

## **Prediction of pyrazinamide resistance in *Mycobacterium tuberculosis* using structure-based machine learning approaches.**

Joshua J Carter<sup>1</sup>, Timothy M Walker<sup>1</sup>, A Sarah Walker<sup>1,2</sup>, Michael G. Whitfield<sup>3,\*</sup>, Timothy EA Peto<sup>1,2</sup>, Derrick W Crook<sup>1,2,4</sup>, Philip W Fowler<sup>1,2</sup>

<sup>1</sup> Nuffield Department of Medicine, University of Oxford, John Radcliffe Hospital, Headley Way, Oxford, OX3 9DU, UK

<sup>2</sup> National Institute of Health Research Oxford Biomedical Research Centre, John Radcliffe Hospital, Headley Way, Oxford, OX3 9DU, UK

<sup>3</sup> Division of Molecular Biology and Human Genetics, SAMRC Centre for Tuberculosis Research, DST/NRF Centre of Excellence for Biomedical Tuberculosis Research, Faculty of Medicine and Health Sciences, Stellenbosch University, Tygerberg, South Africa

<sup>4</sup> NIHR Health Protection Research Unit in Healthcare Associated Infection and Antimicrobial Resistance at University of Oxford in partnership with Public Health England, Oxford, UK

\* for the "EXIT-RIF" investigators: Prof Robin M Warren, Prof Annelies van Rie, Prof Lesley Scott, Prof Wendy Stevens

Correspondence: [philip.fowler@ndm.ox.ac.uk](mailto:philip.fowler@ndm.ox.ac.uk), [@philipwfowler](https://twitter.com/philipwfowler)

## Abstract

Pyrazinamide is one of four first-line antibiotics currently used to treat tuberculosis and has been included in newer treatment regimens undergoing clinical trials due to its unique sterilizing effects and synergy with newer drugs. However, phenotypic antibiotic susceptibility testing for pyrazinamide is problematic. Resistance to pyrazinamide is primarily driven by genetic variation in *pncA*, which encodes PncA, an enzyme that converts pyrazinamide into its active form. We curated a derivation dataset of 291 non-redundant, missense amino acid mutations in PncA with associated high-confidence phenotypes from studies of clinical isolates and *in vitro/in vivo* screening studies and then trained machine learning models to predict pyrazinamide resistance based on sequence- and structure-based features of each missense mutation. The clinical relevance of the models was tested by predicting the binary resistance phenotype of 2,292 clinical isolates harboring missense mutations in PncA to pyrazinamide. The probabilities of resistance predicted by the model were also compared with *in vitro* pyrazinamide minimum inhibitory concentrations of 27 isolates to determine whether the machine learning model could predict the degree of resistance. Finally, we predicted the effect on pyrazinamide resistance of the remaining 814 possible missense mutations caused by single nucleotide polymorphisms in PncA that have not yet been observed in public databases. Overall, this work offers an approach to improve the sensitivity and specificity of pyrazinamide resistance prediction in genetics-based clinical microbiology workflows for tuberculosis, highlights novel mutations for future biochemical investigation, and is a proof of concept for using this approach in other drugs such as bedaquiline.

## Introduction

*Mycobacterium tuberculosis* is an evolutionarily ancient human pathogen that is the leading cause of death by infectious disease worldwide<sup>1</sup>. In 2016, tuberculosis was responsible for 1.7 million deaths and 10.4 million new infections<sup>1</sup>. Tuberculosis control efforts have been hampered by the evolution of resistance to antibiotics, threatening the efficacy of the standard four drug antibiotic therapy consisting of rifampicin, isoniazid, ethambutol, and pyrazinamide<sup>1,2</sup>. Pyrazinamide plays a critical role in tuberculosis treatment through its specific action on slow-growing, “persister” bacteria that often tolerate other drugs due to their reduced metabolism<sup>3-6</sup>. This unique activity has been instrumental in shortening the standard treatment duration to six months, substantially increasing the effectiveness of antibiotic therapy<sup>5,6</sup>. Numerous studies have also found that including pyrazinamide in treatment regimens increases sputum-conversion rates in both pan-susceptible and multi-drug resistant (MDR, defined as resistant to rifampicin and isoniazid) tuberculosis populations<sup>7</sup>. Due to its unique sterilizing effect and its synergy with new tuberculosis drugs such as bedaquiline, pyrazinamide is also included in most new treatment regimens targeting drug-resistant tuberculosis<sup>8-13</sup>. Therefore, accurately and quickly determining whether a clinical isolate is resistant to pyrazinamide is critically important.

The majority of culture-based laboratory methods to determine pyrazinamide resistance are technically challenging, requiring highly-trained technicians. Even then, results are often not reproducible, meaning these methods are rarely employed in low-resource and/or high-burden clinical settings<sup>14</sup>. Even the current WHO-endorsed gold standard, the Mycobacteria Growth Indicator Tube (MGIT), which is relatively simple to use, suffers from low precision, with false-resistance rates of up to 68% reported<sup>15-22</sup>. As the prevalence of multi-drug resistant (MDR) and extensively drug-resistant (XDR) TB increases, this lack of precision becomes more of a problem and hence the WHO has recommended moving toward genetics-based approaches for all first-line antibiotics.

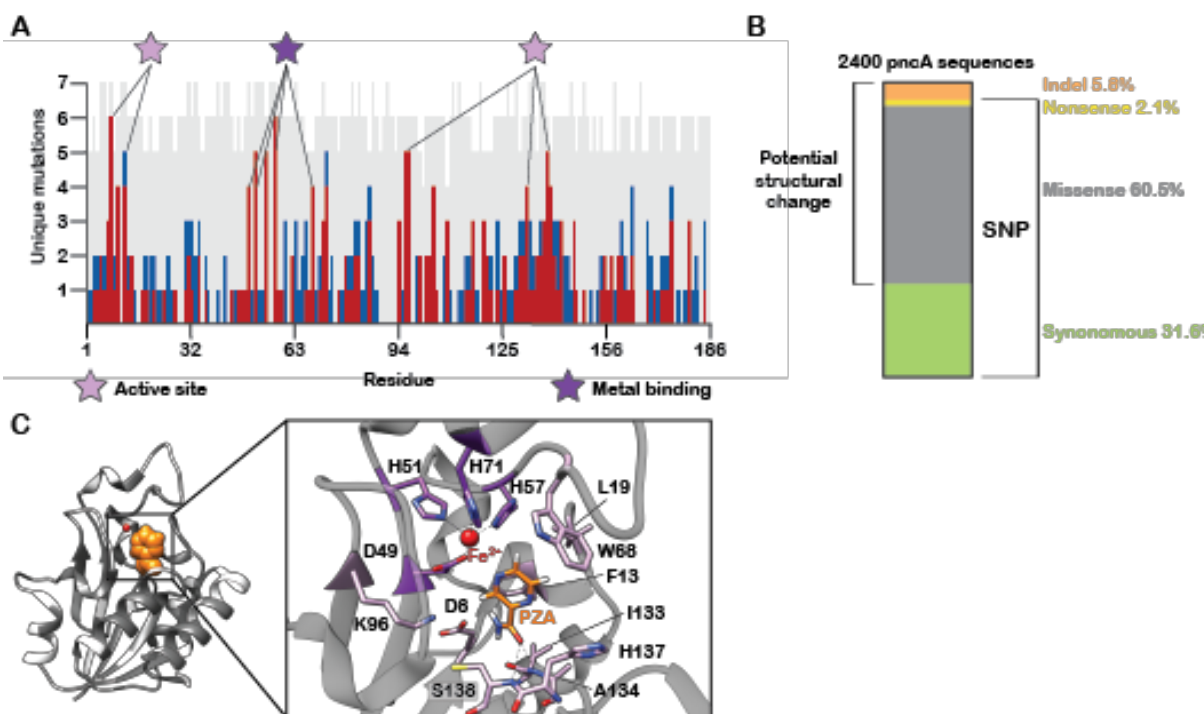
Despite pyrazinamide being used to treat tuberculosis since 1952, less is known about its genetic determinants of resistance compared to other first-line drugs<sup>5</sup>. Resistance to rifampicin or isoniazid is conferred in most isolates (90-95% and 50-97%, respectively) by a small number of highly-penetrant mutations in a small and well-delineated region of a single gene (*rpoB* and *katG*, respectively)<sup>4</sup>. Pyrazinamide is a pro-drug that is converted to its active form of pyranic acid by the action of PncA, a pyrazinamidase/nicotinamidase encoded by the *pncA* gene<sup>24</sup>. While other genetic loci have been implicated in pyrazinamide resistance (notably *rpsA*, *panD*, and the putative efflux pumps *Rv0191*, *Rv3756c*, *Rv3008*, and *Rv1667c*), the majority (70-97%) of pyrazinamide-resistant clinical isolates harbor mutations in either the promoter region or coding sequence of the *pncA* gene<sup>25-33</sup>. In contrast to the well-delineated and relatively restricted “resistance-determining regions” found in *rpoB* (rifampicin, 27 codons) and *katG* (isoniazid, single codon), pyrazinamide-resistant mutations have been identified along the entire length of the *pncA* gene (**Figure 1A**) with no single mutation dominating.

