



Signatures of Selection at Drug Resistance Loci in *Mycobacterium tuberculosis*

 Tatum D. Mortimer,^{a,b} Alexandra M. Weber,^{a,b}  Caitlin S. Pepperell^{a,b}

^aDivision of Infectious Diseases, Department of Medicine, University of Wisconsin—Madison, Madison, Wisconsin, USA

^bDepartment of Medical Microbiology and Immunology, University of Wisconsin—Madison, Madison, Wisconsin, USA

ABSTRACT Tuberculosis (TB) is the leading cause of death by an infectious disease, and global TB control efforts are increasingly threatened by drug resistance in *Mycobacterium tuberculosis*. Unlike most bacteria, where lateral gene transfer is an important mechanism of resistance acquisition, resistant *M. tuberculosis* arises solely by *de novo* chromosomal mutation. Using whole-genome sequencing data from two natural populations of *M. tuberculosis*, we characterized the population genetics of known drug resistance loci using measures of diversity, population differentiation, and convergent evolution. We found resistant subpopulations to be less diverse than susceptible subpopulations, consistent with ongoing transmission of resistant *M. tuberculosis*. A subset of resistance genes (“sloppy targets”) were characterized by high diversity and multiple rare variants; we posit that a large genetic target for resistance and relaxation of purifying selection contribute to high diversity at these loci. For “tight targets” of selection, the path to resistance appeared narrower, evidenced by single favored mutations that arose numerous times in the phylogeny and segregated at markedly different frequencies in resistant and susceptible subpopulations. These results suggest that diverse genetic architectures underlie drug resistance in *M. tuberculosis* and that combined approaches are needed to identify causal mutations. Extrapolating from patterns observed for well-characterized genes, we identified novel candidate variants involved in resistance. The approach outlined here can be extended to identify resistance variants for new drugs, to investigate the genetic architecture of resistance, and when phenotypic data are available, to find candidate genetic loci underlying other positively selected traits in clonal bacteria.

IMPORTANCE *Mycobacterium tuberculosis*, the causative agent of tuberculosis (TB), is a significant burden on global health. Antibiotic treatment imposes strong selective pressure on *M. tuberculosis* populations. Identifying the mutations that cause drug resistance in *M. tuberculosis* is important for guiding TB treatment and halting the spread of drug resistance. Whole-genome sequencing (WGS) of *M. tuberculosis* isolates can be used to identify novel mutations mediating drug resistance and to predict resistance patterns faster than traditional methods of drug susceptibility testing. We have used WGS from natural populations of drug-resistant *M. tuberculosis* to characterize effects of selection for advantageous mutations on patterns of diversity at genes involved in drug resistance. The methods developed here can be used to identify novel advantageous mutations, including new resistance loci, in *M. tuberculosis* and other clonal pathogens.

KEYWORDS *Mycobacterium tuberculosis*, antibiotic resistance, evolution, genomics, positive selection

Received 18 August 2017 **Accepted** 8 January 2018 **Published** 30 January 2018

Citation Mortimer TD, Weber AM, Pepperell CS. 2018. Signatures of selection at drug resistance loci in *Mycobacterium tuberculosis*. mSystems 3:e00108-17. <https://doi.org/10.1128/mSystems.00108-17>.

Editor Jack A. Gilbert, University of Chicago

Copyright © 2018 Mortimer et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

Address correspondence to Caitlin S. Pepperell, cspepper@medicine.wisc.edu.

M*ycobacterium tuberculosis*, the causative agent of tuberculosis (TB), is estimated to have caused 1.4 million deaths in 2015, making it the leading cause of death due to an infectious disease. The proportion of TB due to multidrug-resistant (MDR) *M. tuberculosis* resistant to the first-line antituberculosis drugs isoniazid (INH) and rifampin (RIF) is increasing (1), which poses a major threat to global public health. Unlike most bacteria, *M. tuberculosis* evolves clonally, so resistance cannot be transferred among strains or acquired from other species of bacteria: drug resistance in *M. tuberculosis* results from *de novo* mutation within patients and transmission of drug-resistant clones (2–4). The relative contributions of *de novo* emergence and transmitted drug resistance vary across sampling locations (5–9). Another potential variable in the emergence of drug-resistant TB is *M. tuberculosis*'s lineage structure: seven distinct lineages have been identified among globally extant populations of *M. tuberculosis*. Of these lineages, lineage 2 (L2) has been associated with relatively high rates of drug resistance, and it has been postulated that the acquisition of resistance is a result of higher rates of mutation in this lineage (10). Studies of *M. tuberculosis* evolution within hosts with TB have shown that the emergence of drug resistance is associated with clonal replacements that lead to reductions in genetic diversity of the bacterial population (11, 12).

Many of the methods developed to identify advantageous mutations, such as those conferring antibiotic resistance, depend on recombination to differentiate target loci from neutral variants (13). However, in clonal organisms like *M. tuberculosis*, neutral and deleterious mutations that are linked to advantageous variants will evolve in tandem with them. This linkage among sites can also cause competition between genetic backgrounds with beneficial mutations, decreasing the rate of fixation for beneficial alleles, while deleterious alleles are purged less efficiently (14–16).

While the majority of the *M. tuberculosis* genome is subject to purifying selection (i.e., selection against deleterious mutations) (3), antibiotic pressure exerts strong selection for advantageous variants that confer resistance. *M. tuberculosis* drug resistance has been the focus of extensive investigation, and a variety of resistance mutations have been characterized for commonly used antituberculosis drugs (17). Drug resistance mutations can be associated with fitness costs (18–20), and compensatory mutations that ameliorate these fitness costs have been identified in the context of rifampin resistance (21, 22). Resistance mutations found to have lower fitness costs *in vitro*—as measured by competition assays—are found at higher frequencies among *M. tuberculosis* clinical isolates and appear to be transmitted at higher rates relative to mutations with high *in vitro* fitness costs (18, 23). Candidate loci involved in resistance and compensation for its fitness effects have been identified previously by screening for homoplastic variants (i.e., mutations that emerge more than once in the phylogeny) that are significantly associated with drug-resistant phenotypes (24) and genes with higher *dN/dS* (ratio of nonsynonymous mutations to synonymous mutations) in resistant compared to sensitive isolates (25). Application of these methods to whole-genome sequence data from *M. tuberculosis* clinical isolates has recovered known drug resistance loci, as well as loci associated with cell surface lipids and biosynthesis, DNA replication, and metabolism.

The goal of the present study was to use patterns of genetic diversity at known drug resistance loci to identify the population genomic signatures of positive selection in natural populations of *M. tuberculosis*. Using whole-genome sequence data from two populations for which phenotypic resistance data were available, we have identified several distinct signatures associated with these loci under selection. On the basis of these results, we propose methods of identifying loci under positive selection, including novel drug resistance loci, in clonal bacteria such as *M. tuberculosis*.

RESULTS

We inferred the phylogeny of 1,161 *M. tuberculosis* isolates from Russia and South Africa (see Materials and Methods and see Table S1 in the supplemental material) using the approximate maximum likelihood method implemented in FastTree 2 (Fig. 1). The majority of the isolates belong to lineage 2 (L2) ($n = 667$) and L4 ($n = 481$). *M.*

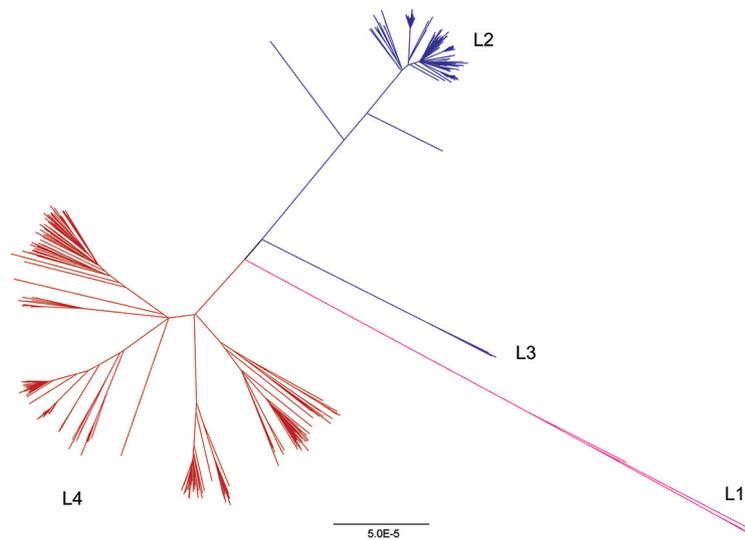


FIG 1 Phylogeny of *Mycobacterium tuberculosis* sample. The phylogeny was inferred using FastTree (60). Lineages are colored as follows: lineage 1 (L1) is shown in pink, lineage 2 (L2) in blue, lineage 3 (L3) in purple, and lineage 4 (L4) in red. Lineage 4 is associated with deeper branching sublineages in comparison with lineage 2. Bar, 5.0E-5 (5.0×10^{-5}) nucleotide substitutions per position.

tuberculosis nucleotide diversity was similar to previous estimates from a globally distributed sample (26). We identified lineage-specific patterns in overall diversity, with L4 having higher diversity than L2 (average number of pairwise differences of L2 [π_{L2}], 3.6×10^{-5} , π_{L4} , 1.5×10^{-4}). Previously published analyses of whole-genome sequence data from L2 indicate that the majority of L2 isolates worldwide belong to a sublineage that has undergone relatively recent expansion (27, 28). In this sample from Russia and South Africa, the majority of L2 isolates belong to this sublineage, while the L4 isolates are associated with deeper branching sublineages. This likely contributes to the observed differences in diversity.

