007).
- Bates, J. H. & Stead, W. W. The history of tuberculosis as a global epidemic. *Med. Clin. North Am.* **77**, 1205–1217 (1993).
- Perry, E. J. & Selden, M. *Chinese Society: Change, Conflict and Resistance* (Routledge, London, 2003).
- Wang, F. & Zuo, X. Inside China's cities: institutional barriers and opportunities for urban migrants. *Am. Econ. Rev.* **89**, 276–280 (1999).
- O'Neill, M. B. et al. Lineage specific histories of *Mycobacterium tuberculosis* dispersal in Africa and Eurasia. Preprint at <https://www.biorxiv.org/content/early/2017/10/27/2010161> (2017).
- Pepperell, C. S. et al. Dispersal of *Mycobacterium tuberculosis* via the Canadian fur trade. *Proc. Natl Acad. Sci. USA* **108**, 6526–6531 (2011).
- Hershberg, R. et al. High functional diversity in *Mycobacterium tuberculosis* driven by genetic drift and human demography. *PLoS Biol.* **6**, e311 (2008).
- Wirth, T. et al. Origin, spread and demography of the *Mycobacterium tuberculosis* complex. *PLoS Pathog.* **4**, e1000160 (2008).
- Gagneux, S. & Small, P. M. Global phylogeography of *Mycobacterium tuberculosis* and implications for tuberculosis product development. *Lancet Infect. Dis.* **7**, 328–337 (2007).
- Stucki, D. et al. *Mycobacterium tuberculosis* lineage 4 comprises globally distributed and geographically restricted sublineages. *Nat. Genet.* **48**, 1535–1543 (2016).
- Linz, B. et al. An African origin for the intimate association between humans and *Helicobacter pylori*. *Nature* **445**, 915–918 (2007).
- Pepperell, C. S. et al. The role of selection in shaping diversity of natural *M. tuberculosis* populations. *PLoS Pathog.* **9**, e1003543 (2013).
- Bjorn-Mortensen, K. et al. Tracing *Mycobacterium tuberculosis* transmission by whole genome sequencing in a high incidence setting: a retrospective population-based study in East Greenland. *Sci. Rep.* **6**, 33180 (2016).
- Lee, R. S. et al. Population genomics of *Mycobacterium tuberculosis* in the Inuit. *Proc. Natl Acad. Sci. USA* **112**, 13609–13614 (2015).
- 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.* **25**, 3260–3266 (2015).
- Holt, K. E. et al. Frequent transmission of the *Mycobacterium tuberculosis* Beijing lineage and positive selection for the EsxW Beijing variant in Vietnam. *Nat. Genet.* **50**, 849–856 (2018).
- Van Soelingen, D. et al. Predominance of a single genotype of *Mycobacterium tuberculosis* in countries of east Asia. *J. Clin. Microbiol.* **33**, 3234–3238 (1995).
- Pang, Y. et al. Spoligotyping and drug resistance analysis of *Mycobacterium tuberculosis* strains from national survey in China. *PLoS ONE* **7**, e32976 (2012).
- Coll, F. et al. A robust SNP barcode for typing *Mycobacterium tuberculosis* complex strains. *Nat. Commun.* **5**, 4812 (2014).
- Comas, I. et al. Out-of-Africa migration and Neolithic coexpansion of *Mycobacterium tuberculosis* with modern humans. *Nat. Genet.* **45**, 1176–1182 (2013).
- Ge, J. *China Population History (Zhongguo Renkou Shi)* (Fudan Univ. Press, Shanghai, 2000).
- Wang, L. et al. Tuberculosis prevalence in China, 1990–2010; a longitudinal analysis of national survey data. *Lancet* **383**, 2057–2064 (2014).
- Neher, R. A. & Hallatschek, O. Genealogies of rapidly adapting populations. *Proc. Natl Acad. Sci. USA* **110**, 437–442 (2013).
- Magiorkinis, G. et al. Integrating phylodynamics and epidemiology to estimate transmission diversity in viral epidemics. *PLoS Comput. Biol.* **9**, e1002876 (2013).
- Guang-Hui, H. *Historical Population Geography of Beijing* [in Chinese] (Peking Univ. Press, Beijing, 1996).
- Hou Ren-Zhi, T. X.-F. *Historical Geography of Beijing City* [in Chinese] (Beijing Yanshan Press, Beijing, 2000).
- Huang, Q.-S. & Yang, G.-H. The placename of immigration in Sichuan and Huguang people migrate into Sichuan. *J. Southwest China Normal Univ.* **3**, 023 (2005).
- Millward, J. A. *Eurasian Crossroads: A History of Xinjiang* (Columbia Univ. Press, New York, 2007).
- Poston, D. L. Jr., Mao, M. X. & Yu, M.-Y. The global distribution of the overseas Chinese around 1990. *Popul. Dev. Rev.* **20**, 631–645 (1994).
- Li, P. S. The rise and fall of Chinese immigration to Canada: newcomers from Hong Kong special administrative region of China and mainland China, 1980–2000. *Int. Migr.* **43**, 9–34 (2005).
- King, H. & Locke, F. B. Chinese in the United States: a century of occupational transition. *Int. Migr. Rev.* **14**, 15–42 (1980).
- Gagneux, S. et al. Variable host-pathogen compatibility in *Mycobacterium tuberculosis*. *Proc. Natl Acad. Sci. USA* **103**, 2869–2873 (2006).