The consequence of this is that whilst line-probe assays have been successfully developed to predict resistance to rifampicin and isoniazid, and these are quicker and more reliably detect resistance than culture-based methods<sup>15,23</sup>, it is much challenging to develop a line probe assay targeting specific regions of *pncA* to predict pyrazinamide resistance. One instead is forced to consider whole-genome sequencing, however here the diversity of resistance-conferring mutations fundamentally limits the sensitivity and specificity of heuristic approaches that aim to predict the effectiveness of pyrazinamide based on a catalogue of previously-observed genetic variants<sup>4,14,30,31,34-36</sup>.

Dataset	Phenotype	Isolate Sources	# Isolates	# Unique Missense Mutations
Exploratory	R/S/F	Ref <sup>37</sup>	2,651	254
Derivation	R/S	Ref <sup>31,37-39</sup>	1,792 + <i>in vitro</i> isolates <sup>39</sup>	291
Clinical	R/S	Ref <sup>31,37,38</sup>	2,292	272
Quantitative	MIC	EXIT-RIF Study	366	27
Prevalence	None	European Nucleotide Archive and ref <sup>37</sup>	25,986	376

**Table 1.** Description of datasets employed in this study. (R=resistant to antibiotic, S=susceptible, F=test failed to return a result)

An in depth analysis of the genetic variation reported in the initial exploratory dataset revealed 2,651 strains with mutations in the open reading frame of *pncA*, a substantial majority (2,400) of which only harbored a single mutation (**Table 1**)<sup>37</sup>. Of these, substitutions represented the majority (94.2%) of the genetic variation, with insertion, deletion, and frameshift mutations making up the remainder (5.8%). Insertions, deletions, and frameshift mutations, along with nonsense substitutions (present in 7.9% of the single mutation strains), were all associated with pyrazinamide resistance, consistent with their likely disruption of the PncA enzyme, thereby preventing pyrazinamide activation. Synonymous substitutions (present in 31.6% of the single mutation strains, **Figure 1B**) were not associated with resistance, suggesting that they can be inferred to not alter the pyrazinamide susceptibility of a particular strain. Thus, nonsynonymous substitution mutations (present in 60.5% of single mutation strains) represent the majority of the potential resistance-causing variation in *M. tuberculosis*.



**Figure 1. Distribution of *pncA* mutations from published datasets. (A)** Barplot of the impact of possible missense mutations in PncA by amino acid position. High confidence resistant (red) and susceptible (blue) mutations are overlaid on the possible missense mutations whose effect on resistance is unknown (grey). **(B)** Distribution of the types of mutations reported by the CRYPTIC consortium *et al*<sup>37</sup>. **(C)** Missense mutations from the dataset plotted onto the PncA structure (PDB ID: 3PL1) in dark grey. A pyrazinamide molecule (orange) has been modeled into the active site.

A genetics-based clinical microbiology for tuberculosis depends on being able to predict or infer the effect of any likely observable *pncA* mutation on pyrazinamide susceptibility. Recent studies to identify pyrazinamide-resistance determining mutations have focused on either classifying mutations from previously observed clinical isolates or discovering novel mutations through *in vitro/in vivo* screening approaches<sup>31,37-39</sup>. However, these strategies are limited, respectively, by the relative paucity of sequenced clinical isolates and the lack of laboratory capacity to systematically generate and test mutants. A computational modelling approach could potentially predict the effect of a significant number of missense mutations before they are

observed in clinical isolates, allowing clinicians to more rapidly make an informed decision in the face of emerging resistance patterns as well as focusing future fundamental *in vitro* studies on the most important mutations to investigate.

In this paper we demonstrate that machine-learning models can robustly and accurately predict the effect of nonsynonymous substitutions in *pncA* on pyrazinamide susceptibility. The models were trained using a derivation dataset of 291 non-redundant, nonsynonymous PncA amino acid mutations collected from pooled MGIT phenotypic studies and a comprehensive *in vitro/in vivo* pyrazinamide-resistance screen (**Methods, Table 1, S1**)<sup>31,37–39</sup>. This dataset reflects the clinically observed distribution of mutations across the *pncA* gene (**Figure 1A**), as well as throughout the protein structure. Since the *pncA* gene is not essential, our hypothesis is that nonsynonymous substitutions mutation confer resistance by abrogating the folding, stability or function of the PncA protein. This led us to consider how each mutation perturbs the local chemistry and overall structure of the protein, Hence, information about the structural and chemical properties and evolutionary context of the wild type and mutant amino acids in question were used as inputs to several different machine-learning models (**Methods**). The predictions from the best performing model were then re-applied to an aggregated clinical dataset to examine their clinical relevance and also checked against a smaller quantitative dataset of *in vitro* pyrazinamide minimum inhibitory concentrations (MICs) to determine their capacity to predict the degree of resistance for specific mutations (**Table 1**). Finally, the model was used to make predictions for *all* (1105) nonsynonymous substitutions possible from single point mutations in *pncA*. These data were then used to predict the occurrence of pyrazinamide resistance in a prevalence dataset derived from public sequence repositories (**Table 1**). This study is a proof of principle for using computational approaches to model and predict antibiotic resistance in other drugs, such as bedaquiline, pretomanid/delamanid, isoniazid, and ethionamide, where some genes implicated in resistance pathways appear to be non-essential.

## Results

### *Structural and evolutionary traits correlate with mutational impact on pyrazinamide susceptibility*

We built a preliminary set of 722 nonsynonymous substitutions in PncA that had either been observed multiple times in clinical isolates for which antibiotic susceptibility testing data was available or were generated in a high-precision laboratory screening study of *pncA* resistance variants (**Methods**). Discarding those mutations where there was uncertainty around whether they conferred resistance (or not) left us with a final derivation dataset of 291 substitutions (**Table 1**). To understand the structural features that determine a mutation's effect on pyrazinamide susceptibility, we mapped our derivation dataset onto the PncA structure. No obvious clustering was revealed, consistent with the previously observed distribution of resistant mutations across the gene sequence and protein structure (**Figure 1A,C**)<sup>14,30,31,35,39</sup>. Interestingly, there were a significant number of PncA codons where mutations associated with resistance and mutations consistent with susceptibility were seen, suggesting that the change in local chemistry introduced by the mutant amino acid is an important factor in determining resistance (**Figure 1A**). The amino acid positions with the highest mutational diversity in the dataset were all residues involved in active site formation or metal binding, suggesting that, consistent with our hypothesis, loss or alteration of these functions is a common mechanism for gaining pyrazinamide resistance. Indeed, previous studies have noted a negative correlation between a mutation's distance from the active site and its tendency to cause resistance (**Figure S1**)<sup>31,39,40</sup>.

Examining the PncA structure also suggested that resistant mutations were more likely to be buried in the protein core, consistent with findings from previous *in vitro* and *in vivo* screens (**Figure 2A**)<sup>31,39</sup>. Mutations in the hydrophobic core of a protein are likely to be destabilizing<sup>41-44</sup>. Indeed, some pyrazinamide-resistant mutations result in reduced production of functional PncA, perhaps due to impaired protein folding/stability<sup>39,45</sup>. To assess a mutation's impact on the stability of PncA, we employed a meta-predictor that calculates the predicted

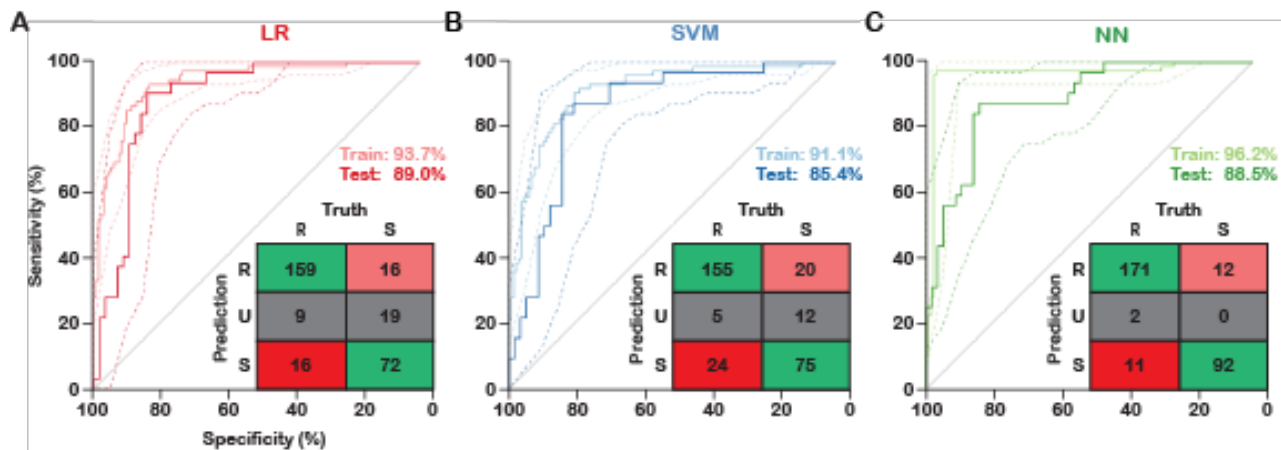


change in free energy of protein unfolding *in silico*<sup>46</sup>. This is a fast, heuristic method; although other more accurate methods exist, these require vastly more computational resource<sup>47</sup>.