The overall diversity of L2 was lower than that of L4 in our sample (Fig. 2, $P < 2.2 \times$

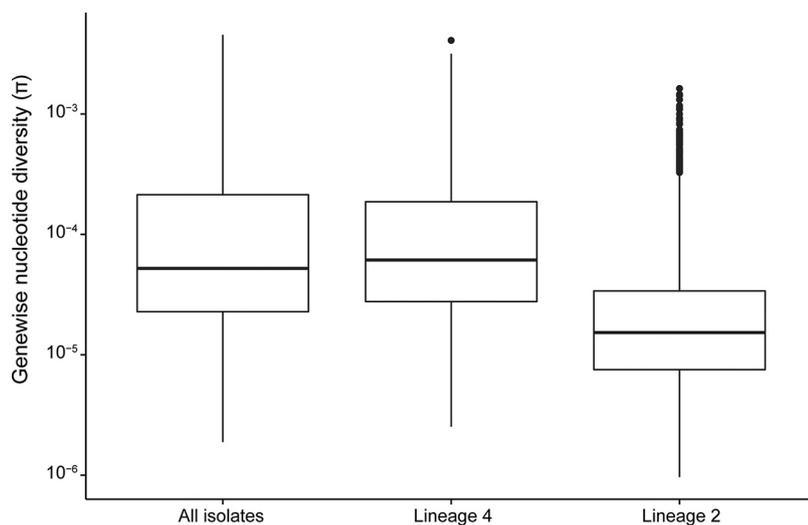


FIG 2 Distributions of gene-wise nucleotide diversity for all isolates, as well as lineages 4 and 2 considered separately. Repetitive regions of the alignment were masked. Sites were included in estimation of π if 95% of isolates in the alignment had a valid nucleotide at the position. We used Egglib to calculate statistics (64). The nucleotide diversity is lower in lineage 2 than in lineage 4 ($P < 2.2 \times 10^{-16}$ by Welch two-sample t test).

TABLE 1 Frequency of resistance in the data set

Drug ^a	Frequency of the following characteristic in the data set:		
	Resistant	Susceptible	Unknown
INH	0.59	0.33	0.08
STR	0.53	0.39	0.07
RIF	0.50	0.43	0.07
EMB	0.30	0.59	0.11
PZA	0.21	0.67	0.12
OFL	0.16	0.39	0.45
PRO	0.15	0.22	0.62
CAP	0.10	0.41	0.49
MOX	0.09	0.41	0.49
ETO	0.06	0.07	0.88
AMI	0.05	0.34	0.61
KAN	0.05	0.12	0.83

^aAMI, amikacin; CAP, capreomycin; EMB, ethambutol; ETO, ethionamide; INH, isoniazid; KAN, kanamycin; MOX, moxifloxacin; OFL, ofloxacin; PRO, prothionamide; PZA, pyrazinamide; RIF, rifampin; STR, streptomycin.

10^{-16}). A total of 762 of the isolates in our sample (66%) are resistant to one or more antituberculosis drugs (Table 1).

Drug-resistant TB can be acquired as a result of *de novo* mutations within a patient or by infection with a resistant strain. When resistance is primarily mediated by *de novo* mutations, diversity should be similar in susceptible and resistant populations, as resistance will arise on multiple genetic backgrounds. In contrast, if resistance develops primarily via transmission of resistant strains, the resistant subpopulation should be less diverse than the susceptible subpopulation. We compared the nucleotide diversity of susceptible and resistant subpopulations and found genome-wide estimates of nucleotide diversity to be higher in isolates susceptible to a range of drugs for which phenotyping data were available ($P = 0.029$ by paired *t* test). In comparisons of gene-wise diversity in susceptible and resistant populations, we found that resistant isolates had a greater number of genes with no diversity, but levels of diversity within genes in which it was measurable were similar between resistant and susceptible populations (Fig. 3).

Of the 3,162 genes included in our analyses, 109 (3%) were invariant across all isolates in our sample. This is likely due to strong purifying selection on these genes. An additional 136 genes harbored variation in the drug-susceptible populations but were invariant across all of the drug-resistant populations. We did not observe the converse, i.e., genes that were invariant in susceptible isolates specifically, which supports the conclusion from genome-wide diversity estimates that resistant isolates represent a subset of the diversity found in susceptible populations and suggests that there may be purifying selection that is specific to the setting of drug resistance. In order to evaluate

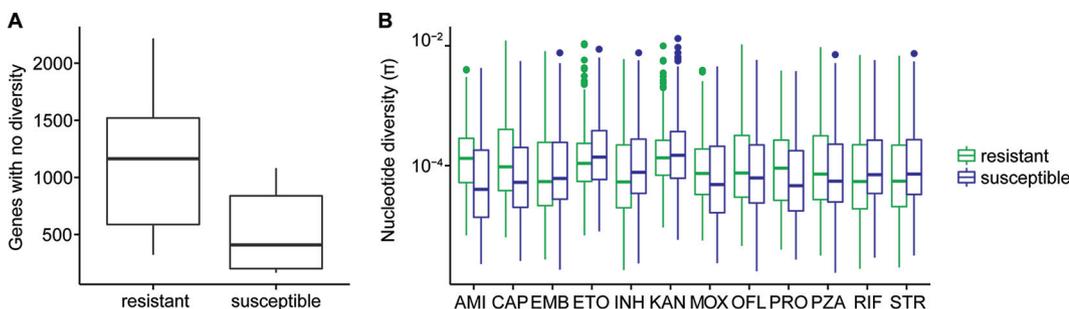


FIG 3 Diversity of resistant and susceptible isolates. (A) Counts of genes with no nucleotide diversity in resistant and susceptible subpopulations. (B) Gene-wise nucleotide diversity (excluding invariant genes) in susceptible and resistant isolates. Among genes in which nucleotide diversity is measurable, it is similar between resistant and susceptible isolates even when drug resistance-associated genes and targets of independent mutation identified by Farhat et al. (24) are removed ($P = 0.13$). Abbreviations: AMI, amikacin; CAP, capreomycin; EMB, ethambutol; ETO, ethionamide; INH, isoniazid; KAN, kanamycin; MOX, moxifloxacin; OFL, ofloxacin; PRO, prothionamide; PZA, pyrazinamide; RIF, rifampin; STR, streptomycin.

whether the observed pattern was likely to arise by chance, we performed weighted random sampling of genes. The weighting was based on diversity in susceptible populations, assuming that genes with low diversity in susceptible populations are more likely to be invariant in resistant populations. After randomly sampling genes in each drug-resistant population 1,000 times, we found that no samples contained shared genes among all resistant populations (first- and second-line drugs). This suggests that specific genes tend to lose diversity in the setting of drug resistance, which could result from purifying selection specific to this setting. A potentially important caveat is that in our data set, drug-resistant populations are not independent, and the same isolates are often resistant to multiple drugs. Since resistance to second-line drugs frequently arises on genetic backgrounds already resistant to one or more first-line drugs, we repeated the sampling with only first-line drugs and found that the maximum number of sampled genes shared across all populations was 2 (median, 0). Overall, these results suggest that certain genes are more likely to lose diversity as drug resistance evolves, but we cannot completely rule out the possibility that the pattern arose as a result of overlapping membership in resistant populations.

We compared the diversity of drug resistance-associated genes (Table 2) with the rest of the genome using two measures of diversity: average number of pairwise differences (π) and number of segregating sites (θ). We found the resistance genes *gid*, *rpsL*, and *pncA* to be in the top fifth percentile of gene-wise π and/or θ values. *rrs* and *ethA* are in the top fifth percentile of θ , but not π . Surprisingly, despite *katG* being a target of multiple drug resistance mutations (Table 2), we did not identify extreme levels of diversity in *katG* (80th and 82nd percentiles of π and θ , respectively).

We also examined gene-wise diversity values within each lineage to look for lineage-specific high-diversity genes. In both L2 and L4, *gid*, *rpsL*, *pncA*, *ethA*, and *thyA* were in the top fifth percentile of diversity (π and/or θ). In L2, *rpoB*, *embB*, *Rv1772*, and *folC* were additionally in the top fifth percentile of gene-wise π and/or θ values. In L4, *Rv0340* was in the top fifth percentile of gene-wise π and/or θ . While *rpoB* and *embB* were not in the top 5th percentile of gene-wise θ in L4, they still had high diversity (91st and 82nd percentiles, respectively). The lineage-specific differences in diversity of *Rv1772*, *folC*, and *Rv0340* suggest that there are interactions between these loci and loci that differentiate L2 and L4.

We used gene-wise estimates of Tajima's D values to investigate gene-specific skews in the site frequency spectrum that could result from selection, where negative values indicate an excess of rare variants and positive values indicate an excess of intermediate frequency variants. We previously identified a relationship between gene length and gene-wise estimates of Tajima's D values for *M. tuberculosis* (26), and this finding was corroborated here ($R^2 = 0.3$ after \log_2 transformation). In order to identify genes with extreme values of Tajima's D values—out of proportion with their length—we performed linear regression on \log_2 -transformed gene lengths and Tajima's D values and identified genes with the largest residuals (Fig. 4). *pncA*, *ethA*, and *embC* all had Tajima's D values lower than expected based on their length (fifth percentile of residual values). This indicates that these genes contain an excess of rare variants compared to other genes in the genome. Excess rare variants can result from a population expansion, a selective sweep, or purifying selection.

We calculated the ratio of π and θ of resistance-associated genes in isolates susceptible and resistant to first-line drugs and identified genes with markedly different diversities in resistant and susceptible subpopulations (Fig. 5A). Among resistance genes in the top fifth percentile of gene-wise π and θ overall, diversity of *pncA* and *ethA* is relatively high among resistant isolates, whereas diversity of *gid* is similar in resistant and susceptible populations. We also examined differences in this ratio between isolates in L2 and L4 (Fig. 5B). *Rv1772* and *embR* were more diverse in resistant isolates in L2, and *kasA* and *tlyA* were more diverse in resistant isolates in L4.