44. McNeill, W. H. Human migration in historical perspective. *Popul. Dev. Rev.* **10**, 1–18 (1984).
45. Gan, F. *Ancient Glass Research Along the Silk Road* (World Scientific, Hackensack, 2009).
46. Kauz, R. *Aspects of the Maritime Silk Road: From the Persian Gulf to the East China Sea* (Harrassowitz, Wiesbaden, 2010).
47. McPherson, K. China and the Maritime Silk Route. In *Proc. of the UNESCO Quanzhou International Seminar on China and the Maritime Routes of the Silk Roads* 55–60 (People's Publishing House, 1991).
48. Lin, J. Y. The Needham puzzle: why the industrial revolution did not originate in China. *Econ. Dev. Cult. Change* **43**, 269–292 (1995).
49. Jones, E. L., Frost, L. & White, C. *Coming Full Circle: An Economic History of the Pacific Rim* (Westview Press, Colorado, 1993).
50. Yusuf, S. & Saich, A. *China Urbanizes: Consequences, Strategies, and Policies* (World Bank, Washington DC, 2008).
51. Millward, J., Dunnell, R. W., Elliott, M. C. & Forêt, P. *New Qing Imperial History. Making of Inner Asian Empire at Qing Chengde* (RoutledgeCurzon, New York, 2004).
52. Mote, F. W., Twitchett, D. & Fairbank, J. K. *The Cambridge History of China: Volume 7, The Ming Dynasty, 1368–1644* (Cambridge Univ. Press, London, 1988).
53. Yang, C. et al. Transmission of *Mycobacterium tuberculosis* in China: a population-based molecular epidemiologic study. *Clin. Infect. Dis.* **61**, 219–227 (2015).
54. Ackley, S. F., Liu, F., Porco, T. C. & Pepperell, C. S. Modeling historical tuberculosis epidemics among Canadian First Nations: effects of malnutrition and genetic variation. *PeerJ* **3**, e1237 (2015).
55. Yang, C. et al. *Mycobacterium tuberculosis* Beijing strains favor transmission but not drug resistance in China. *Clin. Infect. Dis.* **55**, 1179–1187 (2012).
56. De Jong, B. C. et al. Progression to active tuberculosis, but not transmission, varies by *Mycobacterium tuberculosis* lineage in The Gambia. *J. Infect. Dis.* **198**, 1037–1043 (2008).
57. Liu, Q. et al. Genetic features of *Mycobacterium tuberculosis* modern Beijing sublineage. *Emerg. Microbes Infect.* **5**, e14 (2016).
58. Van Laarhoven, A. et al. Low induction of proinflammatory cytokines parallels evolutionary success of modern strains within the *Mycobacterium tuberculosis* Beijing genotype. *Infect. Immun.* **81**, 3750–3756 (2013).
59. Ribeiro, S. C. et al. *Mycobacterium tuberculosis* strains of the modern sublineage of the Beijing family are more likely to display increased virulence than strains of the ancient sublineage. *J. Clin. Microbiol.* **52**, 2615–2624 (2014).
60. Ates, L. S. et al. Mutations in *ppe38* block PE_PGRS secretion and increase virulence of *Mycobacterium tuberculosis*. *Nat. Microbiol.* **3**, 181–188 (2018).
61. Kay, G. L. et al. Eighteenth-century genomes show that mixed infections were common at time of peak tuberculosis in Europe. *Nat. Commun.* **6**, 6717 (2015).
62. Wirth, T. Massive lineage replacements and cryptic outbreaks of *Salmonella* Typhi in eastern and southern Africa. *Nat. Genet.* **47**, 565–567 (2015).
63. Wagner, D. M. et al. *Yersinia pestis* and the plague of Justinian 541–543 AD: a genomic analysis. *Lancet Infect. Dis.* **14**, 319–326 (2014).
64. Mtureja, A. et al. Evidence for several waves of global transmission in the seventh cholera pandemic. *Nature* **477**, 462–465 (2011).
65. Vagene, A. J. et al. *Salmonella enterica* genomes from victims of a major sixteenth-century epidemic in Mexico. *Nat. Ecol. Evol.* **2**, 520–528 (2018).
66. Zhao, Y. et al. National survey of drug-resistant tuberculosis in China. *N. Engl. J. Med.* **366**, 2161–2170 (2012).