In addition to these chemical and structural properties, we also included information on the evolutionary variation at each position obtained from multiple sequence alignments of related orthologs (**Methods**). Unsurprisingly, increased conservation at a position was associated with a higher potential of a mutation at that position to confer resistance (**Figure S1**). Finally, we applied a recent computational method, called MAPP, that quantifies the evolutionary constraints imposed on a given position in a protein. MAPP does this by combining the range of physicochemical amino acid properties observed at a particular position in a multiple sequence alignment with weights generated from the branch lengths of a phylogenetic tree<sup>48</sup>. Resistant mutations had significantly higher MAPP scores, indicating that resistance-conferring mutations in PncA are less conservative in amino acid chemistry and function (**Figure 2B, S1**).

#### *Machine-learning models accurately predict pyrazinamide resistance*

Univariable logistic regression over the derivation dataset revealed that most of the individual predictors were associated with resistant phenotypes (**Table S2, Figure 2F**). The MAPP score and solvent accessibility proved to be the most discriminatory individual features. As PncA can be inactivated through defects in protein folding, reduced stability, distortion of active site geometry, abrogation of metal binding, or some combination of these, we expected a machine-learning approach to be ideally suited to simultaneously consider all these possible mechanisms of PncA inactivation, and hence more accurately predict pyrazinamide resistance/susceptibility.



**Figure 3: Machine learning models predict pyrazinamide resistance from structural features.**

Performance of (A) logistic regression (LR), (B) support vector machine with radial kernel (SVM) and (C) neural network (NN) models for prediction of pyrazinamide resistance. Dotted lines represent 95% confidence intervals from bootstrapping ( $n=10,000$ ) and the area under the curve is reported for training and testing sets. Truth tables are shown for the combined training and test sets.

To evaluate the different models, we randomly divided our derivation dataset mutations (184 resistant, 107 susceptible) into a 70% training set and a 30% validation set, preserving the overall distribution of resistant and susceptible mutations. Models were then trained using repeated 10-fold cross validation (**Methods**). Since the models output a probability of resistance between 0 and 1, we defined three results; resistant (R,  $p < 0.4$ ), susceptible (S,  $p > 0.6$ ) and uncertain (U,  $0.4 \leq p \leq 0.6$ ) (**Figure S2A**). The models were able to call ~87-99% (183-190) of the mutations in the training set using these thresholds. As expected, drops in performance were observed for all models when applied to the independent validation set, however none were statistically significant (**Figure 3, Table S3**). The neural network (NN) model had the highest diagnostic odds ratio (119), followed by logistic regression (LR, 45) and then the support vector machine (SVM, 24, **Figure 3**). As the best performing model, the predictions from the neural network model were used for all further analyses.

### *Analysis of model errors on the derivation set*

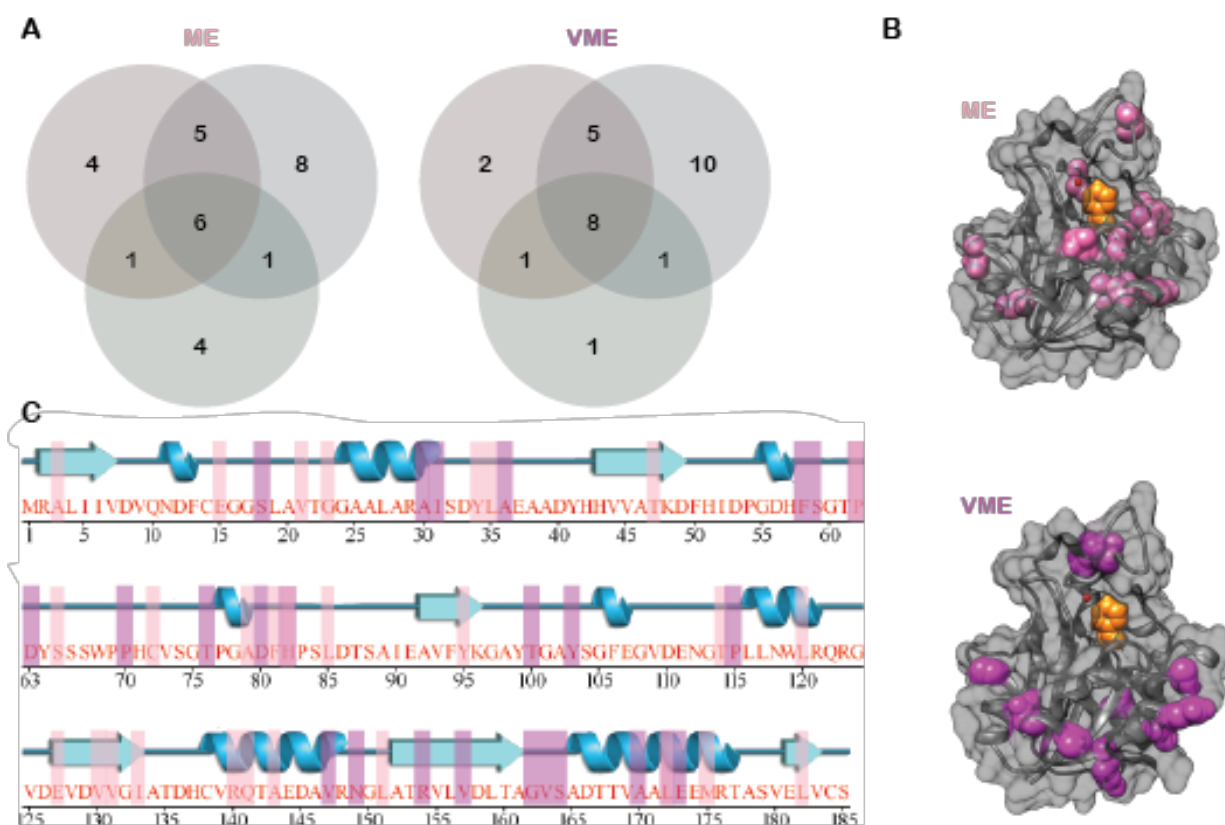
Clinical diagnostic errors for antibiotic resistance are categorized into three classes: very major errors, which represent truly resistant isolates that are called susceptible, major errors (true susceptible cases called resistant), and minor errors, which are not called by the method being tested but are determined as resistant or susceptible by the reference method<sup>49</sup>.

Collectively, the three models made 57 incorrect calls (28 very major errors and 29 major errors); however, only 14 of these were shared between all three models (8 very major errors and 6 major errors, **Figure 4A, Table S4**). The best performing model (neural network) had a sensitivity of 94% (89-97%), a specificity of 88% (81-94%), and a positive predictive value of 93% (89-96%) (**Table S3**).

Although the mutations responsible for the very major errors (predicted susceptible, phenotypically resistant) of the neural network model were dispersed throughout the protein structure, most (10/11) were surface exposed (**Figure 4B**). All these mutations were predicted to either not affect or slightly increase the stability of PncA, suggesting these errors may be due to inaccuracies in the predicted free energy change of unfolding (**Figure S2**). Major errors (predicted resistant, phenotypically susceptible) were typically driven by overestimation of a mutation's potential to effect PncA structure or function. Mapping the mutations responsible for these major errors onto the sequence and structure revealed one cluster near the active site and the coiled turn between the  $\alpha-1$  and  $\alpha-2$  helix (residues 15, 21, 23, 131, and 133) (**Figure 4B, C**). Major errors appear to be caused in part by a combination of overestimation of the MAPP score at invariant positions that are near the active site or buried in the protein core.

Interestingly, several major errors occurred in a region of the active site termed the "oxyanion hole" (residues 131-138) which coordinates the carbonyl group of pyrazinamide in the active site<sup>40</sup>. The effects of mutations that lie in this region could be over-estimated due to their proximity to the active site and relatively lower solvent exposure. As the interaction between the oxyanion hole and pyrazinamide is mediated by the peptide backbone, and is therefore

sidechain-independent, there is likely to be less stringent selection of the residues at these positions as long as the overall peptide backbone structure is maintained. Interestingly however, Gly132 and Ala134 are invariant in the alignment used to generate the MAPP score, which would suggest that these sites are however under strong selective pressures. It has been shown previously that specific residues are favored in the positions surrounding *cis* peptide bonds, so future work could attempt to model the mutations occurring in this functional region more specifically<sup>50</sup>.



**Figure 4. Very major errors are concentrated on the surface of PncA.**

(A) Major (ME) and very major (VME) errors shared between the models. Errors from logistic regression are shown in red, support vector machines in blue, and neural network in green (B) PncA with major (pink) and very major (magenta) errors shown as spheres. (C) Major (pink) and very major (magenta) errors mapped onto the *pncA* primary sequence.

Minor errors are cases where the model could not confidently call a mutation resistant or susceptible and represented 34% of errors made by all three models collectively. The model features for mutations that were called U tended to be intermediate compared to those of resistant and susceptible isolates, which raises the intriguing possibility that these are mutations with an intermediate effect on protein stability and/or enzyme activity (**Figure S2**).

As the neural network model does not clearly indicate which PncA features drive its predictions, we used logistic regression with backwards elimination (**Methods**) on the derivation dataset to gain further insight into the complex interplay between these factors. This analysis revealed that solvent accessibility, distance to the active site, the evolutionary conservation of the wild type amino acid, the number of hydrogen bonds formed by the wild type residue, significant changes in protein stability (measured as a change in free energy of protein unfolding of >2 kcal/mol), and the MAPP score were all independent explanatory factors. In addition, the interactions between both the MAPP score and number of hydrogen bonds were found to moderate the effect of the distance to the active site. A higher MAPP score increased the deleterious effect of a mutation near the active site, while the importance of the number of hydrogen bonds a residue was involved in decreased further from the active site. This may be due to the requirement of hydrogen bonding interactions for proper active site geometry. There was weak evidence for two other interactions (between protein destabilization and either the number of hydrogen bonds or solvent accessibility,  $p=0.073$  and  $p=0.054$  respectively). These results suggest that interactions between model features are important for prediction of resistance, which may be why the neural network model outperforms logistic regression in classifying mutations.

