



26 Charlotte Genestet, PhD

27 Centre International de Recherche en Infectiologie, 7 rue Guillaume Paradin, 69003 Lyon,

28 France. Tel: +33 (0)4 72 07 17 09

29

30

31 **ABSTRACT**

32 Tuberculosis (TB), caused by *Mycobacterium tuberculosis* (Mtb) complex, is still the number  
33 one deadly contagious disease. Mtb infection results in a wide spectrum of clinical presentations  
34 and severity symptoms, but without proven Mtb genetic determinants. Thanks to a collection  
35 of 355 clinical isolates with associated patient's clinical data, we showed that Mtb micro-  
36 diversity within patient isolates is strongly correlated with TB-associated severity scores.  
37 Interestingly, this diversity is driven by a selection pressure to adapt to different lifestyles  
38 related to the infection site. Taken together, these results provide a new insight to better  
39 understand TB pathophysiology. Furthermore, Mtb micro-diversity could be envisioned as a  
40 new prognostic tool to improve the management of TB patients.

41

42

43 **Keywords:** *Mycobacterium tuberculosis*; micro-diversity; tuberculosis severity; pulmonary  
44 tuberculosis; extra-pulmonary tuberculosis; whole genome sequencing; selection pressure;  
45 compartmentalization

46

## 47 INTRODUCTION

48 Tuberculosis (TB) caused by *Mycobacterium tuberculosis* (Mtb) complex remains one of the  
49 most prevalent and deadly infectious diseases, responsible for 10 million new cases and 1.2  
50 million deaths among HIV-negative people worldwide in 2018, and an additional 208 000  
51 deaths among HIV-positive people <sup>1</sup>. Mtb infections result in a wide spectrum of clinical  
52 outcomes, from latent asymptomatic infection to pulmonary or extra- pulmonary manifestations  
53 of disease, with an array of severity symptoms. Such diversity has been historically attributed  
54 to host and environmental factors, while the Mtb complex was previously considered  
55 genetically monomorphic <sup>2</sup>.

56 Since the introduction of the next generation sequencing (NGS) enabling whole genome  
57 sequencing (WGS), outstanding progress has been done in the field of Mtb genomics.  
58 Increasing studies based on NGS have revealed micro-diversity in Mtb clinical isolates: within  
59 hosts, minor variants coexist rather than a clonal colony <sup>3-14</sup>. Although we showed in a previous  
60 study the involvement of initial Mtb micro-diversity in intra-macrophagic persistence and  
61 antibiotic tolerance <sup>15</sup>, the role and the impact of Mtb micro-diversity is still poorly understood.  
62 It was already showed that mixed-strains Mtb infection is associated with poor outcome <sup>16,17</sup>,  
63 but only few reports investigated the link between Mtb micro-diversity and TB severity and  
64 outcome <sup>9,10</sup>. Yet, several major questions remain unanswered: such as the Mtb ability to cause  
65 active disease and various symptoms. While many Mtb virulence factors are well described, to  
66 date, there are no proven genetic determinants associated with virulence, disease progression  
67 or severity of TB <sup>18</sup>.

68 Regarding other bacterial species responsible of chronic infections, such as *Helicobacter pylori*,  
69 *Pseudomonas aeruginosa*, *Staphylococcus aureus*, *Streptococcus pneumoniae*, several reports  
70 highlighted the role of micro-diversity in pathogen adaptation between and within-host.  
71 Bacterial micro-diversity plays a role in the adaptation to immune and treatment pressure, for

72 infecting different body sites and has been suggested to impact outcome and severity of illness  
73 <sup>19–25</sup>.

74 Accordingly, here we explored the relation between Mtb micro-diversity of 355 clinical  
75 isolates, from 311 patients diagnosed at the Lyon University Hospital, and TB clinical  
76 presentation, as well as nutritional and immune status of patients. It revealed a strong correlation  
77 between the detection of Mtb micro-diversity within clinical isolates and TB-associated severity  
78 markers. Furthermore, thanks to a cohort of 42 patients with both microbiologically proven  
79 pulmonary and extra-pulmonary TB, we investigated intra-host micro-evolution of Mtb. We  
80 observed a compartmentalization of variants, driven by a selection pressure to adapt to different  
81 tissues, as shown by dN/dS approach. It should be noted that, besides the detection of canonical  
82 drug resistance determinants, the detection of Mtb micro-diversity within patient's isolates is,  
83 to date, the only other bacterial feature which could be envisioned as a prognostic marker of  
84 poor TB outcome.

85

## 86 **RESULTS**

### 87 *Baseline characteristics of study populations*

88 A total of 355 clinical Mtb isolates from 311 patients were included in this study, including 42  
89 patients with both microbiologically proven pulmonary and extra-pulmonary TB (PTB and  
90 EPTB) and 2 patients with 2 extra-pulmonary localization without PTB. Moreover, 25 patients  
91 with PTB also developed EPTB according to the clinical manifestations and 2 patients with  
92 EPTB may also have PTB but no pulmonary sampling was conducted. The **Table 1** summarize  
93 patients' demographic and descriptive characteristics. As it is a retrospective study, all data  
94 were not available for all patients, then the population size was always specified.  
95 Patient's descriptive characteristics revealed different profiles between patients with PTB or  
96 EPTB. Overall, higher severity scores and lower malnutrition indices were observed for patient

97 with PTB compared to patients with EPTB. Conversely, a lower immune status was observed  
 98 for patient with EPTB compared to patients with PTB. It should be noted that serum lipid and  
 99 serum iron profiles and lymphocyte typing can not be properly compared between the cohorts  
 100 as the data were available only for 20% to 43% of patients.

101

102 **Table 1: Patients' demographic and baseline descriptive characteristics**

Characteristics	Pulmonary TB n=244	Extra-pulmonary TB n=111	<i>p-values</i>
<b>Demographic data</b>			
Age (years), median [IQR]	36 [25-58]	37 [24-60]	0.7505
Gender (male), n (%)	162 (66.4)	57 (51.4)	<b>0.0094</b>
<b>TB outcome</b>			
Fatal outcome, n (%)	8 (3.3)	3 (2.7)	0.6999
Unknown, n (%)	18 (7.4)	11 (9.9)	0.4109
<b>Comorbidity</b>			
HIV, n (%)	15/233 (6.4)	8/107 (7.5)	0.8165
Hepatitis, n (%)	24/189 (12.7)	13/85 (15.3)	0.5701
Diabetes, n (%)	26/234 (11.1)	14/105 (13.3)	0.5867
Immunosuppressive therapy, n (%)	21/235 (8.9)	16/107 (15.0)	0.1317
History of TB, n (%)	19/233 (8.2)	4/106 (3.8)	0.1662
<b>TB-associated severity scores</b>			
Bandim TB score, median [IQR] (n)	4 [3-6] (233)	3 [1.5-5] (105)	<b>0.0002</b>
MUST, median [IQR] (n)	3 [1-5] (227)	2 [0-4] (105)	<b>0.0293</b>
Nutritional risk score, median [IQR] (n)	1 [0.25-3] (56)	1 [0-3] (24)	0.8839
<b>Nutritional status</b>			
BMI (kg/m <sup>2</sup> ), median [IQR] (n)	20.2 [17.7-22.4] (228)	21.63 [19.0-26.2] (107)	<b>&lt;0.0001</b>
Weightloss, median [IQR] (n)	8.0 [4.5-12.3] (227)	6.4 [2.4-11.5] (106)	<b>0.0293</b>
Serum albumin (g/L), median [IQR] (n)	30.0 [26.0-36.7] (209)	31.1 [25.9-36.2] (92)	0.8693
Serum pre-albumin (g/L), median [IQR] (n)	0.13 [0.09-0.18] (152)	0.11 [0.08-0.17] (60)	0.3448
<b>Serum/blood parameters</b>			
Total protein (g/L), median [IQR] (n)	75 [70-81] (235)	75 [70-81] (102)	0.9509
Blood glucose (mmol/L), median [IQR] (n)	5.0 [4.5-6.0] (215)	5.2 [4.6-6.4] (100)	0.3719
Calcium (mmol/L), median [IQR] (n)	2.44 [2.31-2.61] (236)	2.43 [2.32-2.59] (92)	0.7641
Phosphorus (mmol/L), median [IQR] (n)	0.98 [0.85-1.18] (118)	1.07 [0.84-1.21] (57)	0.6397
Magnesium (mmol/L), median [IQR] (n)	0.81 [0.74-0.88] (70)	0.80 [0.71-0.87] (42)	0.2684
Sodium (mmol/L), median [IQR] (n)	137 [134-139] (236)	137 [134-139] (62)	0.1480
Potassium (mmol/L), median [IQR] (n)	4.0 [3.8-4.3] (236)	4.0 [3.7-4.3] (104)	0.3781
Chloride (mmol/L), median [IQR] (n)	103 [100-106] (235)	103 [98-106] (103)	0.2117
Bicarbonate (mmol/L), median [IQR] (n)	25 [23-27] (178)	25 [23-27] (80)	0.8603
Anion gap, median [IQR] (n)	13 [11-14] (165)	13 [11-15] (78)	0.3953
<b>Serum lipid profile</b>			
Total cholesterol (mmol/L), median [IQR] (n)	4.1 [3.3-5.1] (56)	4.2 [3.1-4.9] (24)	1.0000
HDL-C (mmol/L), median [IQR] (n)	1.03 [0.69-1.29] (52)	1.00 [0.54-1.19] (23)	0.7302
LDL-C (mmol/L), median [IQR] (n)	2.58 [1.88-3.33] (50)	2.58 [1.99-3.25] (23)	0.8960
Triglyceride (mmol/L), median [IQR] (n)	1.14 [0.93-1.67] (59)	1.24 [0.99-1.69] (28)	0.4898
<b>Iron balance and inflammatory markers</b>			