67. Liu, Q., Luo, T., Li, J., Mei, J. & Gao, Q. Triplex real-time PCR melting curve analysis for detecting *Mycobacterium tuberculosis* mutations associated with resistance to second-line drugs in a single reaction. *J. Antimicrob. Chemother.* **68**, 1097–1103 (2013).
68. Farhat, M. R. et al. Genomic analysis identifies targets of convergent positive selection in drug-resistant *Mycobacterium tuberculosis*. *Nat. Genet.* **45**, 1183–1189 (2013).
69. Comas, I., Homolka, S., Niemann, S. & Gagneux, S. Genotyping of genetically monomorphic bacteria: DNA sequencing in *Mycobacterium tuberculosis* highlights the limitations of current methodologies. *PLoS ONE* **4**, e7815 (2009).
70. Luo, T. et al. Southern East Asian origin and coexpansion of *Mycobacterium tuberculosis* Beijing family with Han Chinese. *Proc. Natl Acad. Sci. USA* **112**, 8136–8141 (2015).
71. Merker, M. et al. Evolutionary history and global spread of the *Mycobacterium tuberculosis* Beijing lineage. *Nat. Genet.* **47**, 242–249 (2015).
72. Barbier, M. & Wirth, T. The evolutionary history, demography, and spread of the *Mycobacterium tuberculosis* complex. *Microbiol. Spectr.* **4**, TBTB2-0008-2016 (2016).
73. Brudey, K. et al. *Mycobacterium tuberculosis* complex genetic diversity: mining the fourth international spoligotyping database (SpolDB4) for classification, population genetics and epidemiology. *BMC Microbiol.* **6**, 23 (2006).
74. Viegas, S. O. et al. Molecular diversity of *Mycobacterium tuberculosis* isolates from patients with pulmonary tuberculosis in Mozambique. *BMC Microbiol.* **10**, 195 (2010).
75. Zhang, H. et al. Genome sequencing of 161 *Mycobacterium tuberculosis* isolates from China identifies genes and intergenic regions associated with drug resistance. *Nat. Genet.* **45**, 1255–1260 (2013).
76. Joshi, N. A. & Fass, J. N. Sickle: A Sliding-Window, Adaptive, Quality-Based Trimming Tool for FastQ Files (2011); <https://github.com/najoshi/sickle>
77. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
78. Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
79. Koboldt, D. C. et al. VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing. *Genome Res.* **22**, 568–576 (2012).
80. Gan, M., Liu, Q., Yang, C., Gao, Q. & Luo, T. Deep whole-genome sequencing to detect mixed infection of *Mycobacterium tuberculosis*. *PLoS ONE* **11**, e0159029 (2016).
81. Tamura, K., Stecher, G., Peterson, D., Filipski, A. & Kumar, S. MEGA6: Molecular Evolutionary Genetics Analysis version 6.0. *Mol. Biol. Evol.* **30**, 2725–2729 (2013).
82. Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res.* **44**, W242–W245 (2016).
83. RStudio Team. *RStudio: Integrated Development for R* (RStudio, 2015).
84. Paradis, E. pegas: an R package for population genetics with an integrated-modular approach. *Bioinformatics* **26**, 419–420 (2010).
85. Hsieh, T., Ma, K. & Chao, A. iNEXT: an R package for rarefaction and extrapolation of species diversity (Hill numbers). *Methods Ecol. Evol.* **7**, 1451–1456 (2016).
86. Yu, Y., Harris, A. J., Blair, C. & He, X. RASP (Reconstruct Ancestral State in Phylogenies): a tool for historical biogeography. *Mol. Phylogenet. Evol.* **87**, 46–49 (2015).
87. Drummond, A. J., Suchard, M. A., Xie, D. & Rambaut, A. Bayesian phylogenetics with BEAUTi and the BEAST 1.7. *Mol. Biol. Evol.* **29**, 1969–1973 (2012).
88. Gignoux, C. R., Henn, B. M. & Mountain, J. L. Rapid, global demographic expansions after the origins of agriculture. *Proc. Natl Acad. Sci. USA* **108**, 6044–6049 (2011).