		2292 strains		272 mutations		
		MGIT Phenotype		Consistent MGIT Phenotype		
		R	S	R	I	S
Prediction	R	1181 [437]	191 [133]	92 [92]	21 [14]	4 [4]
	U	287 [150]	134 [76]	26 [26]	17 [10]	2 [2]
	S	125 [75]	374 [325]	11 [11]	13 [6]	12 [12]

**Figure 5. Model predictions based on mutations generalize to MGIT phenotypes.** (A) Truth table of the predictions' performance on a dataset of *M. tuberculosis* strains tested by MGIT. Brackets denote predictions based on missense mutations not in the training set. (B) Truth table of model predictions versus average mutation phenotypes. "I" is defined as mutations that are not R or S in >75% of isolates tested (n≥4). Brackets denote phenotypes for mutations for which there was enough clinical evidence to be confident in the assigned phenotype (177/272, 65%, **Methods**).

#### *Neural network predictions generalize to a large clinical dataset*

To assess the generalizability of our best model, the neural network, we applied it to a clinical dataset of 2,292 *pncA* gene sequences with MGIT antibiotic susceptibility results (**Table 1**), each representing a unique isolate collated from published studies of clinical isolates<sup>31,37–39</sup>. In addition to the clinical isolates that were used in the derivation dataset, this dataset also included 500 isolates either where mutations were only seen once or isolates with mutations whose phenotype could not be called confidently. 1,196 isolates (52%) from the clinical dataset harboring mutations were not used in model training. Predicting resistance/susceptibility based on the mutation present, the model correctly predicted 74.1% (1,181/1,593) of MGIT-resistant isolates with *pncA* mutations; however, it performed more poorly for MGIT-susceptible isolates, predicting only 53.5% (374/699) of strains (**Figure 5A**). 68.1% (287/421) of U calls made by the

model were associated with resistance, which suggests that clinically a U call could be interpreted as possible pyrazinamide resistance contingent upon further testing (**Figure 5A**).

Intriguingly, 30 mutations (present in 385 strains, 16.8% of strains classified) had variable MGIT results (defined as having a resistant MGIT phenotype in 25-75% of cases with at least 4 observations). To understand our model performance with these possible “intermediate” mutations, we compared the model predictions with the average phenotypes of the 272 unique mutations found in these strains, which we classified as resistant (MGIT resistant in >75% of isolates with this mutation tested), susceptible (MGIT susceptible in >75% of isolates with this mutation tested) and intermediate (I), the remainder. The model classified 14 (47%) of these high-confidence “intermediate” mutations as resistant, leading to 76 strains (40% of major errors made by the model) with susceptible MGIT phenotypes being misclassified as resistant (**Figure 5B**). While the model predicted U for 10 (33%) of the intermediate mutations, there was no clear relationship between U calls and intermediate mutations, which is consistent with the fact that the model was trained on binary data (**Figure 5B**). Overall, the model more accurately predicted the average phenotypes of mutations than the individual MGIT phenotypes of clinical isolates, with the exception of the 11 very major errors; however, some of these errors may be due to errors in MGIT antibiotic susceptibility testing. Alternatively, these errors could result from resistance that is determined by factors upstream of protein folding and function and is therefore outside the scope of our model.

#### *Model predictions are correlated with pyrazinamide minimum inhibitory concentrations in vitro*

To test the model’s capacity to predict the degree of pyrazinamide resistance conferred by a particular mutation, we compared the calls and predicted probabilities of our model with the minimum inhibitory concentrations (MICs) of pyrazinamide for a set of *M. tuberculosis* isolates collected in South Africa (quantitative dataset, **Table 1, Methods**). Overall, our model correctly predicted the binary (R/S) phenotype for 25 of 27 single missense *pncA* mutations observed in

these isolates (**Figure S3, Table S6**). One of these errors is likely due to a MGIT testing error, as the same mutation (Gly97Asp) is observed in another resistant isolate and has been classified as conferring resistance in other studies. We also compared the predicted probabilities of our model with the MICs to determine if our model was also informative about the relative degree of resistance conferred by a particular mutation. The predicted probabilities were moderately negatively correlated ( $r = -0.49$ ). In future work we will test more isolates over a greater range of pyrazinamide MICs to investigate whether our model can also predict the degree of resistance at the upper end of the resistance spectrum. In addition, several samples with large deletions in *pncA* were also observed; these strains had extremely high (>900  $\mu\text{g/mL}$ ) MICs, which is consistent with the loss of functional PncA protein (**Table S6**). These data taken together with the results for isolates harboring insertions/deletions in the study by CRyPTIC Consortium *et al*<sup>37</sup> confirm it is reasonable to assume all large insertion/deletion mutations and frameshifts in *pncA* confer resistance to pyrazinamide.

#### *Predicting the effect of all possible nonsynonymous pncA mutations on pyrazinamide susceptibility*

Since trained machine learning models require very little computational resource, we applied our model to every possible nonsynonymous SNP in *pncA* (coding for 1,105 unique amino acid changes, 814 previously unobserved missense mutations), thereby estimating the probability that each mutation confers pyrazinamide resistance (**Table S5**). Overall, 22% (244) of missense mutations were predicted to confer resistance, while 63% (691) were predicted to have no effect on the action of pyrazinamide and the remaining 14% (158) were predicted to have an uncertain effect. Interestingly, the proportion of predicted resistant mutations was much lower than that in the derivation set (22% versus 63% respectively). This may be due to an increased likelihood of sequencing pyrazinamide-resistant clinical isolates, leading to an over-representation of resistance-conferring mutations in our derivation dataset as opposed to



susceptible ones. As more unselected studies of whole genome sequencing are conducted, we expect this bias to unwind and consequently more susceptible mutations will be found than perhaps expected for most established drugs. Alternatively, it could be caused by a global underestimation of resistance by our model, which underpredicted resistance by ~10% in the clinical dataset. Finally, this difference could represent the fact that phenotypically intermediate mutations classified as U/I by our model are classified as resistant in the catalogs/screens used to develop the dataset.

In order to understand how these predictions improve our capacity to identify resistant mutations in *pncA*, we queried a bacterial index of the European Nucleotide Archive (ENA) to identify all missense SNPs in the *M. tuberculosis pncA* coding sequence (**Table S5**)<sup>51</sup>. We found ~15,777 *pncA* sequences classified as originating from *M. tuberculosis*, with 4,504 strains harboring missense mutations in *pncA*. We supplemented this ENA dataset with the *pncA* sequences collected by CRYPTIC *et al*<sup>67</sup> as these sequences were not deposited when the index of the ENA was built, bringing the total number of strains with missense mutations to 7,107 (prevalence dataset, **Table 1**). Out of the 376 unique missense mutations observed in the prevalence dataset, 13 were frequent (>100 sequences) and 11 of these were associated with resistance (**Table S5**). 236 (63%) mutations were observed 10 or fewer times and 69 (18%) only once, highlighting the need for approaches capable of predicting the effect of rare missense mutations. We classified the prevalence dataset using a published heuristic catalog<sup>52</sup>, supplemented with our resistant and susceptible model predictions, to quantify how much our machine learning model improves our capacity to screen for potential pyrazinamide resistance using whole genome sequencing. While the heuristic catalog alone was able to classify 5,321 strains (75%, 291 mutations), our model classified an additional 1,147 strains (16%, 47 mutations), allowing us to provisionally classify 91% of the strains with missense mutations in the prevalence dataset.

Of the 4,504 strains with *pncA* missense mutations obtained from the ENA, 77% (3,458 strains) were predicted to be resistant, 14% (640 strains) were predicted susceptible, and 9% (406 strains) were called uncertain. Interestingly, the predicted prevalence of resistance amongst strains harboring missense mutations in *pncA* in the ENA (77%) was higher than the prevalence observed in clinical isolates from CRYPTIC *et al* (70%) and substantially higher than the predicted overall prevalence possible for *pncA* (22%). This may be due to a historic sampling bias that preferentially selects resistant isolates for sequencing; alternatively, it could represent a bona fide enrichment of resistant isolates driven by the selective pressure of antibiotic treatment.

## Discussion

*De novo* prediction of 814 mutations' effects on pyrazinamide resistance constitutes a significant step forward in our ability to predict pyrazinamide resistance from genetics and a proof of concept for using structural approaches to infer the effects of missense mutations on pyrazinamide resistance. While improvements to the model are necessary to achieve the sensitivity and specificity required for routine clinical use, this work increases our ability to identify rare resistance mutations, thereby potentially increasing the capability of whole-genome sequencing based diagnostic susceptibility testing to respond to emerging and rare resistance patterns, as well as prioritizing rare resistance mutations for *in vitro* validation. Additionally, improving the classification of susceptible *pncA* mutations will allow us to begin to disentangle the involvement of other genes in pyrazinamide resistance, including determining the effect of mutations in non-standard pyrazinamide resistance-associated genes such as *panD* and *rpsA*.