Serum iron ( $\mu\text{mol/L}$ ), median [IQR] (n)	5.3 [3.0-9.2] (73)	6.7 [4.5-10.6] (39)	0.1555
Transferrin saturation (%), median [IQR] (n)	13.1 [8.2-20.1] (73)	13.4 [10.1-20.3] (39)	0.2921
Ferritin ( $\mu\text{g/L}$ ), median [IQR] (n)	195 [79-319] (98)	281 [104-759] (48)	<b>0.0336</b>
CRP (mg/L), median [IQR] (n)	48 [13-100] (228)	45 [14-93] (99)	0.8781
<b>Hemogram</b>			
Hemoglobin (g/L), median [IQR] (n)	120 [104-136] (236)	115 [102-130] (106)	0.1644
WBC (G/L), median [IQR] (n)	6.9 [5.3-9.0] (237)	6.3 [4.6-8.1] (105)	<b>0.0294</b>
Neutrophils (G/L), median [IQR] (n)	4.6 [3.2-6.2] (237)	4.4 [2.8-5.5] (105)	0.1218
Eosinophils (G/L), median [IQR] (n)	0.09 [0.03-0.18] (236)	0.075 [0.01-0.148] (104)	<b>0.0446</b>
Basophils (G/L), median [IQR] (n)	0.03 [0.02-0.05] (235)	0.02 [0.01-0.04] (104)	0.0591
Monocytes (G/L), median [IQR] (n)	0.62 [0.46-0.84] (237)	0.55 [0.38-0.75] (105)	<b>0.0179</b>
Lymphocytes (G/L), median [IQR] (n)	1.41 [0.89-2.00] (237)	1.18 [0.75-1.68] (105)	<b>0.0294</b>
<b>Lymphocyte typing</b>			
CD4 T cells (cell/ $\text{mm}^3$ ), median [IQR] (n)	338 [180-613] (81)	362 [198-555] (45)	0.8406
CD8 T cells (cell/ $\text{mm}^3$ ), median [IQR] (n)	274 [168-468] (79)	240 [151-384] (44)	0.1581
CD4/CD8 ratio, median [IQR] (n)	1.43 [0.82-2.07] (79)	1.77 [0.84-2.5] (44)	0.2249

103 Data were expressed as count (percentage, %) for dichotomous variables and as medians  
 104 (interquartile range [IQR]) for continuous values. The number of missing values was excluded  
 105 from the denominator. Fisher exact,  $\chi^2$  test or non-parametric Mann-Whitney U test was used  
 106 to compare groups where appropriate. *p-value* < 0.05 was considered significant. PTB:  
 107 pulmonary tuberculosis; EPTB: extrapulmonary tuberculosis; unknown outcome: loss of follow-  
 108 up or follow-up in another care facility; HIV: human immunodeficiency virus; MUST:  
 109 malnutrition universal screening tool; BMI: body mass index; HDL-C: high-density lipoprotein  
 110 cholesterol; LDL-C: low-density lipoprotein cholesterol; CRP: C-reactive protein; WBC: white  
 111 blood cell.

112

113 ***Mtb micro-diversity is not correlated with Mtb samples characteristics.***

114 WGS data of *Mtb* clinical isolates included in this study were analyzed to detect *Mtb* genetic  
 115 micro-diversity through unfixed mutations (at frequencies between 10 and 90%). These data  
 116 were used to identify the minimum number of variants in each *Mtb* clinical isolates. A variant  
 117 was defined as an assembly of mutations with similar frequencies ( $\pm 10\%$ ). This allowed to  
 118 calculate the  $\alpha$ -diversity of each *Mtb* isolates, which include both the number and the frequency  
 119 of the variants and the genetic distance between these variants.

120 First of all, correlation between Mtb micro-diversity and Mtb samples characteristics, such as  
121 smear results and time to positivity of Mtb culture samples, reflecting bacterial load in clinical  
122 isolates, type of extra-pulmonary or pulmonary samples, Mtb lineages and resistance status,  
123 were explored (**Fig. S1**). Risks of bias, such as delta between Mtb sampling or between  
124 treatment initiation and sampling (only concerning the patients with both microbiologically  
125 proven PTB and EPTB) and samples that underwent freeze / thaw cycle (biobank samples)  
126 compared to those analyzed after a single round of culture (routine practice), were also  
127 evaluated (**Fig. S2**). No significant correlation was observed, except a higher proportion of  
128 isolates with captured micro-diversity in biobank pulmonary samples compared to pulmonary  
129 samples analyzed in routine practice. It may be due to the selection of significant Mtb isolates  
130 from our biobank, such as MDR strains, highly transmissible Mtb strains and strains from  
131 patients with both PTB and EPTB, all factors usually associated with high TB severity.  
132 Interestingly, no significant correlation was established between the bacterial load and Mtb  
133 micro-diversity in clinical isolates.

134

135 ***Mtb micro-diversity is correlated with TB-associated severity indices.***

136 We aimed to investigate the association between Mtb  $\alpha$ -diversity and TB outcome. However,  
137 due to the very low rate of fatal outcome in our setting (8/311; 2.6%), no significant correlation  
138 could be established (**Fig. S3**). Therefore, the correlation between Mtb micro-diversity and 3  
139 TB severity indices was explored: i) the modified Bandim score (the most largely used PTB  
140 prognosis score), which consider 5 symptoms (cough, hemoptysis, dyspnea, chest pain, night  
141 sweats) and 5 clinical findings (anemia, tachycardia, positive finding at lung auscultation, fever,  
142 BMI)<sup>26,27</sup>; ii) the Malnutrition Universal Screening Tool (MUST), based on weight loss, BMI  
143 and anorexia to evaluate malnutrition status of TB patients<sup>28</sup>; iii) the nutritional risk score

144 considering BMI, hypoalbuminemia, hypocholesterolemia and severe lymphocytopenia, then  
145 including both nutritional and immune factors <sup>29</sup>.

146 A significant correlation was observed between the detection of micro-diversity ( $\alpha$ -diversity >  
147 1) in pulmonary Mtb isolates and high Bandim score (**Fig. 1A**,  $p<0.0001$ ), MUST (**Fig. 1B**,  
148  $p<0.0001$ ) and nutritional risk score (**Fig. 1C**,  $p<0.0001$ ). However, no correlation was  
149 observed between the ranges of  $\alpha$ -diversity in Mtb isolates and the TB-associated severity  
150 indices evaluated.

151 As expected, no correlation was observed between the detection of micro-diversity in extra-  
152 pulmonary Mtb isolates and the Bandim score (**Fig. 1D**), as this score is mainly based on PTB  
153 symptoms, nor with the MUST (**Fig. 1E**), as malnutrition is more characteristic of PTB (**Table**  
154 **1**) <sup>30</sup>. Conversely, a correlation between the detection of micro-diversity in extra-pulmonary  
155 Mtb isolates and high nutritional risk score (**Fig. 1F**,  $p=0.0131$ ) was observed, this severity  
156 index also including variables of the patient's immune status. As before, no correlation was  
157 observed between the ranges of  $\alpha$ -diversity in Mtb isolates and the TB severity indices  
158 evaluated.

159

160 *Independent variables associated with micro-diversity in Mtb pulmonary and extra-*  
161 *pulmonary isolates.*

162 The analysis of nutritional indices of TB patients revealed that the detection of micro-diversity  
163 in pulmonary Mtb isolates was correlated with low BMI, severe unintentional weight loss, low  
164 serum pre-albumin, hyponatremia, hypochloremia, low serum bicarbonate (**Fig. S4**), and also  
165 with low total cholesterol, low HDL and high triglyceride (**Fig. S5**), which are all associated  
166 with unfavorable TB outcome <sup>29-34</sup>. The detection of micro-diversity in extra-pulmonary Mtb  
167 isolates was correlated with severe unintentional weight loss, low serum pre-albumin and low  
168 total serum protein (**Figs. S5 and S6**). Furthermore, a correlation was observed between the



169 detection of micro-diversity in both Mtb pulmonary and extra-pulmonary isolates and iron  
170 deficiency, meaning low serum iron and low transferrin saturation (**Fig. 2**).

171 We also explored the correlation between micro-diversity in Mtb isolates and inflammation  
172 (serum CRP and ferritin) and immune markers (hemoglobin, blood white blood cell (WBC),  
173 neutrophil, eosinophil, basophil, monocyte, lymphocyte, CD4 T cells and CD8 T cells count in  
174 peripheral blood and CD4/CD8 ratio). No correlation was observed between these markers and  
175 the detection of micro-diversity in pulmonary Mtb isolates (**Fig. S7**). Conversely, we found a  
176 correlation between the detection of micro-diversity in extra-pulmonary Mtb isolates and high  
177 serum ferritin and low CD4 T cells count in peripheral blood (**Fig. S8**), which is a well-known  
178 risk factor for TB patients, especially for extra-pulmonary TB <sup>35,36</sup>. Finally, in all cases  
179 (nutritional and immune variables), except for phosphorus for pulmonary Mtb isolates and  
180 eosinophils count for extra-pulmonary Mtb isolates, no correlation was observed between the  
181 ranges of Mtb  $\alpha$ -diversity and the host variables studied.