We used fixation index (F_{ST}) outlier analysis to identify single nucleotide polymorphisms (SNPs) and indels that exhibited extreme differences in frequency between susceptible and resistant populations. Our *a priori* expectation was that variants me-

TABLE 2 Signatures of selection in known drug resistance genes

Gene	Locus tag (Rv no.)	Drug(s)	TB Dream ^a	π^b	θ^b	TD ^c	Homoplasmy ^d	F _{ST} ^e	Type(s) ^f
<i>katG</i>	<i>Rv1908c</i>	INH	226	0.80	0.82	0.34	Y	Y	Tight
<i>pncA</i>	<i>Rv2043c</i>	PZA	195	0.97	1.00	0.00	Y	N	Sloppy
<i>embB</i>	<i>Rv3795</i>	EMB	117	0.77	0.89	0.08	Y	Y	Hybrid
<i>ahpC</i>	<i>Rv2428</i>	INH	31	0.20	0.21	0.61	Y	N	
<i>tlyA</i>	<i>Rv1694</i>	CAP	28	0.37	0.89	0.06	N	N	
<i>embC</i>	<i>Rv3793</i>	EMB	28	0.59	0.74	0.01	N	N	
<i>embR</i>	<i>Rv1267c</i>	EMB	25	0.46	0.49	0.28	N	N	
<i>rrs</i>	<i>Rvnr01</i>	STR, KAN, CAP	24	0.89	1.00	0.08	N	N	
<i>ethA</i>	<i>Rv3854c</i>	ETO	23	0.72	1.00	0.00	Y	Y (IG)	Sloppy, tight (IG)
<i>gid</i>	<i>Rv3919c</i>	STR	22	1.00	1.00	0.07	Y	N	Sloppy
<i>gyrB</i>	<i>Rv0005</i>	MOX, OFL	15	0.58	0.91	0.00	Y	N	
<i>fabG1</i>	<i>Rv1483</i>	INH, ETO	13	0.60	0.66	0.30	Y	Y (IG)	Tight
<i>inhA</i>	<i>Rv1484</i>	INH, ETO	13	0.56	0.59	0.32	Y	N	
<i>rpsL</i>	<i>Rv0682</i>	STR	13	0.99	0.95	0.80	Y	Y	Hybrid
<i>gyrA</i>	<i>Rv0006</i>	MOX, OFL	12	0.81	0.94	0.10	Y	Y	Tight
<i>embA</i>	<i>Rv3794</i>	EMB	11	0.77	0.38	0.79	N	N	
<i>kasA</i>	<i>Rv2245</i>	INH	7	0.73	0.18	0.86	N	N	
<i>ndh</i>	<i>Rv1854c</i>	INH	5	0.57	0.52	0.28	N	N	
<i>iniA</i>	<i>Rv0342</i>	EMB, INH	4	0.64	0.33	0.56	N	N	
<i>Rv0340</i>	<i>Rv0340</i>	INH	3	0.89	0.88	0.57	N	N	
<i>iniB</i>	<i>Rv0341</i>	EMB, INH	3	0.07	0.07	0.79	N	N	
<i>fbpC</i>	<i>Rv0129c</i>	INH	3	0.78	0.19	0.89	N	N	
<i>rmID</i>	<i>Rv3266c</i>	EMB	2	0.75	0.36	0.75	N	N	
<i>iniC</i>	<i>Rv0343</i>	EMB, INH	2	0.49	0.67	0.08	N	N	
<i>thyA</i>	<i>Rv2764c</i>	PAS	2	0.84	0.94	0.28	N	N	
<i>nat</i>	<i>Rv3566c</i>	INH	2	0.76	0.55	0.63	N	N	
<i>accD6</i>	<i>Rv2247</i>	INH	1	0.90	0.63	0.90	N	N	
<i>furA</i>	<i>Rv1909c</i>	INH	1	0.80	0.63	0.62	N	N	
<i>Rv1772</i>	<i>Rv1772</i>	INH	1	0.50	0.35	0.54	N	N	
<i>fabD</i>	<i>Rv2243</i>	INH	1	0.26	0.28	0.54	N	N	
<i>fadE24</i>	<i>Rv3139</i>	INH	1	0.36	0.58	0.12	N	N	
<i>rpoB</i>	<i>Rv0667</i>	RIF	1	0.82	0.92	0.18	Y	Y	Hybrid
<i>efpA</i>	<i>Rv2846c</i>	INH	1	0.10	0.11	0.65	N	N	
<i>ethR</i>	<i>Rv3855</i>	ETO		0.58	0.77	0.22	N	N	
<i>Rv0678</i>	<i>Rv0678</i>	BDQ		0.37	0.72	0.25	N	N	
<i>eis</i>	<i>Rv2416c</i>	KAN		0.51	0.28	0.54	N	N	
<i>mshA</i>	<i>Rv0486</i>	ETO		0.86	0.48	0.87	N	N	
<i>rpsA</i>	<i>Rv1630</i>	PZA		0.88	0.62	0.84	N	N	
<i>folC</i>	<i>Rv2447c</i>	PAS		0.66	0.78	0.11	Y	N	
<i>rpIC</i>	<i>Rv0701</i>	LZD		0.57	0.77	0.21	N	N	

^aThe number of distinct entries in the TB Drug Resistance Mutation Database for each gene is reported in the TB Dream column.

^b π and θ are the percentiles for each diversity value, respectively.

^cTD is the percentile of the residual after linear regression of Tajima's D values with gene length.

^dGenes with homoplastic SNPs are indicated with "Y" (for Yes) in the Homoplasmy column. N, no.

^eIf a homoplastic SNP was also an F_{ST} outlier, it is indicated with a "Y" in the F_{ST} column. N, no.

^fGenes are classified as tight, sloppy, or hybrid targets of selection based on diversity, homoplasmy, and F_{ST} results. (IG) indicates an intergenic SNP.

diating resistance would be at markedly higher frequency in the drug-resistant subpopulation and that drug targets would be enriched among genes harboring variants with high F_{ST}. After removing SNPs in regions corresponding to indels and variants at sites missing data for greater than 5% of isolates, the highest F_{ST} SNPs in comparisons of resistant and susceptible subpopulations to first-line drugs are in *katG* (2155168, F_{ST} = 0.89, INH), *rpoB* (761155, F_{ST} = 0.72, RIF), and *rpsL* (781687, F_{ST} = 0.37, streptomycin [STR]). These SNPs were also F_{ST} outliers in the lineage-specific analyses. We used a randomization procedure to assess the significance of observed F_{ST} values and found the maximum F_{ST} values after randomly assigning resistant and susceptible designations to be 0.023 for INH, 0.019 for RIF, and 0.018 for STR. In addition to SNPs within known drug resistance-associated genes, we identified F_{ST} outliers in genes that may be novel targets for drug resistance (Table 3).

Homoplastic SNPs—i.e., SNPs that evolve more than once in a phylogeny—are candidate loci under positive selection and have previously been used to identify

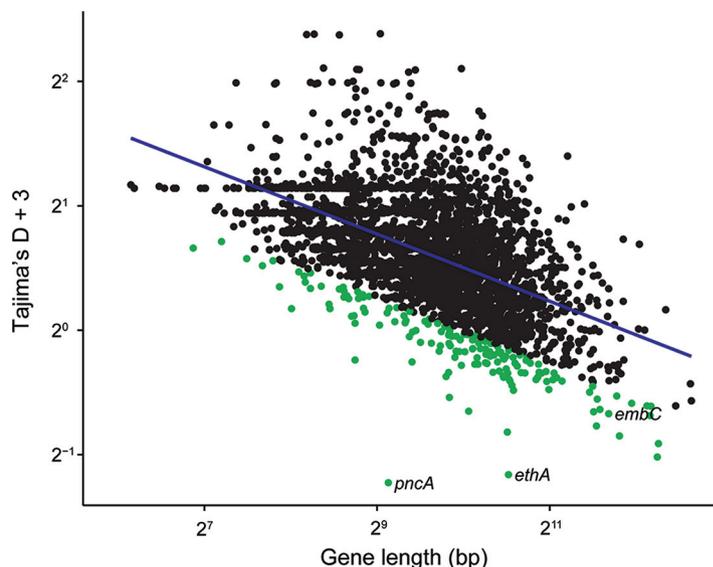


FIG 4 Gene-wise Tajima's D values and gene lengths. Repetitive regions of the alignment were masked. Gene lengths have been log transformed (base 2). We added a constant value of 3 to all Tajima's D values to make them positive and log transformed (base 2) the Tajima's D values. We log transformed (base 2) the gene lengths. The linear regression line is plotted in blue. Genes with regression values in the lower 5% are highlighted in green. Drug resistance-associated genes in this group are labeled. While negative Tajima's D values are normally associated with purifying selection or a recent selective sweep, we find that drug resistance genes with negative Tajima's D values also have high nucleotide diversity. We hypothesize that patterns of diversity at these genes have been affected by relaxation of purifying selection and positive selection for drug resistance.

resistance-associated mutations in *M. tuberculosis* (24). Of the 235 genes with homoplastic SNPs that we identified in our sample, 13 are known to be associated with drug resistance (Fig. 6), and resistance genes were significantly enriched among genes with homoplastic SNPs ($P = 1.2 \times 10^{-4}$ by Fisher's exact test). *pncA* had the largest number of homoplastic SNPs of any gene in the genome (27 distinct SNPs that appear >1 in the phylogeny). The SNPs identified in F_{ST} analysis were also identified as homoplastic (Fig. 6). Our results suggest that complementary approaches based on homoplasy and F_{ST} outlier analysis can be used to identify SNPs associated with a trait of interest (in this case drug resistance). In addition to genic SNPs, we observed homoplastic SNPs that are also F_{ST} outliers in intergenic regions upstream of drug resistance-associated genes (Table 3). These SNPs are candidate resistance and compensatory mutations with a regulatory mechanism of action.

In our analyses of indels, we controlled for the possibility that indels affecting the same gene may not be called in exactly the same position by considering indels within the same gene as identical. We identified four drug resistance-associated genes with homoplastic indels: *gid*, *ethA*, *rpoB*, and *pncA*. F_{ST} values for the deletion in *gid* were in the top fifth percentile for populations resistant to capreomycin (CAP), ethambutol (EMB), ethionamide (ETO), kanamycin (KAN), ofloxacin (OFL), and pyrazinamide (PZA), but interestingly, the deletion was not associated with STR resistance ($F_{ST} = 0.04$). Unlike homoplastic SNPs, homoplastic indels were not significantly enriched for drug resistance-associated loci ($P = 1$).