A principal limitation of this approach is that it can only make predictions for missense mutations in the coding sequence of *pncA*. While we have shown that these represent most (60.5%) of the possible resistant genetic variants in *pncA*, insertions/deletions and nonsense mutations (7.9%) must also be considered, as they are generally associated with resistance.

Likewise, promoter mutations that result in reduced transcription of *pncA* will likely also lead to resistance. While no synonymous mutation has yet been observed to cause pyrazinamide resistance to date, the possibility remains that a synonymous mutation could have an effect on mRNA stability, ribosomal stalling, or codon usage and confer resistance. The model also does not take into account the introduction of protease cleavage sites or other processing abnormalities. Finally, while most pyrazinamide resistance is caused by mutations in *pncA*, recent studies have also implicated other genes, notably *rpsA*, *panD*, and the putative efflux pumps *Rv0191*, *Rv3756c*, *Rv3008*, and *Rv1667c* in pyrazinamide resistance<sup>4,25-32</sup>. Further research is needed to determine if mutations in these genes can be reliably inferred to confer pyrazinamide resistance.

Several predictive features used in the model could be improved upon in future work. The MAPP score relies on the maintenance of function between diverse homologs to determine the evolutionary constraints on each position in a protein. While we selected sequences that contained the residues involved in active site formation and metal binding, we did not experimentally confirm pyrazinamide conversion by each homolog. Additionally, the *in silico* method that we employed to estimate each mutation's effect on protein stability could be improved by comparison and calibration with *in vitro* biochemical data. Finally, as the active site of PncA is formed in part by a *cis* peptide bond between Ile133 and Ala134, more detailed modeling of the evolutionary constraints at this site could more accurately assess the functional impact of a mutation at these positions. Despite the fact that most features we investigated were associated with pyrazinamide resistance, not all were retained as independent predictors in our final logistic regression model. The change in hydrophathy and sidechain volume as well as the Rogov score are all likely encapsulated by the MAPP score, as this takes into account the wild type and mutant amino acids in its calculation.

The accuracy of model predictions based on structural features suggests that the underlying hypothesis of predicting pyrazinamide resistance based on predicted PncA function

is valid. Mapping the potential of each position to harbor a predicted or bona fide resistance mutation onto the structure reveals that many resistance-prone positions are associated with the active site or metal binding, as noted previously (**Figure 1C**). Interestingly, however, most of the other resistance-prone positions are involved in packing interactions between secondary structure elements in PncA, supporting the assertion that a major mechanism of pyrazinamide resistance is loss of protein stability. All susceptibility-prone positions are highly solvent-exposed and many are on flexible loops, consistent with our expectation that these regions experience lower selective pressures and have a lower/negligible effect on PncA stability and function. The effect of perturbing the protein core appears to be more dependent on the specific amino acid chemistries involved, as many codons harbor nearly equal amounts of resistant and susceptible mutations, which is consistent with the ability of a conservative, hydrophobic mutation to be tolerated in a region that relies on non-specific, volume-mediated packing driven by the hydrophobic effect.

One major question that remains is whether the mutations not called by the model (U) represent inaccuracies in the calculation of model features, breakdowns in the model, or mutations with an intermediate effect on protein stability and/or enzyme activity. Most of the un-called mutations have intermediate features that lie in between the resistant and susceptible distributions (**Figure S2**). The MAPP score has been shown to be capable of delineating between mutations that have mild and severely deleterious effects in other genes, suggesting that mutations with intermediate MAPP scores may indeed be intermediate in effect<sup>48</sup>. In addition, some of the mutations that are called U appear to not be reproducible when experimentally tested using the gold-standard culture-based method, MGIT, supporting the possibility of an intermediate class (**Figure 5B**). One previous study has shown associations between reductions in PncA stability/function *in vitro* and outcomes in infected mice, but more work is necessary to fully understand whether this relationship extends to clinical outcomes in patients<sup>39</sup>. While mutations have historically been classified using a binary system, this study

supports the view of mutations as conferring a spectrum of resistance. Future approaches could examine either probabilistic modelling or multi-class classification to attempt to encapsulate the uncertainty in phenotype associated with certain *pncA* mutations.

Predictions made by this model could provide clinicians with an initial estimate of pyrazinamide susceptibility after a novel mutation is observed but before traditional phenotypic testing has been completed. Given the latter can take weeks or even months, this could help guide initial therapy and further antibiotic susceptibility testing. In addition, the putative classification of additional *pncA* mutations potentially enables genetic variants conferring pyrazinamide resistance that do not involve the *pncA* gene to be discovered. The identification of pyrazinamide-susceptible mutations is also crucial, as it has been suggested that any nonsynonymous mutation in *pncA* that is not cataloged as susceptible confers resistance, an incorrect assumption that would lead to overprediction of pyrazinamide resistance<sup>53</sup>.

This study constitutes a proof-of-concept for the computational prediction of pyrazinamide resistance, a critically important drug in the treatment of tuberculosis. However, this approach is not limited to *pncA* but should in theory be extensible to any prodrug system where the converting enzyme is non-essential, such as delamanid, protaminid, or ethionamide as well as to prodrug systems outside of *M. tuberculosis*. Interestingly, a recent study has highlighted similar trends in the features used in this study for resistance-conferring mutations in *katG* (isoniazid), *rpoB* (rifampicin), and *alr* (D-cycloserine), suggesting that this approach may even be applicable to non-prodrug systems<sup>54</sup>. One promising area for future work is in the tuberculosis drug bedaquiline, where resistance is caused in part by mutations in a transcriptional repressor (*Rv0678*) that cause loss of DNA binding and upregulation of efflux pumps<sup>55,56</sup>. *Rv0678* has shown a high degree of mutational promiscuity in published sequencing studies, which would highlight the value of a computational approach<sup>57-61</sup>. The ability of this approach to identify the major mechanisms of resistance to pyrazinamide highlights the need for continued basic research to determine the structures of other bacterial proteins implicated in

antibiotic resistance. Additionally, the efficacy of this approach highlights the value of including evolutionary constraints for prediction of mutational effects. Further understanding of the effect of *pncA* mutations also increases the ability of whole-genome sequencing approaches to move to the forefront of global tuberculosis control efforts.

## Acknowledgements

The study was funded by the National Institute for Health Research Health Protection Research Unit (NIHR HPRU) in Healthcare Associated Infections and Antimicrobial Resistance at Oxford University in partnership with Public Health England (PHE) [HPRU-2012-10041]; the National Institute for Health Research (NIHR) Oxford Biomedical Research Centre (BRC); the CRyPTIC consortium, which is funded by a Wellcome Trust/Newton Fund-MRC Collaborative Award [200205/Z/15/Z] and the Bill and Melinda Gates Foundation Trust [OPP1133541]; and the South African Medical Research Council. The EXIT-RIF project (Prof Annelies Van; Prof Rob Warren, Prof Lesley Scott, Prof Wendy Stevens, Dr Michael Whitfield) is funded by National Institutes of Health grant [#R01 AI099026]. J.J.C. is supported by a fellowship from the Rhodes Trust. T.E.A.P. and D.W.C. are NIHR Senior Investigators. T.M.W. is an NIHR Academic Clinical Lecturer. J.J.C would like to thank Spencer Dunleavy for thoughtful discussions on statistical analysis and modeling. The content is the solely the responsibility of the authors and does not necessarily represent the official views of the South African Medical Research Council. MGW is supported by a fellowship from the Claude Leon Foundation. The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health.

## Author Contributions

JJC, PWF, and TMW designed experiments, JJC carried out the experiments, JJC and ASW performed statistical analyses, JJC, PWF, ASW, and TEAP wrote the manuscript. PWF and DWC supervised the work. MGW contributed samples and edited the manuscript.

## Declaration of Interests

The authors have no interests to declare.

## Methods

Mutations included in the derivation (training and testing) dataset were selected from four large studies/reviews that included phenotypic diagnostic susceptibility testing of clinical isolates for pyrazinamide and one *in vitro/in vivo* phenotypic screening study<sup>31,37-39,52</sup>. Briefly, phenotypes for strains with single missense mutations in *pncA* from the four studies of clinical isolates were aggregated by mutation, tallying the results of the phenotypic testing. Mutations that were resistant or susceptible at least 75% of the time and that had been phenotyped at least 4 times were included as derivation phenotypes. Additionally, mutations that had been phenotyped at least twice with no discrepancies were also included. These mutations were then cross-referenced and supplemented with mutations from Yadon *et al.* that were either enriched (resistant) or depleted (susceptible) in both the *in vitro* and *in vivo* screens performed in that study<sup>39</sup>. Mutations that had conflicting clinical and laboratory phenotypes (n=2) were removed from the derivation dataset. Mutations that were only present in either clinical isolates or *in vitro* isolates but that met the criteria for inclusion from that set were included. This led to a final total of 291 mutations with high-confidence phenotypes of which 184 were resistant and 107 were susceptible to pyrazinamide.