182

### 183 *Models describing the relationship between Mtb micro-diversity and clinical parameters*

184 To go forward, we built a model to relate Mtb micro-diversity in pulmonary (P) and extra  
185 pulmonary (EP) locations to various explanatory variables obtained for most patients in the  
186 cohort. The response variable we chose to model is a binary (0/1) measure of whether Mtb  
187 micro-diversity was detected in clinical isolates.

188 In the case of pulmonary TB infections, model comparisons indicated a best model with only 9  
189 variables (double location of infection, Mtb lineage, weight loss, protein dosage, CRP,  
190 hemoglobin, leukocyte count, neutrophil count and the BANDIM score). 39 models had an  
191 Akaike weight less than two units from the best model. Only eight variables had likely effects  
192 (BANDIM score, Weight loss, Double Location, Hemoglobin, Neutrophil count, CRP, Protein,  
193 and Leukocyte count), and only the effect of Mtb lineage had importance between 0.5 and 0.73,

194 making it a plausible (but not likely) effect (**Fig. 3**). Multi-model inference of parameter values  
 195 ascertained that Mtb micro-diversity in pulmonary isolates was higher with increasing weight  
 196 loss, BANDIM, hemoglobin, and leukocyte counts, and lower with increasing neutrophil  
 197 counts, CRP, and protein (**Table 2**). All lineages had slightly different coefficients, but they  
 198 never differed from the baseline sufficiently for 95% confidence intervals to have the same sign  
 199 (**Table 2**) and asymptotic pairwise tests between lineage coefficients were all non-significant  
 200 in the best model (results not shown).

	Estimate	lower CI	higher CI	Uncond. variance	Nb models	Importance
(Intercept)	2.36E-01	-1.46E+01	1.50E+01	6.71E+01	2000	1.00
WeightLoss	2.73E-01	1.58E-01	3.89E-01	3.49E-03	2000	1.00
BandimScore	5.85E-01	3.15E-01	8.55E-01	1.88E-02	2000	1.00
DoubleLoc1	1.26E+00	2.03E-01	2.31E+00	2.93E-01	1962	0.99
Haemoglobin	2.84E-02	4.43E-03	5.24E-02	1.52E-04	1957	0.99
Neutrophiles	-6.36E-04	-1.52E-03	2.51E-04	2.17E-07	1827	0.93
CRP	-8.97E-03	-1.89E-02	1.00E-03	2.70E-05	1772	0.92
Protein	-4.81E-02	-1.08E-01	1.20E-02	1.00E-03	1654	0.87
Leucocytes	4.06E-04	-3.72E-04	1.18E-03	1.67E-07	1384	0.74
Lineage2	-6.54E-01	-2.86E+00	1.55E+00	1.37E+00	1224	0.65
Lineage3	-2.11E+00	-6.16E+00	1.95E+00	4.26E+00	1224	0.65
Lineage4	-6.93E-01	-2.82E+00	1.44E+00	1.26E+00	1224	0.65
Lineage6	-1.60E-01	-2.52E+00	2.20E+00	2.01E+00	1224	0.65
LineageM. bovis	-3.13E+00	-8.59E+00	2.32E+00	7.82E+00	1224	0.65

201 **Table 2:** Models of Mtb micro-diversity in pulmonary isolates. Multi-model estimates of  
 202 variable coefficients based on the 2000 models pooled in the consensus set of models; only  
 203 models with importance greater than 0.5 were kept in this table. “DoubleLoc1”: presence of  
 204 both P and EP locations of TB infection in the patient, based on microbiological or  
 205 pathophysiological evidence; “Lineage2”: effect of lineage 2 when compared to baseline  
 206 (lineage 1). “Nb models”: number of models in the consensus set (already the result of merging  
 207 four parallel model comparison processes) which incorporated the variable.

208

209 *Mtb micro-diversity within and between pulmonary and extra-pulmonary compartments.*

210 To better understand the role and the impact of Mtb micro-diversity, we focused on paired  
211 isolates from the training cohort to explore micro-evolution of Mtb within individuals (**Fig. 4**).  
212 First of all, the frequency in both pulmonary and extra-pulmonary compartments of each variant  
213 identified was explored. Among the 104 variants identified by WGS, 31/104 (30%) were  
214 specific of pulmonary isolates, 30/104 (29%) were specific of extra-pulmonary isolates, the  
215 frequencies of 22/104 (21%) variants were significantly increased or decreased ( $\geq 10\%$ )  
216 between the compartments and only 21/104 (20%) were found at similar frequencies between  
217 the compartments (**Fig. 4A**). These results suggest a compartmentalization of Mtb variants.  
218 Then we analyzed of the repartition on Mtb genome of the 168 pairwise mutations distances  
219 observed between paired pulmonary and extra-pulmonary isolates, among these 104 variants.  
220 However, it did not reveal any hot spot on Mtb genome (**Fig. 4B**). Among these 168 pairwise  
221 mutations, 22/168 (13%) were intergenic mutations, and 146/168 (87%) were located in coding  
222 regions, among whom 104/146 (71%) were nonsynonymous mutations and 42/146 (29%) were  
223 synonymous SNP.

224 These 104 nonsynonymous mutations were used to perform a principal component analysis  
225 (PCA) for mixed data (**Fig. 4C**). This analysis took into account the delta of frequency of each  
226 pairwise mutation between pulmonary and extra-pulmonary compartment and the functional  
227 categories of these mutations. The first principal component (15.65% of the variance) allowed  
228 a discrimination of pulmonary and extra-pulmonary mutations and of some functional  
229 categories. This revealed an association of pulmonary mutations with cell wall and cell  
230 processes and lipid metabolism categories. Otherwise, extra-pulmonary mutations were  
231 strongly associated with intermediary metabolism and respiration category. This suggests, on  
232 the one hand, an adaptation to pulmonary macrophage infection for pulmonary variants and on  
233 the other hand, a metabolic adaptation to extra-pulmonary tissues.

234 To go forward, we used PAML branch-site models to test for selection pressures associated  
235 with infection localization. No gene function exhibited a significant signal of positive selection  
236 pressure, neither in pulmonary isolates nor in extra-pulmonary ones with our cohort. Still,  
237 relatively low p-values ( $0.1 < p\text{-value} < 0.15$ ) combined with amino acids with significant BEB  
238 posterior P-values were observed for “Virulence, detoxification and adaptation” gene category  
239 in pulmonary Mtb isolates and for “intermediary metabolism and respiration” gene category in  
240 the extra-pulmonary ones. These tendencies require larger samples to be confirmed.

241

## 242 **DISCUSSION**

243 The objective of the present study was to address the impact of Mtb micro-diversity on TB  
244 pathophysiology. For this purpose, we explored the correlation between Mtb micro-diversity  
245 and TB clinical presentation, meaning pulmonary and extra-pulmonary tuberculosis and  
246 severity of illness. A strong correlation was observed between the detection of Mtb micro-  
247 diversity and TB-associated severity markers, in both pulmonary and extra-pulmonary clinical  
248 isolates. Furthermore, this diversity seems to be driven by a selection pressure to adapt to  
249 different tissues.

250 A previous report explored the correlation between baseline Mtb diversity (initiation or major  
251 change to treatment) and 6-month TB outcome and found that Mtb diversity did not affect TB  
252 outcome. However, beside fatal outcome, treatment failure and acquisition of multi-drug  
253 resistance, authors also included loss of follow-up in unfavorable outcome<sup>9</sup>. Nevertheless, the  
254 latter category, representing 33% of poor outcome in the study, is not itself an unfavorable  
255 outcome, therefore we have classified it in the "unknown" category, which could explain the  
256 discrepancies observed between the studies. Otherwise, our results are in accordance with  
257 another study, focusing on the evolution of Mtb micro-diversity from five patients during the  
258 course of TB treatment, which showed that higher TB severity is associated with an increase of

259 Mtb micro-diversity within-host, and more particularly in pre-mortem Mtb isolates of two of  
260 these patients, and without significant impact of TB treatment on Mtb micro-diversity<sup>10</sup>. Yet it  
261 is still unclear whether Mtb micro-diversity is a cause (better Mtb adaptation to treatment, to  
262 immune pressure and/or to various niches) or a consequence (tissue breakdown allowing  
263 sampling of Mtb variants usually inaccessible and/or lower immune response reducing selection  
264 pressure) of the TB severity. The mechanisms driving such diversity remain to be explored.

265 As shown using PAML branch-site models, variants harbored SNP in different functional  
266 categories according to their localization, meaning pulmonary or extra-pulmonary samples.  
267 Even if the tendencies observed require larger samples to be confirmed, the result obtained  
268 suggested an adaptation to pulmonary macrophage infection for pulmonary Mtb isolates and a  
269 metabolic adaptation for extra-pulmonary ones. Moreover, development of *in vitro* models will  
270 be needed to decipher the role and the impact of the identified pairwise variants between  
271 pulmonary and extra-pulmonary compartments.

272 Alongside that, in the present study, no correlation was found between the ranges of Mtb  $\alpha$ -  
273 diversity and TB-associated severity markers. It may be due to the fact that the analysis was  
274 based on minimum number of variants estimated through WGS data to calculate Mtb  $\alpha$ -  
275 diversity, which was a risk to underestimate micro-diversity in Mtb clinical isolate. However,  
276 to be exhaustive regarding the variant composition of an Mtb isolate, this would have required  
277 sequencing of several colonies for each Mtb clinical isolate, which is time- and cost-consuming.  
278 Nevertheless, detection of unfixed mutations at the level of WGS (meaning mutation  
279 frequencies between 10 and 90%) was enough to observe a strong correlation between Mtb  
280 micro-diversity detection and TB-associated severity markers.