We recovered 20 out of 40 known drug targets by identifying genes with extreme values of diversity, homoplastic SNPs, or SNPs that are F_{ST} outliers in comparisons of resistant and susceptible subpopulations. All genes with both extremely high diversity (top fifth percentile) and homoplastic mutations were associated with drug resistance (i.e., *gid*, *ethA*, *pncA*, and *rpsL*). We identified 67 genes with high diversity and Tajima's D values more negative than expected based on gene length; only two of these were associated with drug resistance (i.e., *ethA* and *pncA*). Twenty out of 51 homoplastic SNPs that are also F_{ST} outliers fall within or upstream of known drug resistance-associated

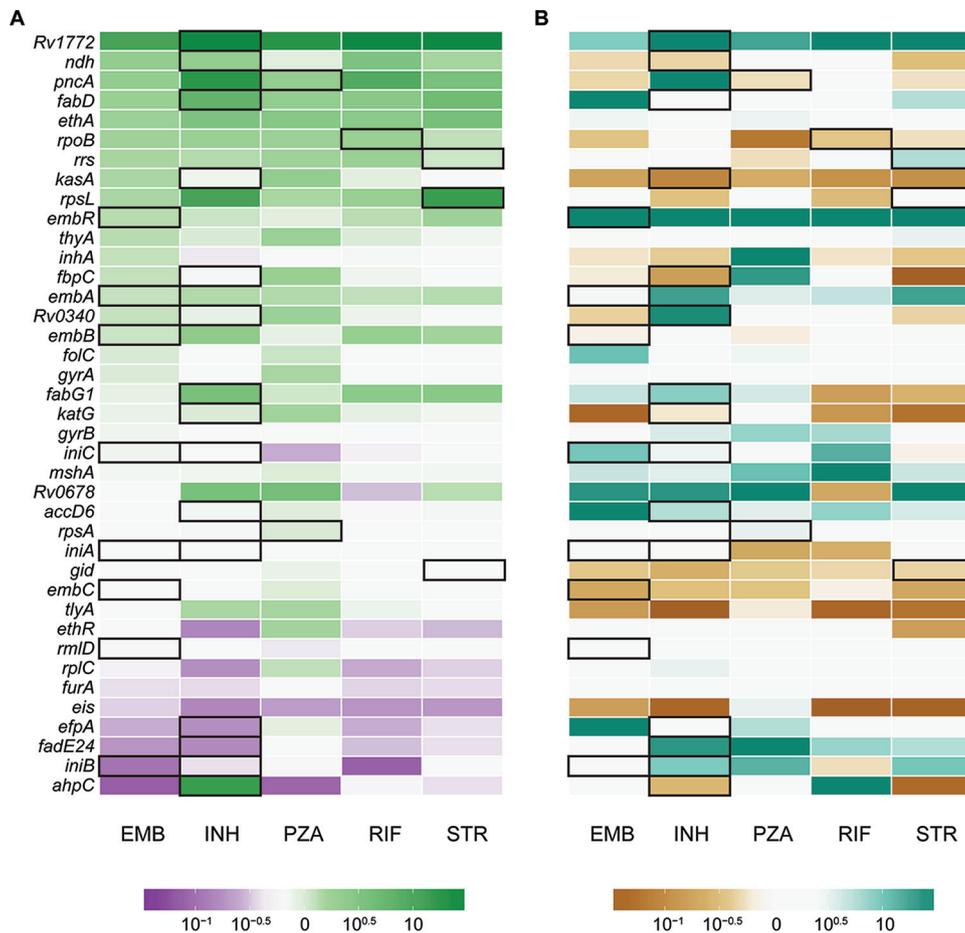


FIG 5 Ratios of nucleotide diversity in resistance-associated genes. Genes with zero diversity were transformed to 1×10^{-16} before calculating ratios. Genes with ratios more extreme than $10^{-1.5}$ or $10^{1.5}$ are all filled with the deepest shade. Genes associated with resistance to each drug are outlined in black. (A) Ratio of nucleotide diversity in resistant and susceptible isolates. Green genes are more diverse in resistant isolates, which could be due to diversifying selection and/or relaxation of purifying selection. Purple genes are more diverse in susceptible isolates, likely due to increased purifying selection. White genes have similar diversities in resistant and susceptible isolates. (B) Comparison of ratios in lineage 2 and lineage 4. Teal genes are more diverse in resistant isolates of lineage 2, suggesting diversifying selection/relaxation of purifying selection specific to this lineage. Brown genes are more diverse in resistant isolates of lineage 4. White genes have similar diversities in lineages 2 and 4.

genes. The remaining SNPs may be false-positive results or novel drug resistance mutations.

DISCUSSION

Highly virulent bacterial pathogens such as *M. tuberculosis*, *Yersinia pestis* (29), *Francisella tularensis* (30), and *Mycobacterium ulcerans* (31) appear to evolve clonally, i.e., with little to no evidence of lateral gene transfer. It is important to identify advantageous mutations in these and other organisms, as they are likely to be associated with phenotypes such as drug resistance, heightened transmissibility, or host adaptation. However, few methods are available for identifying loci under positive selection in the setting of clonal evolution. We adopted an empirical approach to this problem and used natural population data to characterize patterns of diversity at loci known to be under positive selection in *M. tuberculosis*.

In this analysis of clinical isolates from settings where drug resistance is endemic, we found genome-wide diversity to be higher in susceptible *M. tuberculosis* subpopulations than in the subpopulations resistant to first- and second-line drugs (with the exception of populations resistant to prothionamide [PRO] and moxifloxacin [MOX]). The observation of higher diversity in drug-susceptible populations is consistent with

TABLE 3 Homoplastic F_{ST} outliers

Location	Gene ^a	Type	wcF _{ST} ^b value for the following drug:											Known ^c	Lineage			
			AMI	CAP	EMB	ETO	INH	KAN	MOX	OFL	PRO	PZA	RIF			STR		
1821	<i>dnaN</i>	Intergenic				0.43		0.57		0.10							N	All
7570	<i>gyrA</i>	Missense		0.33	0.11	0.46		0.66		0.29		0.18					Y	All
7572	<i>gyrA</i>	Missense							0.06								Y	All
7581	<i>gyrA</i>	Missense							0.07								Y	All
7582	<i>gyrA</i>	Missense							0.35	0.22							Y	All
75233	<i>icd2</i>	Intergenic							0.05								N	All
94388	<i>hycQ</i>	Synonymous			0.07		0.12						0.12	0.13			N	All
230170	<i>Rv0194</i>	Missense					0.12		0.05				0.12	0.13			N	All
332916	<i>vapC25</i>	Missense			0.10					0.09		0.20					N	All
761155	<i>rpoB</i>	Missense			0.31		0.58					0.10	0.72	0.41			Y	All
761161	<i>rpoB</i>	Missense		0.33	0.09	0.51		0.71		0.13		0.16					Y	All
764817	<i>rpoC</i>	Missense	0.19														N	All
781687	<i>rpsL</i>	Missense			0.10		0.32					0.15		0.37			Y	All
922004	<i>Rv0830</i>	Missense		0.30	0.12	0.43				0.10		0.21					N	All
1076880	<i>Rv0965c</i>	Synonymous					0.12						0.12	0.13			N	All
1673425	<i>fabG1</i>	Intergenic									0.11						Y	All
1673432	<i>fabG1</i>	Intergenic				0.52		0.65									Y	All
1722228	<i>pks5</i>	Missense			0.08		0.28			0.07		0.17		0.26			N	All
2122395	<i>lIdD2</i>	Synonymous							0.06								N	All
2155168	<i>katG</i>	Missense			0.36		0.89			0.13		0.32	0.60	0.66			Y	All
2174216	<i>Rv1922</i>	Synonymous										0.08					N	All
2207525	<i>Rv1958c</i>	Intergenic										0.09					N	All
2422824	<i>Rv2161c</i>	Missense		0.30		0.43		0.57		0.10							N	All
2660319	<i>mbtF</i>	Missense			0.06												N	All
2715369	<i>Rv3413c</i>	Intergenic	0.17		0.09		0.28						0.30	0.13			N	All
2866647	<i>lppA</i>	Synonymous					0.12		0.07								N	All
2867298	<i>lppB</i>	Synonymous					0.13										N	All
2867347	<i>lppB</i>	Synonymous					0.13		0.06				0.12	0.14			N	All
2867756	<i>lppB</i>	Synonymous					0.14										N	All
3500149	<i>Rv3134c</i>	Synonymous										0.11					N	All
3550789	<i>Rv3183</i>	Synonymous										0.12	0.13				N	All
3680932	<i>lhr</i>	Synonymous					0.12					0.12	0.13				N	All
4001622	<i>fadA6</i>	Intergenic										0.11					N	All
4247429	<i>embB</i>	Missense			0.25	0.45	0.23		0.05	0.11		0.31	0.21	0.20			Y	All
4247574	<i>embB</i>	Synonymous	0.19		0.07		0.27					0.30					Y	All
4327480	<i>ethA</i>	Intergenic	0.20		0.07		0.27					0.30					Y	All
764948	<i>rpoC</i>	Missense									0.06						Y	L2
4248003	<i>embB</i>	Missense		0.16													Y	L2
698	<i>dnaA</i>	Missense												0.10			N	L4
60185	<i>Rv0057</i>	Missense												0.06			N	L4
761110	<i>rpoB</i>	Missense	0.66														Y	L4
764822	<i>rpoC</i>	Missense												0.06			Y	L4
781822	<i>rpsL</i>	Missense					0.12						0.13	0.14			Y	L4
2123145	<i>lIdD2</i>	Missense												0.06			N	L4
2372550	<i>dop</i>	Missense	0.64														N	L4
2715344	<i>Rv2413c</i>	Intergenic												0.06			N	L4
2986827	<i>Rv2670c</i>	Missense					0.16						0.17	0.15			N	L4
4247431	<i>embB</i>	Missense					0.11						0.11	0.07			Y	L4
4248003	<i>embB</i>	Missense												0.06			Y	L4

^aFor intergenic SNPs, the closest gene is listed.

^bWeir and Cockerham's F_{ST} (wcF_{ST}) values in the top 1% of values genome-wide are reported for each drug.

^cWe identified mutations in genes previously associated with drug resistance (Y for Yes in the Known column) and novel putative resistance or compensatory mutations (N for Novel in the Known column).

a significant role for transmitted resistance in the propagation of drug-resistant *M. tuberculosis*. A recent study of extensively drug-resistant (XDR) *M. tuberculosis* infection in South Africa concluded that XDR cases result primarily from transmission of resistance, rather than *de novo* evolution of resistance mutations during infection (9). The primary studies for the sequence data analyzed here also identified clusters of drug-resistant isolates (5, 6), suggesting that resistant isolates were being transmitted. Our results, along with these previously published observations, suggest that the fitness of drug-resistant isolates can be high enough to allow them to circulate in regions where

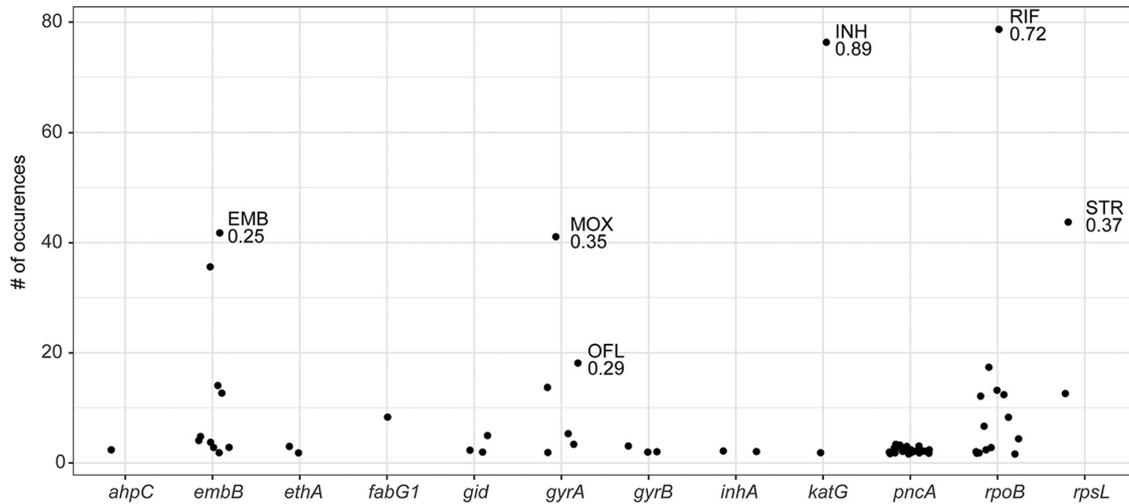


FIG 6 Homoplastic SNPs in drug resistance-associated genes. SNPs with F_{ST} in the top 1% of genome-wide values are labeled with the population (associated drug resistance) and the F_{ST} value. *pncA* is remarkable for harboring diverse homoplastic mutations, each of which occurs relatively infrequently (“sloppy target”). *embB*, *gyrA*, *katG*, *rpoB*, and *rpsL* harbor dominant mutations that occur frequently in the phylogeny and are strongly associated with resistant populations (“tight targets”).

drug-resistant *M. tuberculosis* is endemic. As discussed below, the fitness effects of *M. tuberculosis* drug resistance mutations appear to vary substantially; the finding of transmitted resistance in this and other studies suggests that the fitness of isolates harboring low-cost mutations is comparable to that of susceptible *M. tuberculosis*. The populations in our study have a high burden of drug-resistant TB, and the role of transmitted drug resistance may differ in other settings.

An alternative—but not mutually exclusive—explanation for the observation of higher diversity in susceptible populations is that drug-resistant *M. tuberculosis* is under distinct evolutionary constraints that reduce average genome-wide levels of diversity. In support of this hypothesis, we identified a specific subset of genes that were invariant across drug-resistant populations. Interestingly, while average diversity was lower for resistant subpopulations, the gene-wise diversity distributions had heavier tails, indicating there were more genes with extreme levels of diversity.

We found the genetic architecture of resistance to vary among targets, and resistance-associated genes tended to fall within categories that we term “sloppy,” “tight,” and “hybrid” targets of selection (the latter has a combination of tight and sloppy features and applies to *rpsL*, *embB*, and *rpoB*). “Sloppy” resistance genes are characterized by high levels of diversity. Genes associated with PZA, EMB, ETO, and STR resistance (i.e., *pncA*, *gid*, *rpsL*, *rrs*, and *ethA*) have high levels of diversity; some also had an excess of rare variants (*pncA*, *ethA*, and *embC*). The finding that these genes accumulate multiple, individually rare mutations implies that there is a large target for resistance and/or compensatory mutations within the gene: that is, resistance can result from multiple different variants acting individually or in concert.

In addition to its numerous rare mutations, *pncA* also contains the highest number of homoplastic SNPs (27 SNPs emerged more than once in the phylogeny) of any gene in the data set. Among the 62 nonsynonymous *pncA* mutations in our data set, 55 have been previously reported in association with drug resistance (TB Drug Resistance Mutation Database [32]). The newly described SNPs may mediate drug resistance or compensation for the fitness effects of other variants. Relaxed purifying selection is likely to play a role in concert with selection for diverse advantageous resistance mutations in the accumulation of diversity in *pncA* and other sloppy targets. The fact that numerous mutations are segregating in a natural population suggests that alterations to these genes are generally associated with negligible fitness costs. An *M. tuberculosis* strain harboring a deletion in *pncA* conferring resistance to PZA was estimated to be endemic in Quebec, Canada, by 1800, long before the use of PZA for

the treatment of TB (33–35). This supports the idea that purifying selection on *pncA* is relatively weak, which would contribute to its exceedingly high diversity and broaden the adaptive paths to resistance.

In contrast to *pncA*, *gid*, which is associated with low-level STR resistance (36), does not appear to have the signatures of a “sloppy” target for resistance despite its high diversity. We identified just three homoplastic SNPs within *gid*, and previous studies have found that STR-resistant isolates do not carry the same *gid* mutations (37). This could indicate that a multitude of mutations within *gid* confer resistance, but levels of diversity in the gene were similar in resistant and susceptible isolates. Previous studies of sequence polymorphism in *gid* have identified high diversity in this gene in both resistant and susceptible isolates (37–39): *gid* appears to be subject to relaxed purifying selection in the presence and absence of antibiotic pressure. Since *gid* mutations confer low-level resistance, it is also possible that misclassification of resistance phenotypes contributed to the lack of differentiation we and others have observed between putatively STR-resistant and -susceptible subpopulations. In addition, mutations in *rpsL*, which cause high-level resistance, could mask the contribution of *gid* to STR resistance.

We found some drug targets to be highly diverse in resistant subpopulations of either L2 or L4 (but not both), suggesting that resistance mutations in these genes interact with the genetic background; the fitness effects of mutations in these genes could, for example, vary on different genetic backgrounds. Lineage-specific F_{ST} outliers are another category of candidate locus with lineage-dependent roles in drug resistance (Table 3). Epistatic interactions between drug resistance mutations and *M. tuberculosis* lineage have been reported previously: for example, specific mutations in the *inhA* promoter have been associated with the L1 and *M. africanum* genetic backgrounds (40, 41).

In contrast to “sloppy” targets, we discovered individual homoplastic SNPs associated with drug-resistant subpopulations (i.e., with high F_{ST}) representing “tight” targets of selection in genes conferring resistance to INH, RIF, and STR. Numerous resistance mutations have been reported for *katG*, *rpoB*, *rpsL*, *embB*, and *gyrA*, but we find drug-resistant subpopulations to be defined by a specific subset of mutations in these genes. This suggests that certain mutations are strongly favored relative to others conferring resistance to the same drugs when *M. tuberculosis* is in its natural environment. Antibiotic resistance can impose fitness costs on *M. tuberculosis* during *in vitro* growth, with the range of fitness costs differing among mutations, even within the same gene (18). Mutations can also have different fitness effects depending on the genetic background, but the most fit mutants were the same across *M. tuberculosis* lineages in a study of RIF resistance (18).

In our analyses, we found the dominant INH resistance mutation in *katG* to affect the serine at position 315. This change reduces affinity to INH but preserves catalase activity (42) and is associated with lower fitness costs than other *katG* mutations, both *in vitro* and in a mouse model (43, 44). This mutation was recently shown to precede mutations conferring resistance to other drugs during accumulation of resistance in the evolution of multidrug-resistant *M. tuberculosis* (45). The dominant mutations we identified in *rpoB* (codon 450) and *rpsL* (codon 43) have also been found to have lower fitness costs *in vitro* compared to other mutations conferring resistance to RIF and STR in these genes (18, 44, 46). These results suggest that many of the findings regarding the relative fitness costs of *M. tuberculosis* resistance mutations *in vitro* and in animal models are relevant to the pathogen’s natural environment.

The fitness effects of mutations in *gyrA* (codon 94) and *embB* (codon 306) have not been measured; on the basis of our homoplasy and F_{ST} results, we propose that they have lower fitness costs than other mutations in these genes and that they represent “tight” targets of selection. Mutations at *gyrA* codon 94 were previously found to be the most prevalent in a survey of *gyrA* and *gyrB* mutations in fluoroquinolone-resistant *M. tuberculosis* clinical isolates (47). Interestingly, the mutation in *embB* codon 306 has been previously associated with acquisition of multiple resistances (48), and we find that this position is an F_{ST} outlier for all first-line drugs in L4. This mutation is not an F_{ST}

outlier in L2 (i.e., top fifth percentile), with the percentiles for F_{ST} values ranging from 0.07 to 0.68 for first-line drugs in this lineage. These observations suggest that the genetic background affects interactions among resistance mutations and that *embB* codon 306 is important for acquisition of multidrug resistance in L4, but not in L2.

We searched for indels with the signature of a “tight” target, i.e., homoplastic mutations segregating at markedly different frequencies in drug-susceptible and -resistant subpopulations. Unlike the pattern observed with SNPs, genes associated with drug resistance were not significantly enriched among the genes harboring homoplastic indels. We identified one homoplastic indel that was also an F_{ST} outlier—a deletion in *gid* that causes a frameshift. Patterns of variation in *gid* are complex and suggest a role for relaxation of purifying selection (i.e., in the accumulation of excess SNPs in both resistant and susceptible isolates) and perhaps a tight target associated with multiresistance (i.e., this homoplastic/ F_{ST} outlier deletion that was associated with resistance to CAP, EMB, ETO, KAN, OFL, and PZA).

Our finding that, save for the frameshift mutation in *gid*, indels in resistance genes do not have the signature of “tight” targets suggests that they are generally associated with higher fitness costs than SNPs. Fifteen drug targets have been found in transposon mutagenesis experiments to be essential for *M. tuberculosis* growth *in vitro*, including *rpoB* and *rpsL*; deletions in these genes are likely to interrupt important functions (49). Deletions in nonessential genes could also have fitness costs. Deletions in *katG*, which is nonessential, can result in INH resistance, but they are not observed as frequently in clinical isolates as the KatG S315 SNP is, particularly among transmitted INH-resistant strains (23).

There are several limitations to our study. Resistance to multiple drugs was common in our sample, and in some cases, it was difficult to identify patterns of diversity and population differentiation that were specific to individual drugs. Our results are also limited by the accuracy with which drug resistance phenotypes were determined and a lack of phenotypic data for some drugs (particularly second-line drugs). Our sample was heavily skewed to lineages 2 and 4, and the results are not necessarily applicable to other *M. tuberculosis* lineages. Finally, the data analyzed here were generated with short sequencing read technologies, and we were thus limited to characterizing diversity in regions of the *M. tuberculosis* genome that can be resolved with these methods: regions that were masked from analysis (e.g., due to sequence repeats) may include unknown resistance targets. We also used an L4 genome (H37Rv) as a reference, and gene content specific to L2 may not have been identified.

We were not able to recover all drug resistance-associated genes using the analyses performed here. This is likely a result of limited phenotypic data for some drugs and their associated targets (e.g., *thyA* and *folC*, which are associated with aminosalicic acid [PAS] resistance). Our list of drug targets was dominated by genes associated with INH resistance, and signatures in genes that harbor rare resistance-associated alleles may be subtle compared to the KatG S315 mutation found at high frequency in drug-resistant populations.

We identified 31 SNPs that are not on the list of known drug resistance genes, which both emerged more than once in the phylogeny (homoplasies) and were segregating at markedly different frequencies in resistant and susceptible subpopulations (F_{ST} outliers). These SNPs may be novel resistance determinants; notably, all nonsynonymous SNPs within this group are in genes linked with drug resistance in other studies (i.e., they are in genes encoding efflux pumps, genes differentially regulated in resistant isolates or in response to the presence of drug, potential drug targets, or genes in the same pathways as drug targets or resistance determinants) (50–54). In addition to a direct, previously unrecognized role in resistance, these SNPs could compensate for fitness costs of drug resistance. For example, we identified a homoplastic F_{ST} outlier in *rpoC*, and mutations in *rpoC* have been shown to compensate for RIF resistance in experimental evolution studies (22).

Intriguingly, we found lipid metabolism genes to be enriched among genes harboring homoplastic SNPs ($P = 0.013$). We have previously shown that these genes have

extreme values of diversity in a global sample of *M. tuberculosis* isolates and within individual hosts (26), suggesting that lipid metabolism genes may also be under positive selection in *M. tuberculosis* populations. The results presented here could be extended by phenotypic characterization of lipid profiles and identification of homoplastic variants that are at markedly different frequencies in isolates with distinct lipid profiles.

Here we have used drug resistance loci in *M. tuberculosis* to identify the signatures of positive selection in a clonal bacterium. We found these loci to be associated with distinct patterns of diversity that likely reflect differing genetic architectures underlying the traits under selection. The evolutionary path to resistance is broad for some drugs with “sloppy targets,” whereas for drugs with “tight targets,” the means of acquiring resistance appear more limited. This is likely due to fitness effects of resistance mutations in *M. tuberculosis*'s natural environment, as numerous resistance mutations have been identified in tight target genes. We also found evidence suggesting that there are important interactions among loci during the evolution of resistance. Our results suggest that purifying selection on a subset of genes intensifies in the setting of resistance, which could reflect epistatic interactions and/or a response to the metabolic milieu imposed by antimycobacterial agents. The results presented here can be used to create more realistic models of resistance evolution in *M. tuberculosis* and to develop novel strategies of preventing or mitigating the acquisition of resistance. For example, the narrow path to resistance for drugs with tight targets reveals potentially exploitable vulnerabilities, as does the finding of interdependencies among specific loci and the genetic background in the evolution of resistance and multiresistance. As new TB drugs become available for clinical use, the approach outlined here can be extended to investigate their architectures of resistance.

Efforts are under way to sequence and perform drug susceptibility testing on thousands of *M. tuberculosis* isolates with the goal of creating an exhaustive catalogue of drug resistance mutations and eventually using whole-genome sequencing (WGS) to diagnose drug resistance in clinical settings (CRyPTIC project, <http://modmedmicro.nsms.ox.ac.uk/cryptic/>, last accessed 24 May 2017). We found that loci under positive selection can be identified using relatively simple methods: “tight” targets are highly differentiated in their allele frequencies across phenotypic groups (i.e., F_{ST} outliers) and appear as homoplasies on the phylogeny; “sloppy” targets are characterized by high diversity and/or low Tajima's D values, as well as homoplasies. Extrapolating from patterns observed among known resistance variants, we have discovered new candidate regulatory and genic resistance variants. The methods used in this study are widely available and should scale to analysis of the large collections of genomic and phenotypic data that are currently being generated. This approach can be extended to identify novel resistance loci in bacteria for which drug susceptibility phenotypes are defined, as well as other positively selected loci in clonal bacterial populations.

MATERIALS AND METHODS

Reference guided assembly. We downloaded sequencing read data from two large surveys of drug-resistant *M. tuberculosis* in Russia (5) and South Africa (6). We used FastQC (55) and TrimGalore! (56) for quality assessment and adaptor trimming of the reads. Trimmed reads were mapped to *M. tuberculosis* H37Rv (accession no. NC_000962.3) using BWA-MEM v 0.7.12 (57). We used Samtools v 1.2 (58) and Picard Tools (<https://broadinstitute.github.io/picard/>) for sorting, format conversion, and addition of read group information. Variants were identified using Pilon v 1.16 (59). A detailed description of the reference guided assembly pipeline is available at <https://github.com/pepperell-lab/RGAPepPipe>. We removed isolates with mean coverage less than $20\times$, isolates with percentage of the genome covered at $10\times$ less than 90%, isolates where a majority of reads did not map to H37Rv, and isolates where greater than 10% of sites were unknown after mapping. The final data set contains 1,161 *M. tuberculosis* isolates (see Table S1 in the supplemental material). The alignment was masked to remove repetitive regions, including proline-glutamate (PE)/proline-proline-glutamate (PPE) genes.

Phylogenetic analysis. We estimated the approximately maximum likelihood phylogeny using the masked alignment from reference guided assembly with FastTree 2.1.9 (60). We compiled FastTree using the double precision option to accurately estimate branch lengths of closely related isolates. We used FigTree (<http://tree.bio.ed.ac.uk/software/figtree/>) for tree visualization.

SNP annotation. A VCF (variant call format) file of single nucleotide variants was created from the masked alignment using SNP-sites v 2.3.2 (61). Single nucleotide polymorphisms (SNPs) were annotated

using SnpEff v 4.1j (62) to identify synonymous, nonsynonymous, and intergenic SNPs based on the annotation of *M. tuberculosis* H37Rv.

Indel identification. Insertions and deletions were identified during variant calling with Pilon. We used Emu (63) to normalize indels across multiple isolates. We used a presence/absence matrix for the normalized indels for further analyses of indel diversity.

Population genetics statistics. Whole-genome and gene-wise diversity (π and θ) and neutrality (Tajima's D) statistics were calculated using Egglib v 2.1.10 (64) for whole-genome alignments and gene-wise alignments. Isolates were further divided by lineage and drug resistance phenotype. Sites with missing data due to indels or low-quality base calls in more than 5% of isolates in the alignment were not included in calculation of statistics. Tajima's D values showed a correlation with gene length in our sample. To find genes with extreme values of Tajima's D, we performed linear regression in R (65) on log-transformed Tajima's D values and gene length and identified genes with large residual values. To identify alleles with marked differences in frequency in resistant and susceptible isolates, Weir and Cockerham's F_{ST} (66) was calculated using populations of resistant and susceptible isolates for each drug using vcflib v1.0.0-rc0-262-g50a3 (<https://github.com/vcflib/vcflib>). For nonbiallelic SNPs, we calculated F_{ST} for the two most common variants.

Homoplasy. We used TreeTime (67) to perform ancestral reconstruction and place SNPs and indels in the phylogeny. We identified homoplastic SNPs and indels as those arising multiple times in the phylogeny.

Data availability. Unless otherwise noted, all data and scripts associated with this study are available at <https://github.com/pepperell-lab/mtbDrugResistance>.

SUPPLEMENTAL MATERIAL

Supplemental material for this article may be found at <https://doi.org/10.1128/mSystems.00108-17>.

TABLE S1, TXT file, 0.02 MB.

ACKNOWLEDGMENTS

We thank members of the Pepperell laboratory for their input on analyses and data visualization.

This material is based upon work supported by the National Science Foundation Graduate Research Fellowship Program under grant DGE-1256259 to T.D.M. T.D.M. is also supported by National Institutes of Health National Research Service award (T32 GM07215). C.S.P. is supported by National Institutes of Health (R01AI113287). Funding for this project was provided by the University of Wisconsin School of Medicine and Public Health from the Wisconsin Partnership Program.

Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation or other funding agencies.

REFERENCES

- World Health Organization. 2016. Global tuberculosis report 2016. World Health Organization, Geneva, Switzerland.
- Supply P, Warren RM, Bañuls A-L, Lesjean S, Van Der Spuy GD, Lewis L-A, Tibayrenc M, Van Helden PD, Loch C. 2003. Linkage disequilibrium between minisatellite loci supports clonal evolution of *Mycobacterium tuberculosis* in a high tuberculosis incidence area. *Mol Microbiol* 47: 529–538. <https://doi.org/10.1046/j.1365-2958.2003.03315.x>.
- Pepperell CS, Casto AM, Kitchen A, Granka JM, Cornejo OE, Holmes EC, Birren B, Galagan J, Feldman MW. 2013. The role of selection in shaping diversity of natural *M. tuberculosis* populations. *PLoS Pathog* 9:e1003543. <https://doi.org/10.1371/journal.ppat.1003543>.
- Eldholm V, Balloux F. 2016. Antimicrobial resistance in *Mycobacterium tuberculosis*: the odd one out. *Trends Microbiol* 24:637–648. <https://doi.org/10.1016/j.tim.2016.03.007>.
- Casali N, Nikolayevskyy V, Balabanova Y, Harris SR, Ignatyeva O, Kontsevaya I, Corander J, Bryant J, Parkhill J, Nejentsev S, Horstmann RD, Brown T, Drobniowski F. 2014. Evolution and transmission of drug-resistant tuberculosis in a Russian population. *Nat Genet* 46:279–286. <https://doi.org/10.1038/ng.2878>.
- Cohen KA, Abeel T, Manson McGuire A, Desjardins CA, Munsamy V, Shea TP, Walker BJ, Bantubani N, Almeida DV, Alvarado L, Chapman SB, Mvelase NR, Duffy EY, Fitzgerald MG, Govender P, Gujja S, Hamilton S, Howarth C, Larimer JD, Maharaj K, Pearson MD, Priest ME, Zeng Q, Padayatchi N, Grosset J, Young SK, Wortman J, Mlisana KP, O'Donnell MR, Birren BW, Bishai WR, Pym AS, Earl AM. 2015. Evolution of extensively drug-resistant tuberculosis over four decades: whole genome sequencing and dating analysis of *Mycobacterium tuberculosis* isolates from KwaZulu-Natal. *PLoS Med* 12:e1001880. <https://doi.org/10.1371/journal.pmed.1001880>.
- Coscolla M, Barry PM, Oeltmann JE, Koshinsky H, Shaw T, Cilnis M, Posey J, Rose J, Weber T, Fofanov VY, Gagneux S, Kato-Maeda M, Metcalfe JZ. 2015. Genomic epidemiology of multidrug-resistant *Mycobacterium tuberculosis* during transcontinental spread. *J Infect Dis* 212:302–310. <https://doi.org/10.1093/infdis/jiv025>.
- Guerra-Assunção JA, Crampin AC, Houben R, Mzembe T, Mallard K, Coll F, Khan P, Banda L, Chiwaya A, Pereira RPA, McNerney R, Fine PEM, Parkhill J, Clark TG, Glynn JR. 2015. Large-scale whole genome sequencing of *M. tuberculosis* provides insights into transmission in a high prevalence area. *eLife* 4:e05166. <https://doi.org/10.7554/eLife.05166>.
- Shah NS, Auld SC, Brust JCM, Mathema B, Ismail N, Moodley P, Mlisana K, Allana S, Campbell A, Mthiyane T, Morris N, Mpangase P, van der Meulen H, Omar SV, Brown TS, Narechana A, Shaskina E, Kapwata T, Kreiswirth B, Gandhi NR. 2017. Transmission of extensively drug-resistant tuberculosis in South Africa. *N Engl J Med* 376:243–253. <https://doi.org/10.1056/NEJMoa1604544>.
- Ford CB, Shah RR, Maeda MK, Gagneux S, Murray MB, Cohen T, Johnston JC, Gardy J, Lipsitch M, Fortune SM. 2013. *Mycobacterium tuberculosis* mutation rate estimates from different lineages predict substantial dif-

- ferences in the emergence of drug-resistant tuberculosis. *Nat Genet* 45:784–790. <https://doi.org/10.1038/ng.2656>.
11. Eldholm V, Norheim G, von der Lippe B, Kinander W, Dahle UR, Caugant DA, Mannsåker T, Mengshoel AT, Dyrhol-Riise AM, Balloux F. 2014. Evolution of extensively drug-resistant *Mycobacterium tuberculosis* from a susceptible ancestor in a single patient. *Genome Biol* 15:490. <https://doi.org/10.1186/s13059-014-0490-3>.
 12. Black PA, de Vos M, Louw GE, van der Merwe RG, Dippenaar A, Streicher EM, Abdallah AM, Sampson SL, Victor TC, Dolby T, Simpson JA, van Helden PD, Warren RM, Pain A. 2015. Whole genome sequencing reveals genomic heterogeneity and antibiotic purification in *Mycobacterium tuberculosis* isolates. *BMC Genomics* 16:857. <https://doi.org/10.1186/s12864-015-2067-2>.
 13. Vitti JJ, Grossman SR, Sabeti PC. 2013. Detecting natural selection in genomic data. *Annu Rev Genet* 47:97–120. <https://doi.org/10.1146/annurev-genet-111212-133526>.
 14. Hill WG, Robertson A. 1966. The effect of linkage on limits to artificial selection. *Genet Res* 8:269–294. <https://doi.org/10.1017/S0016672300010156>.
 15. Felsenstein J. 1974. The evolutionary advantage of recombination. *Genetics* 78:737–756.
 16. Birky CW, Walsh JB. 1988. Effects of linkage on rates of molecular evolution. *Proc Natl Acad Sci U S A* 85:6414–6418. <https://doi.org/10.1073/pnas.85.17.6414>.
 17. Almeida Da Silva PE, Palomino JC. 2011. Molecular basis and mechanisms of drug resistance in *Mycobacterium tuberculosis*: classical and new drugs. *J Antimicrob Chemother* 66:1417–1430. <https://doi.org/10.1093/jac/dkr173>.
 18. Gagneux S, Long CD, Small PM, Van T, Schoolnik GK, Bohannon BJM. 2006. The competitive cost of antibiotic resistance in *Mycobacterium tuberculosis*. *Science* 312:1944–1946. <https://doi.org/10.1126/science.1124410>.
 19. Böttger EC, Springer B. 2008. Tuberculosis: drug resistance, fitness, and strategies for global control. *Eur J Pediatr* 167:141–148. <https://doi.org/10.1007/s00431-007-0606-9>.
 20. Strauss OJ, Warren RM, Jordaan A, Streicher EM, Hanekom M, Falmer AA, Albert H, Trollip A, Hoosain E, van Helden PD, Victor TC. 2008. Spread of a low-fitness drug-resistant *Mycobacterium tuberculosis* strain in a setting of high human immunodeficiency virus prevalence. *J Clin Microbiol* 46:1514–1516. <https://doi.org/10.1128/JCM.01938-07>.
 21. Brandis G, Wrande M, Liljas L, Hughes D. 2012. Fitness-compensatory mutations in rifampicin-resistant RNA polymerase. *Mol Microbiol* 85:142–151. <https://doi.org/10.1111/j.1365-2958.2012.08099.x>.
 22. Comas I, Borrell S, Roetzer A, Rose G, Malla B, Kato-Maeda M, Galagan J, Niemann S, Gagneux S. 2011. Whole-genome sequencing of rifampicin-resistant *Mycobacterium tuberculosis* strains identifies compensatory mutations in RNA polymerase genes. *Nat Genet* 44:106–110. <https://doi.org/10.1038/ng.1038>.
 23. Gagneux S, Burgos MV, DeRiemer K, Encisco A, Muñoz S, Hopewell PC, Small PM, Pym AS. 2006. Impact of bacterial genetics on the transmission of isoniazid-resistant *Mycobacterium tuberculosis*. *PLoS Pathog* 2:e61. <https://doi.org/10.1371/journal.ppat.0020061>.
 24. Farhat MR, Shapiro BJ, Kieser KJ, Sultana R, Jacobson KR, Victor TC, Warren RM, Streicher EM, Calver A, Sloutsky A, Kaur D, Posey JE, Plikaytis B, Oggioni MR, Gardy JL, Johnston JC, Rodrigues M, Tang PKC, Kato-Maeda M, Borowsky ML, Muddukrishna B, Kreiswirth BN, Kurepina N, Galagan J, Gagneux S, Birren B, Rubin EJ, Lander ES, Sabeti PC, Murray M. 2013. Genomic analysis identifies targets of convergent positive selection in drug-resistant *Mycobacterium tuberculosis*. *Nat Genet* 45:1183–1189. <https://doi.org/10.1038/ng.2747>.
 25. Zhang H, Li D, Zhao L, Fleming J, Lin N, Wang T, Liu Z, Li C, Galwey N, Deng J, Zhou Y, Zhu Y, Gao Y, Wang T, Wang S, Huang Y, Wang M, Zhong Q, Zhou L, Chen T, Zhou J, Yang R, Zhu G, Hang H, Zhang J, Li F, Wan K, Wang J, Zhang X-E, Bi L. 2013. Genome sequencing of 161 *Mycobacterium tuberculosis* isolates from China identifies genes and intergenic regions associated with drug resistance. *Nat Genet* 45:1255–1260. <https://doi.org/10.1038/ng.2735>.
 26. O'Neill MB, Mortimer TD, Pepperell CS. 2015. Diversity of *Mycobacterium tuberculosis* across evolutionary scales. *PLoS Pathog* 11:e1005257. <https://doi.org/10.1371/journal.ppat.1005257>.
 27. Luo T, Comas I, Luo D, Lu B, Wu J, Wei L, Yang C, Liu Q, Gan M, Sun G, Shen X, Liu F, Gagneux S, Mei J, Lan R, Wan K, Gao Q. 2015. Southern East Asian origin and coexpansion of *Mycobacterium tuberculosis* Beijing family with Han Chinese. *Proc Natl Acad Sci U S A* 112:8136–8141. <https://doi.org/10.1073/pnas.1424063112>.
 28. O'Neill MB, Kitchen A, Zarley A, Aylward W, Eldholm V, Pepperell CS. 2017. Lineage specific histories of *Mycobacterium tuberculosis* dispersal in Africa and Eurasia. *bioRxiv* <https://doi.org/10.1101/210161>.
 29. Achtman M, Zurth K, Morelli G, Torrea G, Guiyoule A, Carniel E. 1999. *Yersinia pestis*, the cause of plague, is a recently emerged clone of *Yersinia pseudotuberculosis*. *Proc Natl Acad Sci U S A* 96:14043–14048. <https://doi.org/10.1073/pnas.96.24.14043>.
 30. Larsson P, Elfsmark D, Svensson K, Wikström P, Forsman M, Brettin T, Keim P, Johansson A. 2009. Molecular evolutionary consequences of niche restriction in *Francisella tularensis*, a facultative intracellular pathogen. *PLoS Pathog* 5:e1000472. <https://doi.org/10.1371/journal.ppat.1000472>.
 31. Vandellannoote K, Meehan CJ, Eddyani M, Affolabi D, Phanuz DM, Eyangoh S, Jordaens K, Portaels F, Mangas K, Seemann T, Marsollier L, Marion E, Chauty A, Landier J, Fontanet A, Leirs H, Steiner TP, de Jong BC. 2017. Multiple introductions and recent spread of the emerging human pathogen *Mycobacterium ulcerans* across Africa. *Genome Biol Evol* 9:414–426. <https://doi.org/10.1093/gbe/evx003>.
 32. Sandgren A, Strong M, Muthukrishnan P, Weiner BK, Church GM, Murray MB. 2009. Tuberculosis Drug Resistance Mutation Database. *PLoS Med* 6:e1000002. <https://doi.org/10.1371/journal.pmed.1000002>.
 33. Nguyen D, Brassard P, Westley J, Thibert L, Proulx M, Henry K, Schwartzman K, Menzies D, Behr MA. 2003. Widespread pyrazinamide-resistant *Mycobacterium tuberculosis* family in a low-incidence setting. *J Clin Microbiol* 41:2878–2883. <https://doi.org/10.1128/JCM.41.7.2878-2883.2003>.
 34. Nguyen D, Brassard P, Menzies D, Thibert L, Warren R, Mostowy S, Behr M. 2004. Genomic characterization of an endemic *Mycobacterium tuberculosis* strain: evolutionary and epidemiologic implications. *J Clin Microbiol* 42:2573–2580. <https://doi.org/10.1128/JCM.42.6.2573-2580.2004>.
 35. Brassard P, Henry KA, Schwartzman K, Jomphe M, Olson SH. 2008. Geography and genealogy of the human host harbouring a distinctive drug-resistant strain of tuberculosis. *Infect Genet Evol* 8:247–257. <https://doi.org/10.1016/j.meegid.2007.11.008>.
 36. Okamoto S, Tamaru A, Nakajima C, Nishimura K, Tanaka Y, Tokuyama S, Suzuki Y, Ochi K. 2007. Loss of a conserved 7-methylguanosine modification in 16S rRNA confers low-level streptomycin resistance in bacteria. *Mol Microbiol* 63:1096–1106. <https://doi.org/10.1111/j.1365-2958.2006.05585.x>.
 37. Spies FS, Ribeiro AW, Ramos DF, Ribeiro MO, Martin A, Palomino JC, Rossetti MLR, da Silva PEA, Zaha A. 2011. Streptomycin resistance and lineage-specific polymorphisms in *Mycobacterium tuberculosis* gidB gene. *J Clin Microbiol* 49:2625–2630. <https://doi.org/10.1128/JCM.00168-11>.
 38. Feuerriegel S, Oberhauser B, George AG, Dafee F, Richter E, Rüscher S, Niemann S. 2012. Sequence analysis for detection of first-line drug resistance in *Mycobacterium tuberculosis* strains from a high-incidence setting. *BMC Microbiol* 12:90. <https://doi.org/10.1186/1471-2180-12-90>.
 39. Jagielski T, Ignatowska H, Bakula Z, Dziewit Ł, Napiórkowska A, Augustynowicz-Kopeć E, Zwolska Z, Bielecki J. 2014. Screening for streptomycin resistance-conferring mutations in *Mycobacterium tuberculosis* clinical isolates from Poland. *PLoS One* 9:e100078. <https://doi.org/10.1371/journal.pone.0100078>.
 40. Homolka S, Meyer CG, Hillemann D, Owusu-Dabo E, Adjei O, Horstmann RD, Browne ENL, Chinbuah A, Osei I, Gyapong J, Kubica T, Ruesch-Gerdes S, Niemann S. 2010. Unequal distribution of resistance-conferring mutations among *Mycobacterium tuberculosis* and *Mycobacterium africanum* strains from Ghana. *Int J Med Microbiol* 300:489–495. <https://doi.org/10.1016/j.ijmm.2010.04.019>.
 41. Fenner L, Egger M, Bodmer T, Altpeter E, Zwahlen M, Jaton K, Pfyffer GE, Borrell S, Dubuis O, Bruderer T, Siegrist HH, Furrer H, Calmy A, Fehr J, Stalder JM, Ninet B, Böttger EC, Gagneux S, Swiss HIV Cohort Study, Swiss Molecular Epidemiology of Tuberculosis Study Group. 2012. Effect of mutation and genetic background on drug resistance in *Mycobacterium tuberculosis*. *Antimicrob Agents Chemother* 56:3047–3053. <https://doi.org/10.1128/AAC.06460-11>.
 42. Wengenack NL, Uhl JR, St Amand AL, Tomlinson AJ, Benson LM, Naylor S, Kline BC, Cockerill FR, III, Rusnak F. 1997. Recombinant *Mycobacterium tuberculosis* KatG(S315T) is a competent catalase-peroxidase with reduced activity toward isoniazid. *J Infect Dis* 176:722–727. <https://doi.org/10.1086/514096>.
 43. Pym AS, Saint-Joanis B, Cole ST. 2002. Effect of katG mutations on the virulence of *Mycobacterium tuberculosis* and the implication for trans-

- mission in humans. *Infect Immun* 70:4955–4960. <https://doi.org/10.1128/IAI.70.9.4955-4960.2002>.
44. Spies FS, von Groll A, Ribeiro AW, Ramos DF, Ribeiro MO, Dalla Costa ER, Martin A, Palomino JC, Rossetti ML, Zaha A, da Silva PEA. 2013. Biological cost in *Mycobacterium tuberculosis* with mutations in the *rpsL*, *rrs*, *rpoB*, and *katG* genes. *Tuberculosis* 93:150–154. <https://doi.org/10.1016/j.tube.2012.11.004>.
 45. Manson AL, Cohen KA, Abeel T, Desjardins CA, Armstrong DT, Barry CE, III, Brand J, TBResist Global Genome Consortium, Chapman SB, Cho S-N, Gabrielian A, Gomez J, Jodals AM, Joloba M, Jureen P, Lee JS, Malinga L, Maiga M, Nordenberg D, Noroc E, Romancenco E, Salazar A, Ssengooba W, Velayati AA, Wingless K, Zalutskaya A, Via LE, Cassell GH, Dorman SE, Ellner J, Farnia P, Galagan JE, Rosenthal A, Crudu V, Homorodean D, Hsueh P-R, Narayanan S, Pym AS, Skrahina A, Swaminathan S, Van der Walt M, Alland D, Bishai WR, Cohen T, Hoffner S, Birren BW, Earl AM. 2017. Genomic analysis of globally diverse *Mycobacterium tuberculosis* strains provides insights into the emergence and spread of multidrug resistance. *Nat Genet* 49:395–402. <https://doi.org/10.1038/ng.3767>.
 46. Billington OJ, McHugh TD, Gillespie SH. 1999. Physiological cost of rifampin resistance induced in vitro in *Mycobacterium tuberculosis*. *Antimicrob Agents Chemother* 43:1866–1869.
 47. Eilertson B, Maruri F, Blackman A, Herrera M, Samuels DC, Sterling TR. 2014. High proportion of heteroresistance in *gyrA* and *gyrB* in fluoroquinolone-resistant *Mycobacterium tuberculosis* clinical isolates. *Antimicrob Agents Chemother* 58:3270–3275. <https://doi.org/10.1128/AAC.02066-13>.
 48. Habzón MH, Bobadilla del Valle M, Guerrero MI, Varma-Basil M, Filliol I, Cavatore M, Colangeli R, Safi H, Billman-Jacobe H, Lavender C, Fyfe J, García-García L, Davidow A, Brimacombe M, León CI, Porras T, Bose M, Chaves F, Eisenach KD, Sifuentes-Osornio J, Ponce de León A, Cave MD, Alland D. 2005. Role of *embB* codon 306 mutations in *Mycobacterium tuberculosis* revisited: a novel association with broad drug resistance and IS6110 clustering rather than ethambutol resistance. *Antimicrob Agents Chemother* 49:3794–3802. <https://doi.org/10.1128/AAC.49.9.3794-3802.2005>.
 49. DeJesus MA, Gerrick ER, Xu W, Park SW, Long JE, Boutte CC, Rubin EJ, Schnappinger D, Ehrt S, Fortune SM, Sasseti CM, Iøerger TR. 2017. Comprehensive essentiality analysis of the *Mycobacterium tuberculosis* genome via saturating transposon mutagenesis. *mBio* 8:e02133-16. <https://doi.org/10.1128/mBio.02133-16>.
 50. Shekar S, Yeo ZX, Wong JCL, Chan MKL, Ong DCT, Tongyoo P, Wong S-Y, Lee ASG. 2014. Detecting novel genetic variants associated with isoniazid-resistant *Mycobacterium tuberculosis*. *PLoS One* 9:e102383. <https://doi.org/10.1371/journal.pone.0102383>.
 51. Danilchanka O, Mailaender C, Niederweis M. 2008. Identification of a novel multidrug efflux pump of *Mycobacterium tuberculosis*. *Antimicrob Agents Chemother* 52:2503–2511. <https://doi.org/10.1128/AAC.00298-08>.
 52. Fu LM, Shinnick TM. 2007. Genome-wide exploration of the drug action of capreomycin on *Mycobacterium tuberculosis* using Affymetrix oligonucleotide GeneChips. *J Infect* 54:277–284. <https://doi.org/10.1016/j.jinf.2006.05.012>.
 53. Phong TQ, Ha DTT, Volker U, Hammer E. 2015. Using a label free quantitative proteomics approach to identify changes in protein abundance in multidrug-resistant *Mycobacterium tuberculosis*. *Indian J Microbiol* 55:219–230. <https://doi.org/10.1007/s12088-015-0511-2>.
 54. Phelan J, Coll F, McNeerney R, Ascher DB, Pires DEV, Furnham N, Coeck N, Hill-Cawthorne GA, Nair MB, Mallard K, Ramsay A, Campino S, Hibberd ML, Pain A, Rigouts L, Clark TG. 2016. *Mycobacterium tuberculosis* whole genome sequencing and protein structure modelling provides insights into anti-tuberculosis drug resistance. *BMC Med* 14:31. <https://doi.org/10.1186/s12916-016-0575-9>.
 55. Andrews S. 2012. FastQC. Babraham Bioinformatics, Babraham Institute, Cambridge, United Kingdom.
 56. Kreuger F. 2013. TrimGalore! Babraham Bioinformatics, Babraham Institute, Cambridge, United Kingdom.
 57. Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. ArXiv:1303.3997 [q-bio.GN]
 58. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>.
 59. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman J, Young SK, Earl AM. 2014. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One* 9:e112963. <https://doi.org/10.1371/journal.pone.0112963>.
 60. Price MN, Dehal PS, Arkin AP. 2010. FastTree 2 – approximately maximum-likelihood trees for large alignments. *PLoS One* 5:e9490. <https://doi.org/10.1371/journal.pone.0009490>.
 61. Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, Harris SR. 2016. SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments. *Microb Genom* 2:e000056. <https://doi.org/10.1099/mgen.0.000056>.
 62. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w11118; iso-2; iso-3. *Fly* 6:80–92. <https://doi.org/10.4161/fly.19695>.
 63. Salazar A, Earl A, Desjardins C, Abeel T. 2015. Normalizing alternate representations of large sequence variants across multiple bacterial genomes. *BMC Bioinformatics* 16:A8. <https://doi.org/10.1186/1471-2105-16-S2-A8>.
 64. De Mita S, Siol M. 2012. EggLib: processing, analysis and simulation tools for population genetics and genomics. *BMC Genet* 13:27. <https://doi.org/10.1186/1471-2156-13-27>.
 65. R Core Team. 2015. R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.
 66. Weir BS, Cockerham CC. 1984. Estimating F-statistics for the analysis of population structure. *Evolution* 38:1358–1370. <https://doi.org/10.1111/j.1558-5646.1984.tb05657.x>.
 67. Sagulenko P, Puller V, Neher R. 2017. TreeTime: maximum likelihood phylodynamic analysis. *bioRxiv* <https://doi.org/10.1101/153494>.