The change in mass, volume, charge, hydrophobicity, distance from the Fe<sup>2+</sup> atom and pyrazinamide molecule, solvent accessibility, MAPP score, Rogov score, degree of hydrogen

bonding, and predicted change in the free energy of protein unfolding were determined for each mutation. Hydrophobicity was estimated using the Kyte-Doolittle scale. Distances were calculated as the minimum distance between each residue and the Fe<sup>2+</sup> atom or pyrazinamide molecule using UCSF Chimera. Solvent accessibility and predicted number of hydrogen bonds were calculated in UCSF Chimera. *In silico* calculation of the change in free energy of protein unfolding was calculated using a meta-predictor as described in Broom *et al*<sup>46</sup>. The MAPP score was calculated using software available at (<http://mendel.stanford.edu/SidowLab/downloads/MAPP/index.html>) using related orthologs. The PncA amino acid sequence alignment used to generate the MAPP score can be found in the Supplementary Material.

Logistic regression, support vector machines with radial kernels, and multi-layer perceptron neural networks were implemented using the R caret package (Supplemental Code 1). Briefly, 70% of the derivation dataset was randomly selected (maintaining the ratio of resistant to susceptible phenotypes) as a training set and 30% was reserved as a test set. All three model types were trained using repeated (n=10) using 10-fold cross validation with class weights to compensate for the class imbalance. Performance was then estimated on the test set. In order to select the variables used for logistic regression, backwards stepwise elimination (exit p=0.15) was performed on the entire derivation dataset to select the main effects and then interactions between the significant terms were manually investigated, retaining any with heterogeneity p<0.01. Two additional weak interactions (between the protein destabilization factor and either number of hydrogen bonds (p=0.073) or solvent accessibility (p=0.054)) were not included in the final model. The final logistic regression model was trained using the identified main effects with the two significant interactions on the training set using 10-fold cross validation to select hyperparameters before being applied to the test set.

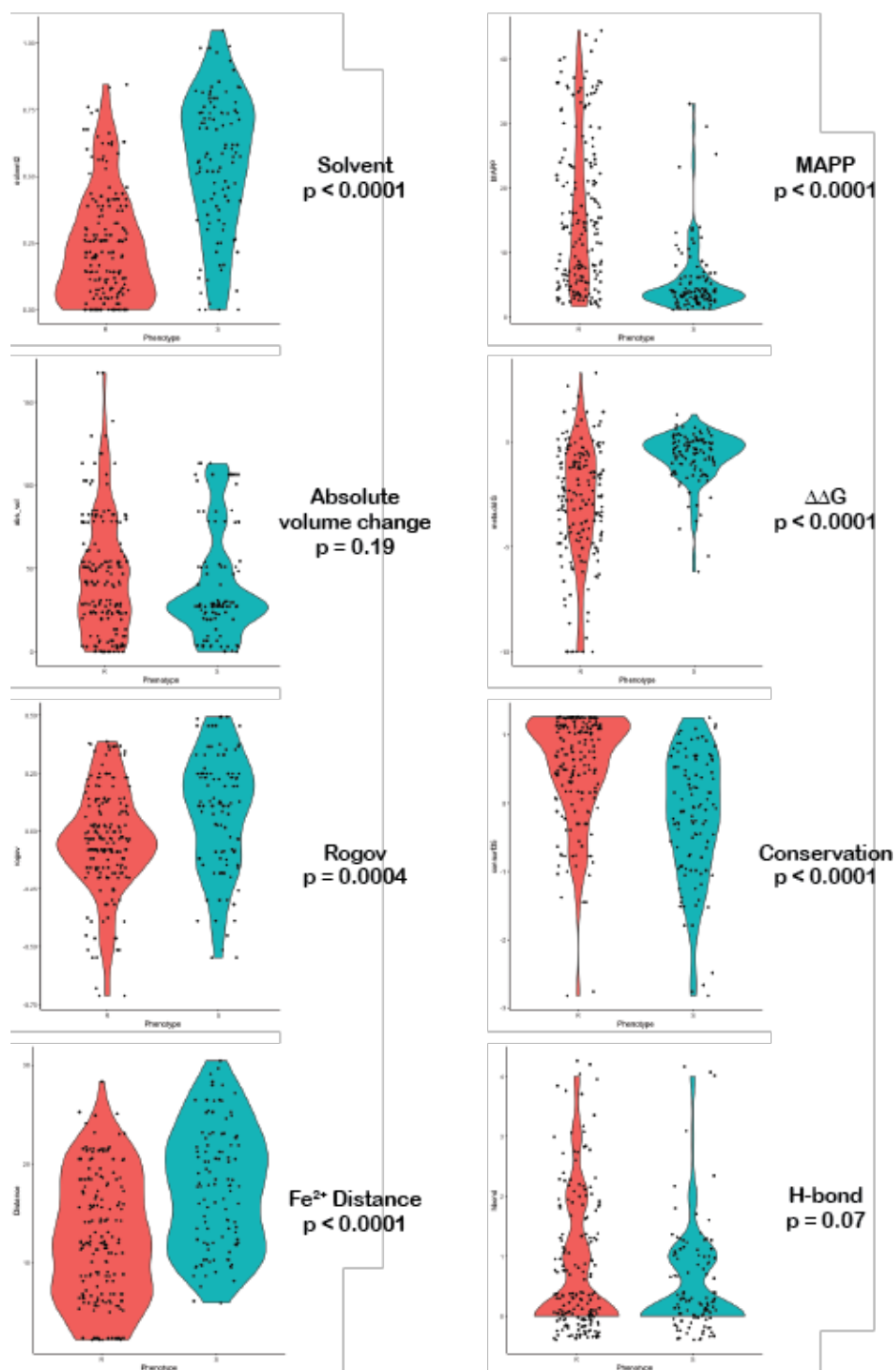
All isolates in the quantitative MIC test set were collected in South Africa. Of the 366 *Mycobacterium tuberculosis* clinical isolates, 333 were collected as part of a prospective cohort



study (“EXIT-RIF”) aimed at comparing the outcome of patients diagnosed with rifampicin resistant tuberculosis by MTBDR*plus* (Hain LifeSciences) or Xpert MTB/RIF between November 2012 and December 2013 in three South African provinces (Free State, Eastern Cape and Gauteng). A *Mycobacterium tuberculosis* databank housed at the SAMRC Centre for Tuberculosis Research, consisting of ~45,000 drug resistant isolates collected in the Western Cape province since 2001, was queried to identify isolates containing both PZA MIC data and *pncA* genotypic data, this produced the remaining 33 *Mycobacterium tuberculosis* clinical isolates. 27 isolates were selected from this set that harbored single amino acid mutations in PncA.

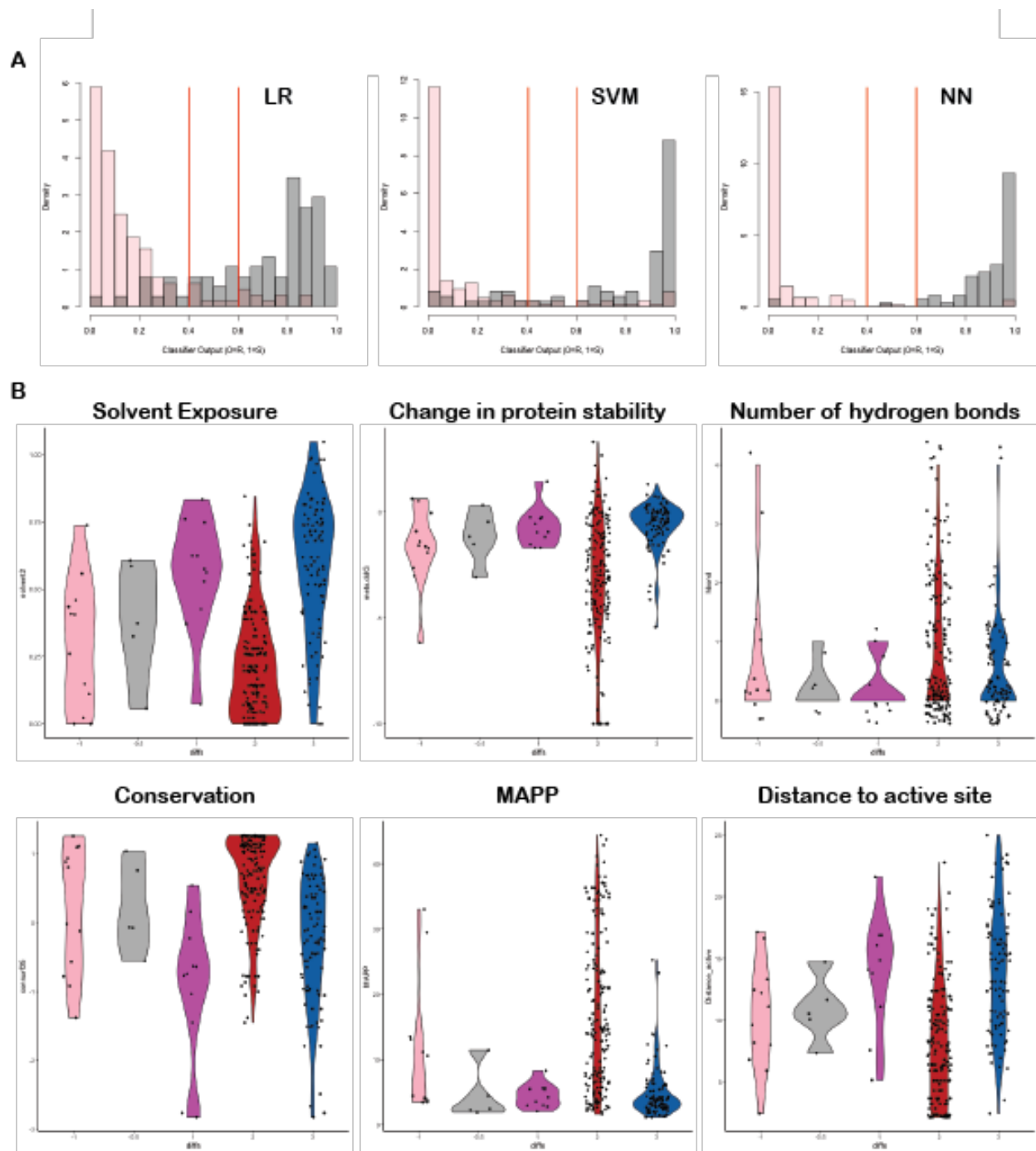
All MICs were determined using the non-radiometric BACTEC MGIT 960 method (BD Diagnostic Systems, NJ, USA) with manufactured supplied pyrazinamide medium/supplement as previously described<sup>62</sup>. This system makes use of modified test media which supports the growth of mycobacteria at a pH of 5.9. The MICs were determined at 900, 600 and 300 µg/ml for the large deletion isolates and 100, 75, 50, 25 µg/ml for the rest. A fully susceptible MTB laboratory strain H37Rv (ATCC 27294) was included as a control.

## Supplemental Figures



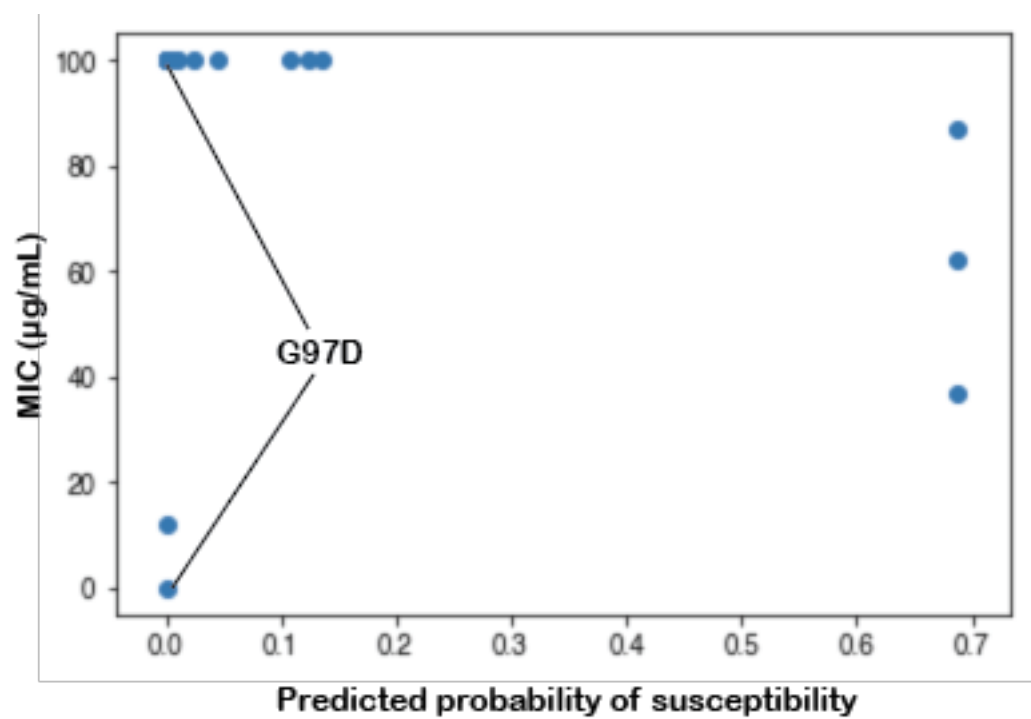
**Figure S1. Distributions of structural features vary with resistance.**

Resistant (red) and susceptible (blue) features are shown with a p-value calculated by a Mann-Whitney U test.



**Figure S2. Models effectively predict resistant and susceptible PncA mutations**

(A) Distributions of model predictions on the training sets. Samples where  $p < 0.4$  were defined as resistant, whilst any sample with  $p > 0.6$  was defined as susceptible. (B) Distributions of major errors (pink), minor errors (grey), very major errors (purple), resistant (red), and susceptible (blue) mutations for the neural network model.



**Figure S3. The predicted probabilities of susceptibility are negatively correlated with pyrazinamide MIC.**

## References

1. WHO. *Global Tuberculosis Report 2017*. World Health Organization (2017).  
doi:WHO/HTM/TB/2017.23
2. Who & The World Health Organization. *Treatment of tuberculosis: guidelines. 4Th Ed.* 160 (2010). doi:10.1164/rccm.201012-1949OC
3. Njire, M. *et al.* Pyrazinamide resistance in Mycobacterium tuberculosis: Review and update. *Adv. Med. Sci.* **61**, 63–71 (2016).
4. Zhang, Y. & Yew, W. W. Mechanisms of drug resistance in Mycobacterium tuberculosis: Update 2015. *Int. J. Tuberc. Lung Dis.* **19**, 1276–1289 (2015).
5. Zhang, Y. & Mitchison, D. The curious characteristics of pyrazinamide: A review. *Int. J. Tuberc. Lung Dis.* **7**, 6–21 (2003).
6. Mitchison, D. A. The action of antituberculosis drugs in short-course chemotherapy. *Tubercle* **66**, 219–225 (1985).
7. Fox, W., Ellard, G. A. & Mitchison, D. A. Studies on the treatment of tuberculosis undertaken by the British Medical Research Council Tuberculosis Units, 1946-1986, with relevant subsequent publications. *Int. J. Tuberc. Lung Dis.* **3**, (1999).
8. Dawson, R. *et al.* Efficiency and safety of the combination of moxifloxacin, pretomanid (PA-824), and pyrazinamide during the first 8 weeks of antituberculosis treatment: A phase 2b, open-label, partly randomised trial in patients with drug-susceptible or drug-resistant pul. *Lancet* **385**, 1738–1747 (2015).
9. Chang, K. C. *et al.* Pyrazinamide may improve fluoroquinolone-based treatment of multidrug-resistant tuberculosis. *Antimicrob. Agents Chemother.* **56**, 5465–5475 (2012).
10. Zumla, A. I. *et al.* New antituberculosis drugs, regimens, and adjunct therapies: Needs, advances, and future prospects. *Lancet Infect. Dis.* **14**, 327–340 (2014).
11. Nuermberger, E. *et al.* Powerful bactericidal and sterilizing activity of a regimen

- containing PA-824, moxifloxacin, and pyrazinamide in a murine model of tuberculosis. *Antimicrob. Agents Chemother.* **52**, 1522–1524 (2008).
12. Rosenthal, I. M. *et al.* Daily dosing of rifapentine cures tuberculosis in three months or less in the murine model. *PLoS Med.* **4**, 1931–1939 (2007).
  13. Veziris, N. *et al.* A once-weekly R207910-containing regimen exceeds activity of the standard daily regimen in murine tuberculosis. *Am. J. Respir. Crit. Care Med.* **179**, 75–79 (2009).
  14. Whitfield, M. G. *et al.* A global perspective on pyrazinamide resistance: Systematic review and meta-analysis. *PLoS One* **10**, 1–16 (2015).
  15. World Health Organization. *Guidelines for surveillance of drug resistance in tuberculosis.* (2010).
  16. Chang, K. C., Yew, W. W. & Zhang, Y. Pyrazinamide susceptibility testing in mycobacterium tuberculosis: A systematic review with meta-analyses. *Antimicrob. Agents Chemother.* **55**, 4499–4505 (2011).
  17. Jr., H. D., Horn, D. L. & Alfalla, C. Drug-Resistant Tuberculosis: Inconsistent Results of Pyrazinamide Susceptibility Testing. *JAMA J. Am. Med. Assoc.* **273**, 916–917 (1995).
  18. Miller, M. A., Thibert, L., Desjardins, F., Siddiqi, S. H. & Dascal, A. Testing of susceptibility of Mycobacterium tuberculosis to pyrazinamide: Comparison of Bactec method with pyrazinamidase assay. *J. Clin. Microbiol.* **33**, 2468–2470 (1995).
  19. Hoffner, S. *et al.* Proficiency of drug susceptibility testing of Mycobacterium tuberculosis against pyrazinamide: the Swedish experience. *Int J Tuberc Lung Dis* **17**, 1486–1490 (2013).
  20. Pandey, S., Newton, S., Upton, A., Roberts, S. & Drinkovi, D. Characterisation of pncA mutations in clinical Mycobacterium tuberculosis isolates in New Zealand. *Pathology* **41**, 582–584 (2009).
  21. Simons, S. O. *et al.* Validation of pncA Gene Sequencing in Combination with the

- Mycobacterial Growth Indicator Tube Method To Test Susceptibility of Mycobacterium tuberculosis to Pyrazinamide. 428–434 (2012). doi:10.1128/JCM.05435-11
22. Chedore, P., Bertucci, L., Wolfe, J., Sharma, M. & Jamieson, F. Potential for erroneous results indicating resistance when using the bactec MGIT 960 system for testing susceptibility of Mycobacterium tuberculosis to pyrazinamide. *J. Clin. Microbiol.* **48**, 300–301 (2010).
  23. World Health Organization. WHO Guideline: The use of molecular line probe assays for the detection of resistance to isoniazid and rifampicin. 21–23 (2016).
  24. Scorpio, A. *et al.* Characterization of *pncA* mutations in pyrazinamide-resistant *Mycobacterium tuberculosis*. *Antimicrob. Agents Chemother.* **41**, 540–543 (1997).
  25. Ramirez-Busby, S. M. *et al.* A Multinational Analysis of Mutations and Heterogeneity in PZase, RpsA, and PanD Associated with Pyrazinamide Resistance in M/XDR *Mycobacterium tuberculosis*. *Sci. Rep.* **7**, 1–9 (2017).
  26. Sheen, P. *et al.* A multiple genome analysis of *Mycobacterium tuberculosis* reveals specific novel genes and mutations associated with pyrazinamide resistance. (2017). doi:10.1186/s12864-017-4146-z
  27. Gopal, P. *et al.* Pyrazinamide resistance is caused by two distinct mechanisms: Prevention of coenzyme a depletion and loss of virulence factor synthesis. *ACS Infect. Dis.* **2**, 616–626 (2016).
  28. Zhang, Y., Zhang, J., Cui, P., Zhang, Y. & Zhang, W. crossm Identification of Novel Efflux Proteins. **61**, 1–10 (2017).
  29. Hirano, K., Takahashi, M., Kazumi, Y., Fukasawa, Y. & Abe, C. Mutation in *pncA* is a major mechanism of pyrazinamide resistance in *Mycobacterium tuberculosis*. *Tuber. Lung Dis.* **78**, 117–122 (1997).
  30. Stoffels, K., Mathys, V., Fauville-Dufaux, M., Wintjens, R. & Bifania, P. Systematic analysis of pyrazinamide-resistant spontaneous mutants and clinical isolates of

- Mycobacterium tuberculosis. *Antimicrob. Agents Chemother.* **56**, 5186–5193 (2012).
31. Degano, M. *et al.* Mycobacterium tuberculosis Pyrazinamide Resistance Determinants :  
a. *Mol. Microbiol.* **5**, 1–10 (2014).
  32. Scorpio, A. & Zhang, Y. Mutations in *pncA*, a gene encoding  
pyrazinamidase/nicotinamidase, cause resistance to the antituberculous drug  
pyrazinamide in tubercle bacillus. *Nat. Med.* **2**, 662–667 (1996).
  33. Kim, N., Petingi, L. & Schlick, T. Network theory tools for RNA modeling. *WSEAS Trans.  
Math.* **12**, 941–955 (2013).
  34. Driesen, M. *et al.* Evaluation of a novel line probe assay to detect resistance to  
pyrazinamide, a key drug used for tuberculosis treatment. *Clin. Microbiol. Infect.* **24**, 60–  
64 (2018).
  35. Ramirez-Busby, S. M. & Valafar, F. Systematic review of mutations in pyrazinamidase  
associated with pyrazinamide resistance in mycobacterium tuberculosis clinical isolates.  
*Antimicrob. Agents Chemother.* **59**, 5267–5277 (2015).
  36. Kalokhe, A. S. *et al.* Multidrug-resistant tuberculosis drug susceptibility and molecular  
diagnostic testing: a review of the literature. *Am J Med Sci.* **345**, 143–148 (2013).
  37. The CRyPTIC Consortium and the 100, 000 Genomes Project. Prediction of Susceptibility  
to First-Line Tuberculosis Drugs by DNA Sequencing.  
<https://doi.org/10.1056/NEJMoa1800474> 1–14 (2018). doi:10.1056/NEJMoa1800474
  38. Whitfield, M. G. *et al.* Mycobacterium tuberculosis *pncA* polymorphisms that do not confer  
pyrazinamide resistance at a breakpoint concentration of 100 micrograms per milliliter in  
MGIT. *J. Clin. Microbiol.* **53**, 3633–3635 (2015).
  39. Yadon, A. N. *et al.* A comprehensive characterization of PncA polymorphisms that confer  
resistance to. *Nat. Commun.* doi:10.1038/s41467-017-00721-2
  40. Petrella, S. *et al.* Crystal structure of the pyrazinamidase of mycobacterium tuberculosis:  
Insights into natural and acquired resistance to pyrazinamide. *PLoS One* **6**, (2011).



41. Faure, G. & Koonin, E. V. Universal distribution of mutational effects on protein stability, uncoupling of protein robustness from sequence evolution and distinct evolutionary modes of prokaryotic and eukaryotic proteins. *Phys. Biol.* **12**, (2015).
42. Lim, W. A., Farruggio, D. C. & Sauer, R. T. Structural and Energetic Consequences of Disruptive Mutations in a Protein Core. *Biochemistry* **31**, 4324–4333 (1992).
43. Guo, H. H., Choe, J. & Loeb, L. A. Protein tolerance to random amino acid change. *Proc. Natl. Acad. Sci.* **101**, 9205–9210 (2004).
44. Chen, H. & Zhou, H. X. Prediction of solvent accessibility and sites of deleterious mutations from protein sequence. *Nucleic Acids Res.* **33**, 3193–3199 (2005).
45. Yoon, J. H., Nam, J. S., Kim, K. J. & Ro, Y. T. Characterization of pncA mutations in pyrazinamide-resistant Mycobacterium tuberculosis isolates from Korea and analysis of the correlation between the mutations and pyrazinamidase activity. *World J. Microbiol. Biotechnol.* **30**, 2821–2828 (2014).
46. Broom, A., Jacobi, Z., Trainor, K. & Meiering, E. M. Computational tools help improve protein stability but with a solubility tradeoff. *J. Biol. Chem.* **292**, 14349–14361 (2017).
47. Fowler, P. W. *et al.* Robust Prediction of Resistance to Trimethoprim in Staphylococcus aureus Resource Robust Prediction of Resistance to Trimethoprim in Staphylococcus aureus. *Cell Chem. Biol.* 1–11 (2018). doi:10.1016/j.chembiol.2017.12.009
48. Stone, E. A. & Sidow, A. Physicochemical constraint violation by missense substitutions mediates impairment of protein function and disease severity. 978–986 (2005). doi:10.1101/gr.3804205.
49. U.S. Department of Health and Human Services Food and Drug Administration. Guidance for Industry and FDA. Class II Special Controls Guidance Document : Antimicrobial Susceptibility Test Systems. 1–42 (2009).
50. Pal, D. & Chakrabarti, P. Cis peptide bonds in proteins: residues involved, their conformations, interactions and locations. *J. Mol. Biol.* **294**, 271–288 (1999).

51. Bradley, P., Bakker, H. den, Rocha, E., McVean, G. & Iqbal, Z. Real-time search of all bacterial and viral genomic data. *bioRxiv* 234955 (2017). doi:10.1101/234955
52. Walker, T. M. *et al.* Whole-genome sequencing for prediction of Mycobacterium tuberculosis drug susceptibility and resistance: A retrospective cohort study. *Lancet Infect. Dis.* **15**, 1193–1202 (2015).
53. Zignol, M. *et al.* Population-based resistance of Mycobacterium tuberculosis isolates to pyrazinamide and fluoroquinolones: results from a multicountry surveillance project. *Lancet Infect. Dis.* **16**, 1185–1192 (2016).
54. Portelli, S., Phelan, J. E., Ascher, D. B., Clark, T. G. & Furnham, N. Understanding molecular consequences of putative drug resistant mutations in Mycobacterium tuberculosis. *Sci. Rep.* 1–12 (2018). doi:10.1038/s41598-018-33370-6
55. Milano, A. *et al.* Azole resistance in Mycobacterium tuberculosis is mediated by the MmpS5-MmpL5 efflux system. *Tuberculosis* **89**, 84–90 (2009).
56. Nguyen, T. V. A., Anthony, R. M., Bañuls, A. L., Vu, D. H. & Alffenaar, J. W. C. Bedaquiline Resistance: Its Emergence, Mechanism, and Prevention. *Clin. Infect. Dis.* **66**, 1625–1630 (2018).
57. Zhang, S. *et al.* Identification of novel mutations associated with clofazimine resistance in Mycobacterium tuberculosis. *J. Antimicrob. Chemother.* **70**, 2507–2510 (2015).
58. Villellas, C. *et al.* Unexpected high prevalence of resistance-associated Rv0678 variants in MDR-TB patients without documented prior use of clofazimine or bedaquiline. *J. Antimicrob. Chemother.* **72**, 684–690 (2017).
59. Xu, J. *et al.* Primary clofazimine and bedaquiline resistance among isolates from patients with multidrug-resistant tuberculosis. *Antimicrob. Agents Chemother.* **61**, 1–8 (2017).
60. Hartkoorn, R. C., Uplekar, S. & Cole, S. T. Cross-resistance between clofazimine and bedaquiline through upregulation of mmp15 in mycobacterium tuberculosis. *Antimicrob. Agents Chemother.* **58**, 2979–2981 (2014).

61. Somoskovi, A., Bruderer, V., Hömke, R., Bloemberg, G. V. & Böttger, E. C. A mutation associated with clofazimine and bedaquiline cross-resistance in MDR-TB following bedaquiline treatment. *Eur. Respir. J.* **45**, 554–557 (2015).
62. Piersimoni, C. *et al.* Prevention of false resistance results obtained in testing the susceptibility of *Mycobacterium tuberculosis* to pyrazinamide with the bactec MGIT 960 system using a reduced inoculum. *J. Clin. Microbiol.* **51**, 291–294 (2013).