281 As WGS is performed in routine practice in our lab, as well as in other TB diagnosis lab, it  
282 could be envisioned as an all-in-one solution, to detect antibiotic resistance<sup>37</sup>, to infer Mtb  
283 transmission chains and to perform epidemiological monitoring<sup>38,39</sup> and now as a prognosis

284 tool. In the frame of cancer and microbiological research, calling algorithm for low frequency  
285 variants were developed <sup>40-42</sup> and may be adapt to Mtb WGS data. Therapeutic drug monitoring  
286 and implementation of additional management measures could be performed for patients with  
287 detectable Mtb micro-diversity in clinical isolates. It would ensure optimal anti-TB drug doses  
288 and prevent slow response to treatment, which would reduce risks of treatment failure and of  
289 drug resistance acquisition <sup>43,44</sup>.

290 In conclusion, although these results need to be confirmed in an independent prospective  
291 validation study, Mtb micro-diversity within clinical isolate could be a useful prognosis tool to  
292 ensure optimal management of TB patients.

293

## 294 **METHODS**

### 295 *Ethical considerations*

296 For this study we recorded demographical (age, sex), clinical (extrapulmonary and/or  
297 pulmonary TB), microbiological (smear sputum results, growth delay, antibiotic resistance,  
298 lineage) and nutritional and immune data. All data were implemented in a database, in  
299 accordance with the decision 20-216 of the ethics committee of the Lyon University Hospital,  
300 France and the French Bioethics laws (Reference methodology MR-004 that covers the  
301 processing of personal data for purposes of study, evaluation or research that does not involve  
302 the individual). Relevant approval regarding access to patient-identifiable information are  
303 granted by the French data protection agency (Commission Nationale de l'Informatique et des  
304 Libertés, CNIL).

305

### 306 *Mtb samples and data collection*

307 In this retrospective study, a total of 355 Mtb clinical isolates were included, from 311 patients  
308 diagnosed with microbiologically proven TB at the Lyon University Hospital. We included all

309 Mtb clinical isolates for which WGS was performed in routine practice at the Lyon University  
310 Hospital from January 2017 to January 2020 and significant isolates of our collection (MDR  
311 Mtb, representative strains of large previously identified clusters, samples from patients with  
312 both pulmonary and extra-pulmonary TB) which were captured and implemented in the  
313 database during sequencing development in the lab <sup>38,39</sup>. It should be noted that we excluded  
314 pleural TB as depending of the clinical presentation it can be considered as PTB or EPTB.  
315 Microbiological characteristics, such as Mtb lineage, smear results, time to positivity and drug  
316 resistance, were recorded (**Table S1**), as was patients' demographic and baseline descriptive  
317 characteristics (**Table 1**). Regarding patients' descriptive characteristics, only data available  
318 between 2 weeks before TB diagnosis and 1 week after initiation of anti-TB treatment or  
319 nutritional supplementation were considered.

320

### 321 *TB-associated severity indices*

322 Three TB-associated severity indices were assessed in this study.

323 The modified Bandim score is the most largely used PTB prognosis score. It considers 5  
324 symptoms (cough, hemoptysis, dyspnea, chest pain, night sweats) and 5 clinical findings  
325 (anemia, tachycardia, positive finding at lung auscultation, fever, BMI<18 and <16), with one  
326 point for each. One clinical finding was excluded, the mid upper arm circumference (MUAC)  
327 as this data was not available in the Lyon University Hospital. Accordingly, patients were  
328 stratified into two severity classes, mild (Bandim score  $\leq 4$ ) and moderate or severe ( $\geq 5$ ) <sup>26,27,45</sup>.

329 The nutritional status of TB patients was evaluated thanks to the Malnutrition Universal  
330 Screening Tool (MUST), which include three variables: unintentional weight loss score (weight  
331 loss < 5% = 0, weight loss 5–10% = 1, weight loss > 10% = 2), BMI score (BMI>20.0 = 0, BMI  
332 18.5–20.0 = 1, BMI<18.5 = 2) and anorexia (if yes = 2). Malnutrition is frequently observed in

333 patients with PTB and a previous study showed a poorer prognosis for PTB patients with MUST  
334  $\geq 4$ <sup>28</sup>.

335 The nutritional risk score is a four-points score including both nutritional and immune  
336 characteristics: low BMI (<18.5), hypoalbuminaemia (<30.0 g/L), hypocholesterolaemia  
337 (<4mmol/L) and severe lymphocytopenia (<0.7 G cells/L). A high nutritional risk score ( $\geq 3$ )  
338 has been shown to be associated to poor prognosis in PTB<sup>29,46</sup>.

339

#### 340 *Culture of Mycobacterium tuberculosis*

341 The routine laboratory diagnostic workflow consisted of treatment of pulmonary samples with  
342 the modified Kubica's digestion-decontamination method<sup>47</sup>, followed by inoculation in  
343 Mycobacteria Growth Indicator Tubes (MGITs) incubated in a BD BACTEC™ MGIT™  
344 960 instrument (BD, Sparks, MD, USA). Extrapulmonary samples were inoculated using the  
345 same medium without prior decontamination. Mtb genomic DNA extractions were performed  
346 after a single round of culture. Biobank Mtb isolates were inoculated in MGIT until exponential  
347 phase before Mtb genomic DNA extraction.

348

#### 349 *Whole genome sequencing*

350 Genomic DNA of Mtb positive cultures was purified from cleared lysate using a QIAamp DNA  
351 mini Kit (Qiagen). DNA libraries were prepared with Nextera XT kit (Illumina, San Diego,  
352 USA). Samples were sequenced on NextSeq or MiSeq system (Illumina) to produce 150 or 300  
353 base-pair paired-end reads at the Bio-Genet NGS facility of Lyon University Hospital, as  
354 previously described<sup>38</sup>. Reads were mapped with BOWTIE2 to the Mtb H37Rv reference  
355 genome (Genbank NC000962.2) and variant calling was made with SAMtools mpileup, as  
356 previously described<sup>38</sup>.

357



358 ***Illumina data analysis***

359 A valid nucleotide variant was called if the position was covered by a depth of at least 10 reads  
360 and supported by a minimum threshold rate of 10%. Regions with repetitive or similar  
361 sequences were excluded, i.e. regions of PE, PPE, PKS, PPS, ESX genes. The reference genome  
362 coverage breadth was at least 93% with a mean depth of coverage of at least 50x.

363

364 ***Variant assignment***

365 In a previous study, we showed no significant differences in variant detection and frequencies  
366 between sequencing on direct samples and after subculture on media used in routine practice <sup>6</sup>.  
367 Moreover, for this study, 10 isolates were extracted and sequenced twice to evaluate the  
368 variability in mutation frequencies between sequencing experiments. In both sequencing  
369 experiments, 52 unfixed mutations were detected at similar frequencies ( $\pm 10\%$ ), ranging from  
370 10 to 90% (**Fig. S9A**). Accordingly, to identify the minimum number of variant in each Mtb  
371 clinical isolate, a variant was defined as an assembly of mutations at frequencies of  $\pm 10\%$  as  
372 illustrated in **Fig. S9B**.

373

374 ***Mtb  $\alpha$ -diversity indices***

375 The alpha diversity index was calculated by using the software R statistic, Package (vegan)  
376 version 2.5-7.

377

378 ***Selection analysis***

379 We used PAML package to test for selection pressures in our dataset. This method uses  
380 powerful statistics to test for heterogeneous dn/ds ratios at different positions and/or in different  
381 branches. It has already been successfully used in MTC genome to confirm that positions  
382 involved in drug resistance and some positions in membrane proteins are subjected to positive

383 selection<sup>48</sup>. Shortly, this method allows to compare scenarios including (H1) or not (H0) higher  
384 dn/ds ratio on some residues. This means that two categories of sites are set up that have  
385 different distributions of their dn/ds (options can force some of the characteristics of these  
386 distributions, here both M1 and M2 models distributions were explored). The method allows to  
387 identify whether scenario H1 is more likely than H0 using a Likelihood Ratio Test (true if p-  
388 value<0.05).

389 Here we explored whether genes having the same functions (as characterized by available  
390 annotation) have a higher probability to be under selection in the evolutionary branches leading  
391 to one type of localization or the other (pulmonary versus extra-pulmonary). These tests were  
392 performed using the branch-site model of PAML.

393 While compiling the statistics enabling LRT tests described above, PAML package also  
394 associates a probability to each codon that underwent positive selection pressure: the Bayesian  
395 Empirical Bayes (BEB) posterior probability that this codon has a dn/ds ratio higher than 1 as  
396 compared to the alternative scenario, starting from a (true if P>0.95).

397 To do so, we built the matrix including all SNPs for all variants (intra et inter patients  
398 variations). For each SNP, we reconstituted the corresponding amino-acid using annotation  
399 available from mycobrowser Release 2 (<https://mycobrowser.epfl.ch/releases>). The variant  
400 phylogeny was reconstructed using RAxML using a GTRCAT model. All terminal branches  
401 leading to pulmonary isolates were labelled for identifying selection in the lungs (“Pulmonary”  
402 analysis). All terminal branches leading to any extra-pulmonary isolate were labelled for  
403 identifying selection in microaerophilic organs *i.e.* organs other than the lungs (“Extra-  
404 Pulmonary” analysis).

405

#### 406 *Statistical analysis*

#### 407 *Univariate analysis*

408 The study variables were expressed as count (percentage, %) for dichotomous variables and as  
409 medians (interquartile range [IQR]) for continuous values. The number of missing values was  
410 excluded from the denominator. Non-parametric statistical methods Fisher exact test,  $\chi^2$  test,  
411 Mann-Whitney U test and Kruskal-Wallis analysis, using Dunn's Multiple Comparison Test  
412 were used to compare groups, where appropriate. Statistical analyses were performed with  
413 Graph Pad Prism 5. \* $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

414

#### 415 *Response variable and subsetting explanatory variables*

416 We built a model to relate Mtb micro-diversity in pulmonary (P) and extra pulmonary (EP)  
417 isolates to various explanatory variables obtained for most patients in the cohort. The response  
418 variable we chose to model is a binary (0/1) measure of whether Mtb displayed or not some  
419 diversity as observed through counts of unique SNP profiles.

420 Among all the potential explanatory variables recorded in the initial dataset, we retained a  
421 chosen subset, following a series of filters:

422 1. We first removed 15 variables based on three non-exclusive conditions: (i) relatively low  
423 coverage in the cohort, (ii) little to no relationship with the response variable based on the  
424 literature, (iii) important redundancy with other variables in the dataset (e.g. percentage of  
425 neutrophil when the raw neutrophil count was also included);

426 2. Among the remaining variables, we removed all those that had less than 90% coverage (i.e.  
427 those that had missing values for 10% or more of the patients).

428 Apart from the location of the infection (P vs. EP), the final consolidated dataset included 31  
429 variables: 10 categorical variables and 21 quantitative ones. This dataset consisted of 282 rows  
430 (patient:TBlocation), 200 for pulmonary locations and 82 for EP locations. 42 patients had both  
431 microbiologically proven P and EP locations, so the final dataset included 240 patients (158  
432 with only P location, 40 with only EP location, and 42 patients with both locations).

433

434 *Model comparison*

435 In a first exploratory phase, we looked for all models explaining the response variable using a  
436 limited number of explanatory variables. We used generalized linear modelling, assuming that  
437 the binary response variable could be modelled through a binomial distribution and using a logit  
438 transformation linking response to explanatory variables. With 31 possible explanatory  
439 variables, the number of potential models to test is very high and impossible to tackle ( $2^{31}$ , i.e.  
440 more than 2 billion models). In order to reduce this complexity, we chose to restrict our search  
441 to models incorporating between 0 and 12 variables. As presented in the result, this approach  
442 was sufficient to obtain “best models” that had fewer than 12 variables, hence hinting at the  
443 uselessness of pursuing our search further into models of higher complexity.

444

445 To rank the different tested models, we used the Akaike Information Criterion corrected for  
446 small sample sizes (AICc) <sup>49</sup> because we expected many tapered effects of the different  
447 variables <sup>50</sup>. Models with the lowest AICc values were the ones that had the best goodness-of-  
448 fit. Models were computed and compared using the ‘glmulti’ package version 1.0.8 in the R  
449 software version 3.6.3 <sup>51</sup>. To optimize computation time, we first looked for all models with 0  
450 – 9 variables among the 31 present in the dataset, and then looked for models with 10, 11 and  
451 12 variables (thus, 4 parallel uses of glmulti). For each of these four chunks of model  
452 comparison, we retained only the 500 best models (sensu AICc) and gathered all these models  
453 in a consensus pool of good models. The importance of model variables was then assessed using  
454 Akaike weights of each of the 31 variables by summing the Akaike weights of all models  
455 incorporating the focal variable <sup>50,52</sup>. Since each variable was incorporated in exactly half of the  
456 tested models, the importance is expected to be  $\frac{1}{2}$  for variables that did not modify model fit  
457 better than the null expectation. Following Massol et al. (2007)<sup>52</sup>, we thus considered that all

458 variables that were *likely* to have an effect on the response variable were those with an  
459 importance larger than  $1/(1+e^{-1})$ , i.e. 0.73, *plausible* or *implausible* on either side of 0.5, and  
460 unlikely when variable importance was lower than 0.27. For prediction purposes, we retained  
461 all models that were within 2 units of the best model's AICc and used multi-model inference  
462 based on this set of models<sup>50</sup>. We computed unconditional variance using the method of  
463 Johnson and Omland (2004)<sup>53</sup> and obtained confidence intervals on model predictions using  
464 the method suggested by Burnham and Anderson<sup>50</sup>, using function 'predict' in the R package  
465 'glmulti'.

466

#### 467 *Parameter value comparisons*

468 When categorical variables had an effect on the response variable, we tested for pairwise  
469 differences in the coefficients associated to the different levels of the categorical variables. To  
470 do so, we used the R package 'emmeans' version 1.4.5, which tested pairwise differences  
471 between marginal means averaged over all values of other categorical variables (e.g. differences  
472 between "Mtb lineages" were assessed by averaging the effect of "double location"), using  
473 asymptotic test on the z-score obtained from the pair of coefficients<sup>54</sup>.

474

#### 475 *Estimating multi-model errors*

476 Once the set of good models had been determined, we re-sampled the dataset in order to obtain  
477 estimates of the multi-model prediction errors (i.e. false negatives and positives). For a given  
478 target proportion of initial dataset rows to be included in the training set, we drew a random  
479 sample of rows stratified by unique combinations of categorical variables used by the multi-  
480 model, i.e. we made sure that all combinations of factors were included in the training set. This  
481 resulted in actual fraction of data rows in the training dataset slightly higher than the target

482 proportion because some combinations of categorical variable modalities were quite rare and  
483 thus systematically added to the training set.

484 With a given training dataset, we considered the rest of the dataset as validation dataset. We  
485 fitted the whole set of good models using only the training dataset, with possibly estimates of  
486 model coefficients different from those obtained with the whole dataset due to sampling. Based  
487 on the probability of observing some Mtb micro-diversity for all samples in the training set, we  
488 looked for a threshold on this probability that would maximize the true skill statistic (TSS) if  
489 the multi-model were to predict diversity when model predictions were above this threshold  
490 and no diversity otherwise. TSS is a simple statistic (also called informedness or Youden's J  
491 statistic) equal to the sum of sensitivity and specificity minus one. This threshold optimization  
492 procedure was performed using function 'optim.thresh' in R package 'SDMTools' version 1.1-  
493 221.2<sup>55</sup>.

494 The performance of the multi-model inferred from the training dataset was assessed by  
495 predicting the response variable in the validation dataset, using the above-mentioned threshold  
496 for predicting presence/absence of Mtb micro-diversity. This prediction yielded a confusion  
497 matrix (observed vs. predicted absence/presence of Mtb micro-diversity) which was then  
498 analysed using standard statistics, i.e. its sensitivity, specificity, accuracy (1 - error rate) and  
499 TSS.

500

## 501 **ACKNOWLEDGMENTS**

502 This work was supported by the LABEX ECOFECT (ANR-11-LABX-0048) of Université de  
503 Lyon, within the program "Investissements d'Avenir" (ANR-11-IDEX-0007) operated by the  
504 French national research agency (Agence nationale de la recherche, ANR).

505 The funders had no role in study design, data collection and analysis, decision to publish, or  
506 preparation of the manuscript.

507

508 **Competing interests**

509 The authors declare that they have no conflict of interest.

510

511 **REFERENCES**

512 1. WHO, G. Global tuberculosis report 2020. [https://www.who.int/publications-detail-](https://www.who.int/publications-detail-redirect/9789240013131)  
513 [redirect/9789240013131](https://www.who.int/publications-detail-redirect/9789240013131) (2020).

514 2. Gagneux, S. & Small, P. M. Global phylogeography of *Mycobacterium tuberculosis* and  
515 implications for tuberculosis product development. *The Lancet Infectious Diseases* **7**,  
516 328–337 (2007).

517 3. Brown, A. C. *et al.* Rapid Whole-Genome Sequencing of *Mycobacterium tuberculosis*  
518 Isolates Directly from Clinical Samples. *Journal of Clinical Microbiology* **53**, 2230–2237  
519 (2015).

520 4. Casali, N. *et al.* Whole Genome Sequence Analysis of a Large Isoniazid-Resistant  
521 Tuberculosis Outbreak in London: A Retrospective Observational Study. *PLoS Med* **13**,  
522 (2016).

523 5. Doyle, R. M. *et al.* Direct Whole-Genome Sequencing of Sputum Accurately Identifies  
524 Drug-Resistant *Mycobacterium tuberculosis* Faster than MGIT Culture Sequencing.  
525 *Journal of Clinical Microbiology* **56**, (2018).

526 6. Genestet, C. *et al.* Subcultured *Mycobacterium tuberculosis* isolates on different growth  
527 media are fully representative of bacteria within clinical samples. *Tuberculosis (Edinb)*  
528 **116**, 61–66 (2019).

529 7. Ley, S. D., de Vos, M., Van Rie, A. & Warren, R. M. Deciphering Within-Host  
530 Microevolution of *Mycobacterium tuberculosis* through Whole-Genome Sequencing: the  
531 Phenotypic Impact and Way Forward. *Microbiol Mol Biol Rev* **83**, (2019).

- 532 8. Lieberman, T. D. *et al.* Genomic diversity in autopsy samples reveals within-host  
533 dissemination of HIV-associated *Mycobacterium tuberculosis*. *Nature Medicine* **22**,  
534 1470–1474 (2016).
- 535 9. Nimmo, C. *et al.* Dynamics of within-host *Mycobacterium tuberculosis* diversity and  
536 heteroresistance during treatment. *EBioMedicine* **55**, 102747 (2020).
- 537 10. O’Neill, M. B., Mortimer, T. D. & Pepperell, C. S. Diversity of *Mycobacterium*  
538 *tuberculosis* across Evolutionary Scales. *PLOS Pathogens* **11**, e1005257 (2015).
- 539 11. Pérez-Lago, L. *et al.* Whole genome sequencing analysis of inpatient microevolution in  
540 *Mycobacterium tuberculosis*: potential impact on the inference of tuberculosis  
541 transmission. *J. Infect. Dis.* **209**, 98–108 (2014).
- 542 12. Shockey, A. C., Dabney, J. & Pepperell, C. S. Effects of Host, Sample, and in vitro  
543 Culture on Genomic Diversity of Pathogenic *Mycobacteria*. *Front. Genet.* **10**, (2019).
- 544 13. Vargas, R. *et al.* In-host population dynamics of *M. tuberculosis* during treatment failure.  
545 *bioRxiv* 726430 (2019) doi:10.1101/726430.
- 546 14. Votintseva, A. A. *et al.* Same-Day Diagnostic and Surveillance Data for Tuberculosis via  
547 Whole-Genome Sequencing of Direct Respiratory Samples. *Journal of Clinical*  
548 *Microbiology* **55**, 1285–1298 (2017).
- 549 15. Genestet, C. *et al.* *Mycobacterium tuberculosis* micro-diversity promotes intra-  
550 macrophagic persistence and antibiotic tolerance. *bioRxiv* 2020.05.22.110775 (2020)  
551 doi:10.1101/2020.05.22.110775.
- 552 16. Cohen, T. *et al.* Within-Host Heterogeneity of *Mycobacterium tuberculosis* Infection Is  
553 Associated With Poor Early Treatment Response: A Prospective Cohort Study. *J Infect*  
554 *Dis* **213**, 1796–1799 (2016).



- 555 17. Shin, S. S. *et al.* Mixed Mycobacterium tuberculosis-Strain Infections Are Associated  
556 With Poor Treatment Outcomes Among Patients With Newly Diagnosed Tuberculosis,  
557 Independent of Pretreatment Heteroresistance. *J. Infect. Dis.* **218**, 1974–1982 (2018).
- 558 18. Gagneux, S. Ecology and evolution of *Mycobacterium tuberculosis*. *Nature Reviews*  
559 *Microbiology* **16**, 202–213 (2018).
- 560 19. Ailloud, F. *et al.* Within-host evolution of *Helicobacter pylori* shaped by niche-specific  
561 adaptation, intragastric migrations and selective sweeps. *Nat Commun* **10**, 1–13 (2019).
- 562 20. Azarian, T., Ridgway, J. P., Yin, Z. & David, M. Z. Long-Term Intra-host Evolution of  
563 Methicillin Resistant *Staphylococcus aureus* Among Cystic Fibrosis Patients With  
564 Respiratory Carriage. *Front. Genet.* **10**, (2019).
- 565 21. Chaguza, C. *et al.* Within-host microevolution of *Streptococcus pneumoniae* is rapid and  
566 adaptive during natural colonisation. *Nat Commun* **11**, 3442 (2020).
- 567 22. Harkins, C. P. *et al.* The Microevolution and Epidemiology of *Staphylococcus aureus*  
568 Colonization during Atopic Eczema Disease Flare. *J. Invest. Dermatol.* **138**, 336–343  
569 (2018).
- 570 23. Levade, I. *et al.* *Vibrio cholerae* genomic diversity within and between patients. *Microbial*  
571 *Genomics* **3**, e000142 (2017).
- 572 24. Winstanley, C., O'Brien, S. & Brockhurst, M. A. *Pseudomonas aeruginosa* Evolutionary  
573 Adaptation and Diversification in Cystic Fibrosis Chronic Lung Infections. *Trends*  
574 *Microbiol.* **24**, 327–337 (2016).
- 575 25. Wongphutorn, P., Chomvarin, C., Sripan, B., Namwat, W. & Faksri, K. Detection and  
576 genotyping of *Helicobacter pylori* in saliva versus stool samples from asymptomatic  
577 individuals in Northeastern Thailand reveals intra-host tissue-specific *H. pylori* subtypes.  
578 *BMC Microbiol.* **18**, 10 (2018).

- 579 26. Dewi, D. N. S. S., Mertaniasih, N. M. & Soedarsono. SEVERITY OF TB CLASSIFIED  
580 BY MODIFIED BANDIM TB SCORING ASSOCIATES WITH THE SPECIFIC  
581 SEQUENCE OF ESXA GENES IN MDR-TB PATIENTS. *Afr J Infect Dis* **14**, 8–15  
582 (2020).
- 583 27. Rudolf, F. *et al.* The Bandim tuberculosis score: reliability and comparison with the  
584 Karnofsky performance score. *Scand. J. Infect. Dis.* **45**, 256–264 (2013).
- 585 28. Miyata, S., Tanaka, M. & Ihaku, D. The prognostic significance of nutritional status using  
586 malnutrition universal screening tool in patients with pulmonary tuberculosis. *Nutr J* **12**,  
587 42 (2013).
- 588 29. Kim, D. K. *et al.* Nutritional deficit as a negative prognostic factor in patients with miliary  
589 tuberculosis. *European Respiratory Journal* **32**, 1031–1036 (2008).
- 590 30. Casha, A. R. & Scarci, M. The link between tuberculosis and body mass index. *Journal of*  
591 *Thoracic Disease* **9**, E301-E303–E303 (2017).
- 592 31. Cegielski, J. P. & McMurray, D. N. The relationship between malnutrition and  
593 tuberculosis: evidence from studies in humans and experimental animals. *Int. J. Tuberc.*  
594 *Lung Dis.* **8**, 286–298 (2004).
- 595 32. van Crevel, R. *et al.* Decreased plasma leptin concentrations in tuberculosis patients are  
596 associated with wasting and inflammation. *J. Clin. Endocrinol. Metab.* **87**, 758–763  
597 (2002).
- 598 33. Dao, C. N. *et al.* Hyponatremia, hypochloremia, and hypoalbuminemia predict an  
599 increased risk of mortality during the first year of antiretroviral therapy among HIV-  
600 infected Zambian and Kenyan women. *AIDS Res. Hum. Retroviruses* **27**, 1149–1155  
601 (2011).
- 602 34. Vrieling, F. *et al.* Patients with Concurrent Tuberculosis and Diabetes Have a Pro-  
603 Atherogenic Plasma Lipid Profile. *EBioMedicine* **32**, 192–200 (2018).

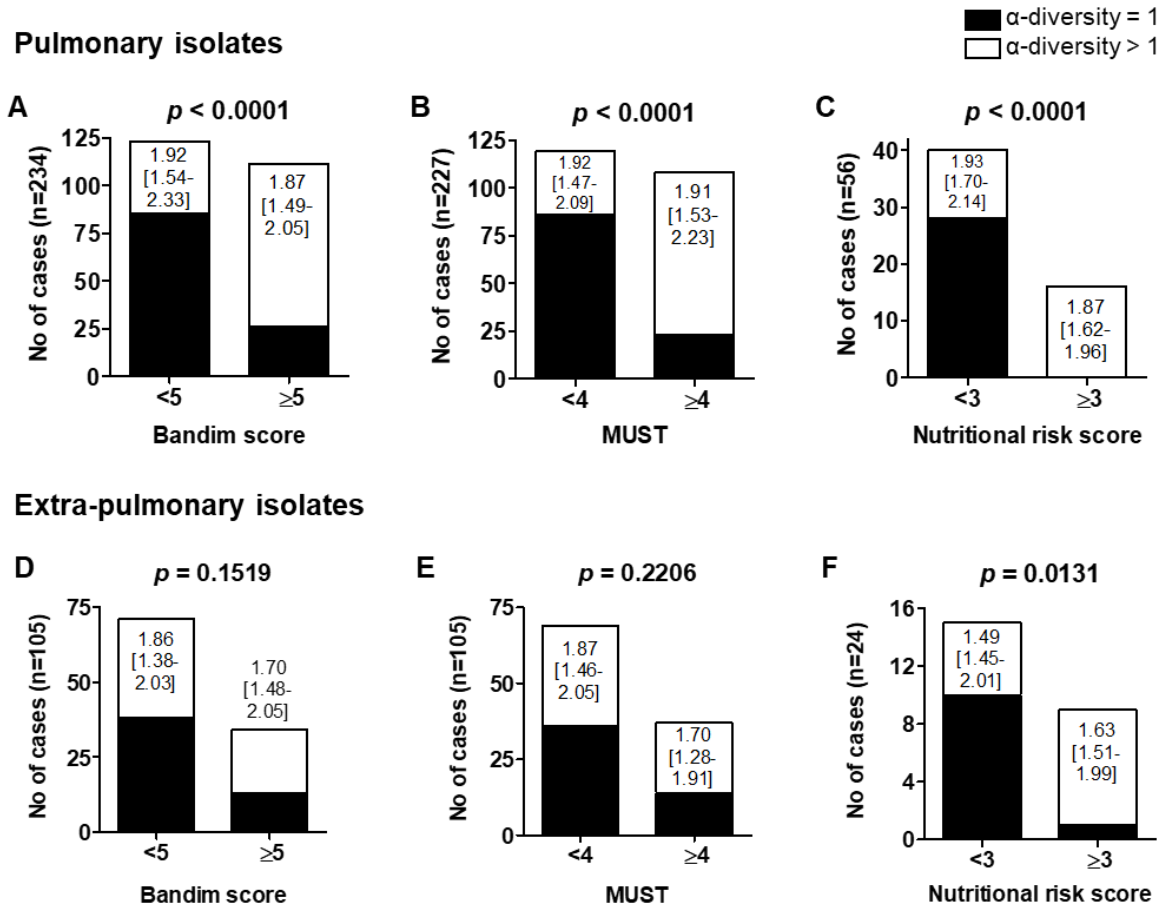
- 604 35. Shivakoti, R., Sharma, D., Mamoon, G. & Pham, K. Association of HIV infection with  
605 extrapulmonary tuberculosis: a systematic review. *Infection* **45**, 11–21 (2017).
- 606 36. Skogmar, S. *et al.* CD4 Cell Levels during Treatment for Tuberculosis (TB) in Ethiopian  
607 Adults and Clinical Markers Associated with CD4 Lymphocytopenia. *PLOS ONE* **8**,  
608 e83270 (2013).
- 609 37. Genestet, C. *et al.* Whole-genome sequencing in drug susceptibility testing of  
610 *Mycobacterium tuberculosis* in routine practice in Lyon, France. *Int. J. Antimicrob.*  
611 *Agents* **55**, 105912 (2020).
- 612 38. Genestet, C. *et al.* Prospective Whole-Genome Sequencing in Tuberculosis Outbreak  
613 Investigation, France, 2017–2018. *Emerging Infectious Diseases* **25**, 589–592 (2019).
- 614 39. Genestet, C. *et al.* Routine survey of *Mycobacterium tuberculosis* isolates reveals  
615 nosocomial transmission. *Eur. Respir. J.* **55**, (2020).
- 616 40. Eliseev, A. *et al.* Evaluation of haplotype callers for next-generation sequencing of  
617 viruses. *Infect. Genet. Evol.* **82**, 104277 (2020).
- 618 41. Spencer, D. H. *et al.* Performance of Common Analysis Methods for Detecting Low-  
619 Frequency Single Nucleotide Variants in Targeted Next-Generation Sequence Data. *J Mol*  
620 *Diagn* **16**, 75–88 (2014).
- 621 42. Xu, C. A review of somatic single nucleotide variant calling algorithms for next-  
622 generation sequencing data. *Computational and Structural Biotechnology Journal* **16**, 15–  
623 24 (2018).
- 624 43. Alsultan, A. & Peloquin, C. A. Therapeutic drug monitoring in the treatment of  
625 tuberculosis: an update. *Drugs* **74**, 839–854 (2014).
- 626 44. Choi, R., Jeong, B.-H., Koh, W.-J. & Lee, S.-Y. Recommendations for Optimizing  
627 Tuberculosis Treatment: Therapeutic Drug Monitoring, Pharmacogenetics, and  
628 Nutritional Status Considerations. *Ann Lab Med* **37**, 97–107 (2017).

- 629 45. Wejse, C. *et al.* TBscore: Signs and symptoms from tuberculosis patients in a low-  
630 resource setting have predictive value and may be used to assess clinical course. *Scand. J.*  
631 *Infect. Dis.* **40**, 111–120 (2008).
- 632 46. Kim, H.-J. *et al.* The impact of nutritional deficit on mortality of in-patients with  
633 pulmonary tuberculosis. *Int. J. Tuberc. Lung Dis.* **14**, 79–85 (2010).
- 634 47. Kent, P. T. & Kubica, G. P. Public Health Mycobacteriology: A Guide for the Level III  
635 Laboratory. | National Technical Reports Library - NTIS.  
636 <https://ntrl.ntis.gov/NTRL/dashboard/searchResults/titleDetail/PB86216546.xhtml> (1985).
- 637 48. Osório, N. S. *et al.* Evidence for Diversifying Selection in a Set of Mycobacterium  
638 tuberculosis Genes in Response to Antibiotic- and Nonantibiotic-Related Pressure.  
639 *Molecular Biology and Evolution* **30**, 1326–1336 (2013).
- 640 49. Hurvich, C. M. & Tsai, C.-L. Regression and time series model selection in small  
641 samples. *Biometrika* **76**, 297–307 (1989).
- 642 50. Burnham, K. P. & Anderson, D. R. *Model Selection and Multimodel Inference: A*  
643 *Practical Information-Theoretic Approach*. (Springer-Verlag, 2002). doi:10.1007/b97636.
- 644 51. Calcagno, V. & Mazancourt, C. de. glmulti: an R package for easy automated model  
645 selection with (generalized) linear models. *Journal of Statistical Software* **34**, 1 (2010).
- 646 52. Massol, F., David, P., Gerdeaux, D. & Jarne, P. The influence of trophic status and large-  
647 scale climatic change on the structure of fish communities in Perialpine lakes. *J Anim*  
648 *Ecol* **76**, 538–551 (2007).
- 649 53. Johnson, J. B. & Omland, K. S. Model selection in ecology and evolution. *Trends in*  
650 *Ecology & Evolution* **19**, 101–108 (2004).
- 651 54. Lenth, R. V. *et al.* *emmeans: Estimated Marginal Means, aka Least-Squares Means*.  
652 (2021).

653 55. VanDerWal, J., Falconi, L., Januchowski, S. & Storlie, L. S. and C. *SDMTools: Species*  
 654 *Distribution Modelling Tools: Tools for processing data associated with species*  
 655 *distribution modelling exercises.* (2014).

656

657 **FIGURES AND LEGENDS**



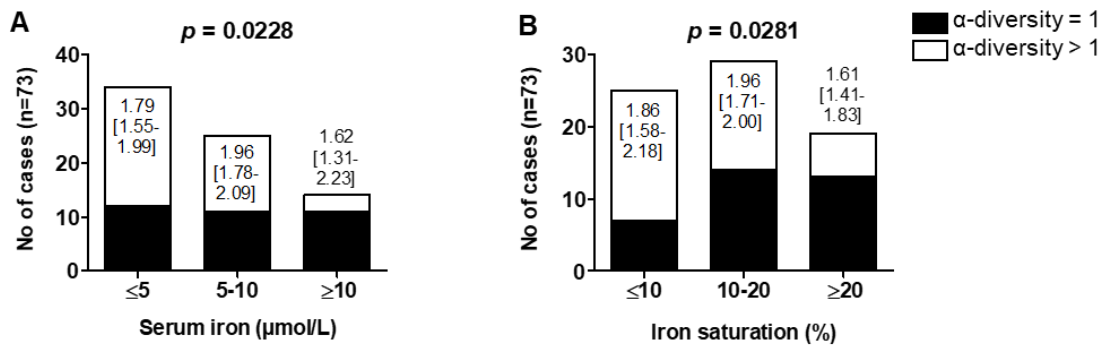
658

659 **Figure 1: The detection of genetic micro-diversity in Mtb isolates is associated with TB-**  
 660 **associated severity indices**

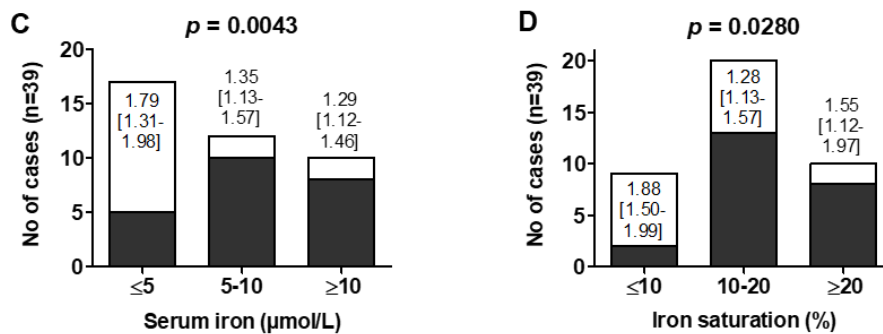
661 Association between the detection and the range of Mtb  $\alpha$ -diversity from pulmonary (A-C) or  
 662 extra-pulmonary samples (D-F) and TB-associated severity indices, the Bandim score (A, D),  
 663 the MUST (malnutrition universal screening tool, B and E), and the nutritional risk score (C,  
 664 F). Black bar:  $\alpha$ -diversity=1 no diversity detected by WGS; white bar:  $\alpha$ -diversity>1 at least  
 665 two variants detected by WGS.  $p = x.xxxx$ : non-parametric statistical method Fisher exact test

666 was used to compare groups. *p-value* < 0.05 was considered significant. x.xx [y.yy-z.zz]:  
667 median [IQR] of Mtb  $\alpha$ -diversity. Mann-Whitney U test was used to compare ranges of Mtb  $\alpha$ -  
668 diversity between groups (no statistical differences observed).  
669

### Pulmonary isolates



### Extra-pulmonary isolates

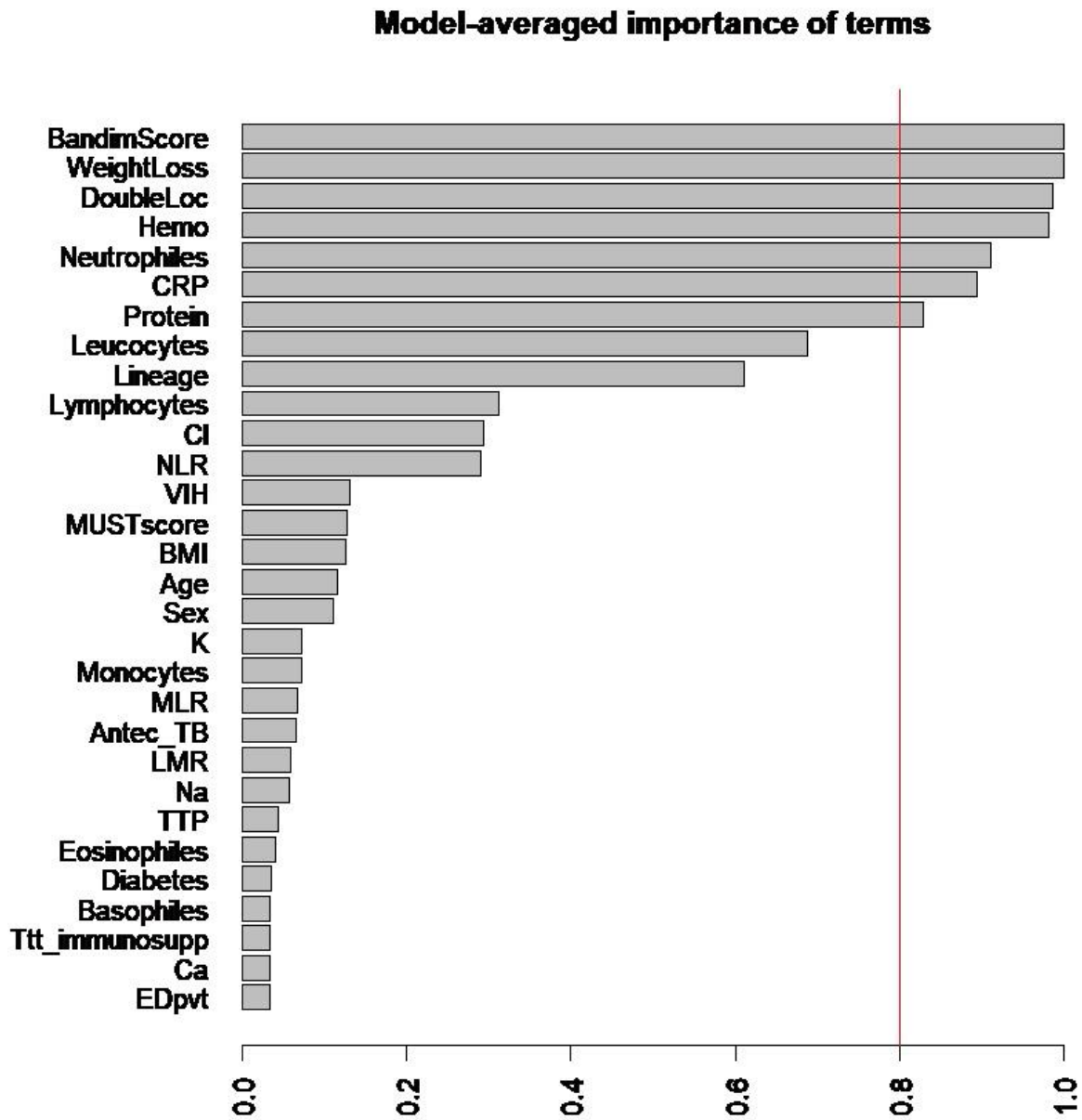


670

671 **Figure 2: The detection of genetic micro-diversity in pulmonary and extra-pulmonary**  
672 **Mtb isolates is associated with iron deficiency**

673 Association between the detection and the range of Mtb  $\alpha$ -diversity from pulmonary (A and B)  
674 and extra-pulmonary (C and D) samples and serum iron level (A and C) and transferrin  
675 saturation (B and D). Black bar:  $\alpha$ -diversity=1 no diversity detected by WGS; white bar:  $\alpha$ -  
676 diversity>1 at least two variants detected by WGS.  $p = \text{x.xxxx}$ : non-parametric statistical  
677 method  $\chi^2$  test was used to compare groups. *p-value* < 0.05 was considered significant. x.xx  
678 [y.yy-z.zz]: median [IQR] of  $\alpha$ -diversity. Kruskal-Wallis analysis, using Dunn's Multiple  
679 Comparison Test was used to compare ranges of  $\alpha$ -diversity between groups (no statistical  
680 differences observed).

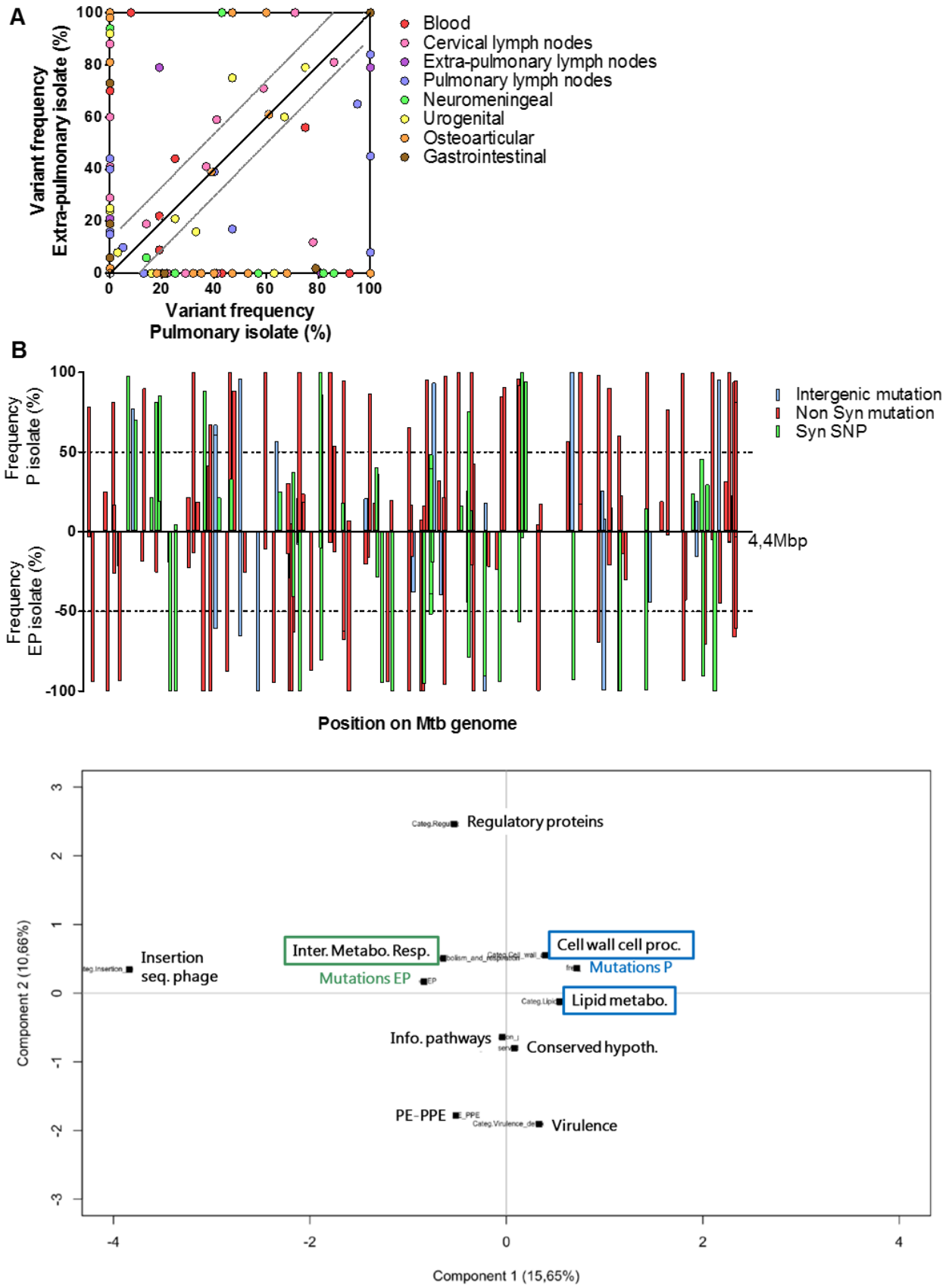
681



682

683 **Figure 3: Model-averaged importance of terms for Mtb micro-diversity in pulmonary**  
684 **isolates.**

685 Model-averaged importance of each term in the model (**Table 2**), which is defined as the  
686 proportion of the 2000 best models in which a given term appears. Red line indicates 80%  
687 support. Terms with an importance above the red line are included in our final model.



688

689 **Figure 4: Compartmentalization of Mtb variants between pulmonary and extra-**  
 690 **pulmonary compartments**



691 (A) Frequencies in pulmonary and extra-pulmonary isolates of the 104 variants identified in  
692 paired isolates from the 42 patients of the training cohort with both microbiologically proven  
693 pulmonary and extra-pulmonary TB. (B) Repartition on Mtb reference genome of the 168  
694 pairwise mutations distance identified between the 42 paired isolates from the training cohort.  
695 Each bar represents a pairwise mutation distance between paired pulmonary and extra-  
696 pulmonary isolates. *y-axis* positive value: mutation frequency in pulmonary isolate; *y-axis*  
697 negative value: mutation frequency in extra-pulmonary isolate. Bleu bar: intergenic mutation;  
698 red bar: nonsynonymous mutation; green bar: synonymous single nucleotide polymorphism  
699 (SNP). (C) Principal component analysis (PCA) for mixed data was performed with the delta  
700 of frequency of the 104 nonsynonymous pairwise mutation observed between paired pulmonary  
701 and extra-pulmonary isolates and the functional categories of these mutations. The first  
702 principal component allowed a discrimination of pulmonary and extra-pulmonary mutations.