Acknowledgements

We thank T. M. Walker for sharing the geographic information of 3,651 MTBC isolates of multiple continents origin. We also thank Y.-X. Fu and X. Liu for advice on effective population size calculation and fruitful discussions, and D. Brites and C. Wang for help with clarifying technical details of data analysis during this work. This work was supported by the Natural Science Foundation of China (91631301 and 81661128043 to Q.G. and 81701975 to Q.L.). C.S.P. was supported by the National Institutes of Health (grant 1R01AI113287-01A1). S.G. was supported by the Swiss National Science Foundation (grants IZJRZ3_164171, 3100030_166687, IZLSZ3_170834 and CRSIIS_177163). This work was also supported by MINECO research grant SAF2016-77346-R (to I.C.), the European Research Council (638553-TB-ACCELERATE to I.C.), the National Science and Technology Major Project of China (2017ZX10201302 to Q.G., 2018ZX10103001 to Y.Z.), the Sanming Project of Medicine in Shenzhen (SZSM201611030 to Q.G.), JSGG20170413142559220 (to Q.G.), and National Basic Research programme of China (2014CB744403 to Y.Z.).

Author contributions

Q.L., Y.Z., C.S.P. and Q.G. designed and implemented the study. Y.P., B.W., Y.Z. and Q.G. collected and contributed the MTBC isolates analysed in this study. Q.L., A.M. and Y.Z. conducted the SNP genotyping work. M.L. and C.Y. conducted the MIRU-VNTR typing and analysis. Q.L., T.L., M.G. and T.Z. analysed the sequencing reads and performed the genetic analysis. L.W., H.-X.Z. and L.J. participated in the analysis of integrating tuberculosis history with Chinese human population history. Q.J. performed the statistical analysis. Q.L., I.C., S.G., C.S.P. and Q.G. drafted the manuscript. All authors critically reviewed and approved the final version of the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41559-018-0680-6>.

Reprints and permissions information is available at www.nature.com/reprints.

Correspondence and requests for materials should be addressed to Y.Z. or C.S.P. or Q.G.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

Statistical parameters

When statistical analyses are reported, confirm that the following items are present in the relevant location (e.g. figure legend, table legend, main text, or Methods section).

n/a Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- An indication of whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistics including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated
- Clearly defined error bars
State explicitly what error bars represent (e.g. SD, SE, CI)

Our web collection on [statistics for biologists](#) may be useful.

Software and code

Policy information about [availability of computer code](#)

Data collection

The MTBC isolates were genotyped or whole-genome sequenced in this study were collected by Chinese Center for Disease Control and Prevention, and no software or computer code were applied here.

Data analysis

All commercial codes/softwares used in this study have been described with the exact versions and source links in Methods section. These include Tableau (v10.4), Sickle (), Bowtie2 (v2.2.9), SAMtools (v1.3.1), VarScan (v2.3.9), MEGA (v6.0), FigTree (v1.4.3), iTOL (v3), RStudio (v3.4.0), Sufer (v12), BEAST (v1.8.0). All custom code including R and Perl scripts have been uploaded to GitHub(https://github.com/StopTB/China_TB_Evolutionary_History).

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

Sequencing reads have been submitted to the European Nucleotide Archive (EMBL-EBI) under study accession PRJEB23157. The geographic Information for individual isolates are listed in Supplementary Table 3. The analysis scripts used in this study are available online at GitHub (https://github.com/StopTB/China_TB_Evolutionary_History).

Field-specific reporting

Please select the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/authors/policies/ReportingSummary-flat.pdf

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	A total of 4,578 MTBC isolates were included in this study. The number of MTBC isolates collected from each province was proportional to its TB prevalence.
Data exclusions	702 MTBC isolates that failed in re-culture process were excluded in this study, and these isolates were not within the 4,578 isolates addressed above.
Replication	All phylogenetic trees reconstructed in this study were tested with 500 bootstrap repeats. For each BEAST analysis, at least 30 million chains were generated.
Randomization	Cluster-randomized sampling method was used to collect MTBC isolates in each province in China.
Blinding	Blinding was not relevant here as there is no group allocation in this study.

Reporting for specific materials, systems and methods

Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Unique biological materials
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants

Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging