

1 **Novel SNP-based assay for genotyping *Mycobacterium avium* subsp. *paratuberculosis***

2

3 Celia Leão^{a,b,#}, Robert J Goldstone^c, Josephine Bryant^d, Joyce McLuckie^a, João Inácio^{e,f,#},
4 David GE Smith^{a,c} and Karen Stevenson^a

5

6 Moredun Research Institute, Pentlands Science Park, Penicuik, United Kingdom^a; Unidade
7 Estratégica de Investigação e Serviços em Produção e Saúde Animal, Instituto Nacional de
8 Investigação Agrária e Veterinária, I.P. (INIAV, I.P.), Lisboa, Portugal^b; Institute of
9 Infection, Inflammation & Immunity, University of Glasgow, Glasgow, United Kingdom^c;
10 Division of Infection and Immunity, University College London, London, United Kingdom^d;
11 Unidade de Microbiologia Médica, Instituto de Higiene e Medicina Tropical, Universidade
12 Nova de Lisboa, Lisboa, Portugal^e; School of Pharmacy and Biomolecular Sciences,
13 University of Brighton, Brighton, United Kingdom^f

14

15 **Running Head: SNP assay for genotyping *M. a. paratuberculosis***

16

17 #Address correspondence to:

18 Célia Leão, celia.leao@iniav.pt

19 Instituto Nacional de Investigação Agrária e Veterinária, I.P. (INIAV, I.P.), Unidade

20 Estratégica de Investigação e Serviços em Produção e Saúde Animal, Lisboa, Portugal

21 João Inácio, J.Inacio@brighton.ac.uk

22 School of Pharmacy and Biomolecular Sciences, University of Brighton, 815 Cockcroft

23 building, Lewes road, Brighton BN2 4GJ, United Kingdom

24

25 **ABSTRACT**

26 Typing of *Mycobacterium avium* subspecies *paratuberculosis* strains presents a challenge
27 since they are genetically monomorphic and traditional molecular techniques have limited
28 discriminatory power. The recent advances and availability of whole genome sequencing has
29 extended possibilities for the characterization of *Mycobacterium avium* subspecies
30 *paratuberculosis* and it can provide a phylogenetic context to facilitate global epidemiology
31 studies. In this study we developed a SNP assay based on polymerase chain reaction and
32 restriction enzyme digestion or sequencing of the amplified product. The SNP analysis was
33 performed using genome sequence data from 133 *Mycobacterium avium* subspecies
34 *paratuberculosis* isolates with different genotypes from eight different host species and
35 seventeen distinct geographic regions around the world. A total of 28402 SNPs were
36 identified among all the isolates. The minimum number of SNPs required to distinguish
37 between all the 133 genomes was 93 and between only the Type C isolates was 41. To reduce
38 the number of SNPs and PCRs required we adopted an approach based on sequential
39 detection of SNPs and a decision tree. By the analysis of 14 SNPs *Mycobacterium avium*
40 subspecies *paratuberculosis* isolates can be characterized within 14 phylogenetic groups with
41 a higher discriminatory power compared to MIRU-VNTR assay and other typing methods.
42 Continuous updating of genome sequences are needed in order to better characterize new
43 phylogenetic groups and SNP profiles. The novel SNP assay is a discriminative, simple,
44 reproducible method and requires only basic laboratory equipment for the large-scale global
45 typing of *Mycobacterium avium* subspecies *paratuberculosis* isolates.

46

47

48 INTRODUCTION

49 *Mycobacterium avium* subspecies *paratuberculosis* (*Map*) causes Johne's disease, a chronic
50 infectious enteritis principally of ruminants. The disease occurs worldwide and is responsible
51 for significant losses to the livestock industry. *Map* also has been detected in a subset of
52 human patients with Crohn's disease (1) although the zoonotic role of the bacterium remains
53 controversial.

54 Strain typing is a prerequisite for tracing the sources of infection and for studying the
55 epidemiology, population structure and evolutionary relationships between isolates. It can
56 also reveal the genetic diversity underlying important phenotypic characteristics such as host
57 specificity, pathogenicity, antibiotic resistance and virulence. Typing of *Map* strains presents
58 a challenge since *Map*, like *Mycobacterium tuberculosis*, is genetically monomorphic (2).

59 Genetic diversity among *Map* strains has been investigated using molecular techniques such
60 as restriction fragment length polymorphism and IS900 analysis (IS900 RFLP) (3), pulsed-
61 field gel electrophoresis (PFGE) (4), amplified fragment length polymorphism (AFLP)
62 analysis (5), random amplified polymorphic DNA (RAPD) analysis (6), mycobacterial
63 interspersed repetitive unit – variable number tandem repeat (MIRU-VNTR) analysis (7) and
64 short-sequence repeat (SSR) analysis (8). However, these techniques have limited
65 discriminatory power when applied to *Map* and although this can be increased by combining
66 complementary genotyping techniques, it is often insufficient for accurately determining
67 relationships among isolates or global epidemiological studies (9; 10).

68 Whole genome sequencing (WGS) provides the ultimate resolution of isolates and, unlike the
69 techniques above, it can provide a phylogenetic context to facilitate global epidemiology
70 studies and affirm epidemiological connections (10; 11; 12). Although WGS is becoming
71 cheaper, it is still too expensive to be used for routine genotyping of *Map* isolates and
72 requires robust data handling and analysis processes. Single nucleotide polymorphisms
73 (SNPs) have been used successfully to type several genetically monomorphic pathogens,

74 including *M. tuberculosis* (13), *Mycobacterium bovis* (14), *Salmonella enteritica* Typhi (15)
75 and *Yersinia pestis* (16). SNP assays have been used to discriminate between *Map* strain
76 Types I, II and III (17; 18) and between strains derived from animal and human hosts (19).
77 However, these assays were based on a limited number of SNPs, many of which were not
78 informative when applied to a wider wild-type population. The purpose of this study was to
79 develop a SNP assay that is discriminative, practicable and reproducible for the large-scale
80 global typing of *Map* isolates. Hence we developed a SNP typing method based on
81 polymerase chain reaction (PCR) and restriction enzyme digestion, which would minimize
82 costs and require only basic laboratory equipment. Additionally, we adopted an approach
83 based on sequential detection of SNPs and a decision tree to reduce the number of SNPs and
84 PCRs required.

85

86 **MATERIALS AND METHODS**

87 **Selection of genome sequences and *Map* strains.** Genome sequences from 133 *Map* isolates
88 generated in a previous study (10) were selected for SNP and phylogenetic analyses. This
89 panel was chosen to maximise genetic diversity and reduce phylogenetic discovery bias (2)
90 since it comprised isolates with different genotypes (determined by multiplex PFGE and
91 MIRU-VNTR), which were selected from 17 different countries and isolated from eight
92 different host species. The panel also included isolates representing the major strain types that
93 have been described previously (reviewed by 20; 21). The sequence reference numbers are
94 given in Figures 1 and 2. The field isolates (n=26) for which genome sequence data was not
95 available and that were used to validate the SNP assay are shown in Table 1.

96

97 **Preparation of genomic DNA.** Genomic DNA from the UK field isolates (n=10) (Table 1)
98 was extracted from plugs used to perform the PFGE analysis. Briefly, half of each plug was
99 washed three times with 2 ml of Tris EDTA buffer pH 8 (10 mM Tris.Cl, 1 mM EDTA) for

100 10 min with shaking. After washing the plugs, 0.5-1 ml of sterile deionised water was added
101 and the agarose was melted at 70°C. The suspension was stored at -20°C until required for
102 PCR. DNA from the Portuguese field isolates (n=16) (Table 1) was extracted from pure
103 cultures grown in Middlebrook 7H9 medium supplemented with 10% OADC (Oleic
104 Albumin Dextrose Catalase) and 2 mg/l of mycobactin as previously described (10) using the
105 QIAamp DNA Mini Kit (Qiagen) according to manufacturer's instructions with minor
106 modifications. Briefly, the bacterial suspension was centrifuged at 5000 × g for 10 min and
107 the pellet was resuspended in 180 µl of ATL buffer. Zirconium beads (0.1 mm) were added
108 to the mixture and mechanical disruption of the cells was performed twice with a FastPrep
109 FP120 Bio101 bead shaker (Savant Instruments Inc., Holbrook, NY) at 6.5 msec⁻¹ for 45
110 seconds. The disrupted mixture was cooled on ice for 15 min and 20 µl of proteinase K was
111 added. The remaining procedure was completed according to manufacturer's instructions.
112 The DNA was eluted with 200 µl of AE buffer and stored at -20°C until required for PCR.

113

114 **SNP analysis and phylogenomics.** The genome sequence of *Map* K10 (22; Accession
115 number NC_002944.2) was used as the reference genome for the phylogenetic analyses. Raw
116 genome sequence data for the *Map* isolates are available from the European Nucleotide
117 Archive under accession PRJEB2204. The assembly of the reads into contigs was performed
118 using Velvet assembler program (freely available at <https://www.ebi.ac.uk/~zerbino/velvet/>)
119 and the alignment of the sequences and positioning of SNPs was executed utilising MUMmer
120 package (freely available at <http://mummer.sourceforge.net/>) using *Map* K10 as reference
121 sequence. SNPs were then extracted and concatenated and a phylogenetic tree was calculated
122 using phyML (freely available at <http://atgc.lirmm.fr/phyml>). This phylogenetic tree was then
123 imported into R using the ade4 package to explore the 'paths' which exist between the
124 genomes. These 'paths' are a description of the strains which group on the same branch as the
125 structure of the phylogenetic tree descends (i.e. all strains are included on the 'root branch',

126 which then bifurcates to include a subset of strains on one branch and the remaining strains
127 on the other, and so on). Using these ‘paths’ to assign groups of genomes, SNPs which were
128 shared by all strains on each branch were compiled. This data then allowed comparison of the
129 SNPs present in reducing sized groups of genome sequences (as the tree descends through
130 specific branches, the number of genomes remaining in each subsequent group becomes
131 progressively smaller). This data permitted the detection of discriminating SNPs which are
132 present in one group of sequences yet absent in others. A ‘set cover’ analysis was then
133 performed on this dataset to calculate the minimum number of SNPs which were required to
134 discriminate between the groups of strains at each selected level of the tree (see below).
135 These SNPs were then selected and taken forward for validation as below. SNPs were named
136 according to their base position in the revised *Map* K10 genome sequence.

137

138 **Set cover approximation and decision tree construction.** The set cover problem asks,
139 given a universe of elements to be covered (in this case, the universe was a series of pair-wise
140 combinations of all genome groups), and a number of ‘sets’ which contain varying
141 combinations of the elements to be covered (in this case, these ‘sets’ are specific SNPs which
142 are found to discriminate between one or more of the pair-wise combinations of genome
143 groups), what is the minimum number of sets required to cover all the elements of the
144 universe. In this case, the set cover asks what is the minimum number of SNPs required to
145 ‘fill’ all the pairwise combinations of genome groups. In this way, the defined set of SNPs
146 would discriminate between every possible pairwise combination of genome groups. The set
147 cover problem is a class of decision problem known to be NP-complete, meaning the time
148 taken to reach an exact solution increases exponentially with problem size. Due to the large
149 number of total SNPs and genome groups to consider, we opted to use an approximation
150 algorithm known as the greedy solution to estimate the set cover. Greedy algorithms work by
151 iteratively selecting the set which covers the largest number of uncovered universe elements.

152 The greedy approach is considered a suitable polynomial time approximation for the set
153 cover problem. This analysis resulted in the identification of a set of SNPs which could
154 distinguish between all pairwise combinations of genome groups.
155 The chosen SNPs which comprised the set cover formed the basis of the decision tree for
156 determining group assignment of genomes, since every bifurcation of the phylogenetic tree
157 could be discriminated by at least one SNP present in the set cover. The structure of the
158 phylogenetic tree was then manually investigated in light of the SNPs within the set cover to
159 minimise the number of SNPs required to determine between phylogenetically relevant
160 branches. These SNPs were then validated as outlined below.

161

162 **SNP validation and primer design.** Primers for PCR amplification of the selected SNPs
163 were designed using the online software Primer3 (<http://bioinfo.ut.ee/primer3-0.4.0/>) and the
164 revised *Map* K10 genome sequence (23; Accession number AE016958.1) (Table 2). Primers
165 were designed to be 18-20 mer with a melting temperature between 63 and 67°C.
166 PCR reactions were carried out in 50 µl containing 200 µM of each deoxynucleotide
167 triphosphate (Invitrogen), 0.5 µM of each primer (Table 2), 1 U of Phusion® High- Fidelity
168 DNA polymerase, 1× Phusion GC buffer (New England Biolabs), and 4 µl of extracted DNA
169 (5 ng/µl). Amplification was performed in a TC-PLUS Thermal cycler (Techne) with an
170 initial step at 98 °C for 3 min, followed by 35 cycles at 98 °C for 30 sec, 63-67 °C for 30 sec
171 (annealing temperatures provided in Table 2) and 72 °C for 40 sec, ending with a step at 72
172 °C for 10 min. The amplified products were electrophoretically analysed in a 1.5% (w/v)
173 agarose gel stained with SYBR® Safe DNA Gel Stain (Life Technologies) in 0.5× Tris-
174 Borate-EDTA (TBE) buffer. Gel electrophoresis images were acquired with an
175 Alphaimager™ 2200 (Alpha Innotech). DNA ladder IV (Bioline) and *Map* K10 (positive
176 control) were included on each gel.

177 Restriction endonuclease analysis of PCR products was performed according to the
178 manufacturer's instructions using one unit of restriction enzyme and 10 µl of amplified
179 product in a total reaction volume of 25 µl. All restriction endonucleases were purchased
180 from New England Biolabs (Table 2). Restricted products were detected by electrophoresis
181 on 1.5% (w/v) agarose gels as described above.

182 Confirmation of the presence of the SNPs was obtained by sequencing the PCR products.
183 PCR product (40µl) was purified using QIAquick PCR purification Kit (Qiagen) according to
184 manufacturer's instructions. Sequencing of PCR products was carried out by Eurofins
185 Genomics (MWG-Biotech) using the same primers used to the amplification of the
186 fragments. Confirmation of the presence or absence of the SNP in the expected position of
187 the genome was achieved using Basic Local Alignment Search Tool (Blast -
188 <http://blast.ncbi.nlm.nih.gov/Blast.cgi>). The phylogenetic profile for each isolate was
189 obtained by combining the results for all SNPs.

190

191 **Discriminatory power of SNP-based genotyping assay.** The discriminatory power of the
192 assay was calculated using the Hunter-Gaston Discriminatory Index (HGDI), according to
193 Hunter and Gaston (24).

194

195 **RESULTS**

196 The genome sequence data from 133 *Map* strains were compared to the reference *Map* strain
197 K10 to identify SNPs. A total of 28402 SNPs were identified among all the isolates. A
198 phylogenetic tree was generated based on the SNP analysis and distinct phylogenetic groups
199 were identified (Figures 1 and 2), which conformed to the broadly recognised phylogenetic
200 structure of *Map* (10). By using an adaptation of the 'set cover' problem, the minimum
201 number of SNPs required to discriminate between all the isolates was calculated to be 93. To

202 refine the number of SNPs to a number manageable for routine laboratory procedures, we
203 considered a strategy based on sequential detection of SNPs and a decision tree (Figure 3).
204 The phylogenetic analysis distinguished two major strain groups corresponding to those
205 previously designated Type C and Type S (Figure 1). We identified a SNP (snp3842359),
206 which could be detected using BsmB1 (Table 2) that would discriminate between these two
207 groups. Type C strains have an 'A' and Type S strains either a 'G' or 'C' at base position
208 3842359. This constituted the first step in the decision tree, the next step being determined
209 according to whether the isolate was Type C or Type S (Figure 3).

210 For further analysis of Type S strains, we identified a SNP (snp343677), which could be
211 detected using AvaII that discriminated the Type S sub-groups Type I and Type III (Table 2,
212 Figures 1 and 3). Additionally, snp3842359 also could be used to distinguish Type I and
213 Type III strains since Type I have a "G" and Type III have a "C" at base position 3842359
214 (Figure 3), which could be detected by PCR amplification and sequencing of the product.

215 The Type C group comprised the majority of the isolates and the high homogeneity within
216 this group posed a challenge for identification of clade specifying SNPs. Firstly, SNP
217 analysis was repeated with only the sequence data from the 115 Type C isolates (Figure 2)
218 and the minimum number of SNPs required to distinguish between these 115 isolates was
219 determined to be 41. We then considered three principal sub-groups designated Bison (as
220 reported previously, 10) and A and B as shown in Figures 1 and 2. We identified a SNP
221 (snp50173), which could distinguish the Bison group from both sub-groups A and B using
222 ApoI (Table 2, Figures 1 and 3). This constituted step 2 in the decision tree (Figure 3).

223 Within the Bison group, Indian Bison type could be differentiated from US Bison type using
224 snp2327379 and further differentiation of the Indian bison type was possible using
225 snp305277 (Table 2, Figure 3). Due to the limited number of Bison type strains available
226 extensive verification of these SNPs was not possible in this study.

227 To further discriminate between Type C isolates in sub-groups A and B, we identified a SNP
228 (snp4111202), which could be detected using FatI (Table 2, Figures 1 and 3). Due to the
229 small number of isolates in sub-group B, we did not, at this stage, seek additional SNPs to
230 further discriminate the isolates within this group. This constituted step 3 in the decision tree
231 (Figure 3).

232 Within the larger subgroup A, SNPs were identified that could discriminate ten groups
233 (snp3879247, snp2939977, snp1932058, snp1327872, snp3844632, snp1966028, snp305277,
234 snp4339946, snp2087274 and snp1686154, Tables 2 and 3, Figures 2 and 3). These SNPs
235 were verified by sequencing of the PCR products and comparison of the sequences with the
236 reference *Map* K10 strain. It was not possible to identify SNPs with specific restriction
237 endonuclease sites for the clades differentiated using snp2939977 and snp1932058 but all
238 other SNPs could be detected by restriction endonuclease analysis (Table 2).

239

240 **SNP validation and genotyping.** For the validation of the selected SNPs comprising the
241 decision tree (Figure 3), DNA from isolates belonging to Type C (Bison group:
242 MAPMRI029, MAPMRI031, MAPMRI127, MAPMRI034, MAPMRI117 and MAPMRI026;
243 Sub-group A: MAPMRI110, MAPMRI120 and MAPMRI027; Sub-group B: MAPMRI059,
244 MAPMRI136 and MAPMRI091) and Type S (Type I: MAPMRI007 and MAPMRI001; Type
245 III: MAPMRI051, MAPMRI045 and MAPMRI047) (Figure 2) were used for amplifying
246 products containing SNPs (as indicated in Figure 3) and the PCR products were digested with
247 the correspondent restriction enzymes (Table 2). PCR products were also purified and
248 sequenced to confirm SNPs. To assess the validity of the ten SNPs for discriminating the
249 clades within Sub-group A, thirty isolates from the original panel of 115 sequenced strains
250 were re-tested. These were selected to be representative of the ten different phylogenetic
251 clades as shown in Figure 2 and were subjected to analyses for all 14 SNPs (Figure 3), all of

252 which were confirmed to be present. The *Map* isolates were grouped into SNP profiles as
253 shown in Table 3 and Figure 3.

254 A further 26 *Map* isolates (Table 1) not previously sequenced or typed using this SNP assay
255 were genotyped using the 14 SNPs. These isolates belonged to four phylogenetic groups
256 within the Sub-group A. Significantly, ten isolates from different geographic regions of
257 United Kingdom with the same MIRU-VNTR and PFGE profile were classified into three
258 different SNP profiles, one of which was not identified in the original phylogenetic analysis
259 and therefore represented a new SNP profile (SNP11) (Table 3). Two UK isolates were
260 identified to belong to the SNP1 profile, six isolates to SNP9 and two isolates to profile
261 SNP11. The 16 Portuguese isolates were distributed among three phylogenetic groups: all the
262 isolates from Azores were present in the same group identified by profile SNP3; two isolates
263 from the same region in the north of Portugal were identified in the same phylogenetic group
264 as six isolates from the United Kingdom (profile SNP2); and, the remaining four Portuguese
265 isolates from the same region were found to belong to the new phylogenetic group together
266 with the two DNAs from the United Kingdom (profile SNP11) (Table 1).

267

268 **Discriminatory power of SNP-based genotyping assay.** The discriminatory power of the
269 SNP-based assay was 0.8390 for the 56 isolates that were used for the validation of the assay.
270 In order to compare the discriminatory power of the SNP assay with MIRU-VNTR analysis,
271 we used the typing results for 46 isolates, which had been typed using both methods. The
272 HGDI was calculated to be 0.8135 for the SNP assay and 0.6386 for MIRU VNTR.

273

274 **DISCUSSION**

275 Several methods have been used to characterize *Map* strains but they have some limitations.
276 Techniques based on the analysis of total genomic DNA such as RFLP and PFGE require
277 culture of the isolates to prepare moderate amounts of high quality DNA and are therefore

278 slow, can be technically demanding, labour intensive, hard to standardize and expensive.
279 Furthermore, RFLP and PFGE can clearly distinguish between Type C and S but do not give
280 sufficient discrimination within these strain types for detailed epidemiological studies.
281 Techniques such as AFLP and RAPD employ PCR to detect smaller genomic DNA
282 fragments but are less utilised for epidemiological studies due to difficulties in
283 standardisation, reproducibility and limited discriminative power. Other typing methods
284 based on repetitive sequences such as SSR and MIRU-VNTR are popular due to their ease of
285 use and rapidity but are again limited with respect to their ability to discriminate within the
286 two major strain types and the typing results may not reflect the evolutionary relationships
287 between isolates (10; 12; 25; 26).

288 In this study a novel typing assay based on SNP analysis by PCR and restriction or
289 sequencing of the amplified products was developed. This technique is easy to perform, is
290 applicable to a small quantity of genomic DNA and is based on standard PCR and restriction
291 endonuclease analysis. It was possible to refine the number of SNPs to a number manageable
292 for routine laboratory procedures by adopting an approach based on sequential detection of
293 SNPs via a decision tree. The SNP assay was highly discriminative, possessing a higher
294 discriminatory power than MIRU-VNTR when applied to 46 *Map* isolates.

295 SNP-based typing assays are particularly useful for monomorphic pathogens that exhibit
296 limited genetic diversity. Furthermore, they have the advantage that they can be used to
297 determine phylogenetic relationships, unlike techniques based on mobile or repetitive DNA
298 elements, which interrogate a relatively small proportion of the mycobacterial genome and
299 can exhibit homoplasy. However, SNP discovery is subject to phylogenetic discovery bias
300 (2), a phenomenon well described for *M. tuberculosis* (27) and *Bacillus anthracis* (28), and is
301 most likely to be encountered where information is missing on strains geographically-
302 restricted or belonging to rare phylogenetic groupings. For this reason, we utilised a large
303 collection of global isolates, which had been previously genotyped by classical molecular

304 tools (PFGE and MIRU-VNTR) to maximise genetic diversity and include representatives of
305 all previously reported strain types. The SNPs identified in this study should provide the
306 necessary means to unambiguously classify *Map* strains within this global framework. Even
307 so, the composition of any panel of SNPs needs to be reviewed or augmented once additional
308 groups of strains that were not included in the initial analysis are discovered. This has been
309 illustrated in this study with the discovery of a new phylogenetic group represented by profile
310 SNP 11 comprising six isolates when an additional uncharacterised 26 isolates were screened
311 using the SNP assay. WGS needs to be performed and the sequence comparisons and SNP
312 analysis repeated to determine the positions of these isolates within the phylogenetic tree and
313 determine any additional SNPs that could be used to define the group. In this study, we
314 concentrated on finding SNPs to differentiate within Type C strains in sub-group A. The SNP
315 assay needs to be expanded to differentiate between strains within Type S, Bison type and
316 Type C subgroup B, but WGS data from more strains belonging to these phylogenetic groups
317 are required for SNP discovery to provide a phylogenetically robust framework for strain
318 differentiation combined with sufficient discriminatory power for detailed genetic studies.
319 SNPs have been described in previous studies that differentiate between the two major
320 phylogenetic groups, Type S and Type C, and between *Map* strain Types I and III. A PCR-
321 REA assay described by Whittington et al. (29) based on a SNP at base position 223 in the
322 *IS1311* insertion sequence has been used extensively for discriminating between Type S and
323 Type C strains. However, in a recent study (10) the *IS1311*-REA incorrectly identified strain
324 MAPMRI074 as a Type S strain when WGS and SNP analysis clearly confirmed it to be
325 Type C, suggesting that the C to T allelic variation at base pair position 223 in *IS1311*
326 occurred after the initial divergence of Type C from Type S strains. The SNP identified in
327 this study, snp3842359, and the corresponding restriction endonuclease BsmBI for PCR-REA
328 SNP detection could provide a more reliable alternative assay for differentiating Type S and
329 C strains.

330 *IS1311* SNP analysis has also been used to distinguish Bison type strains from non-Bison
331 Type C and Type S strains (30). In Bison-type strains, all copies of *IS1311* have a “T” at base
332 pair position 223, whereas the non-Bison Type C strains have one or more copies with a “C”
333 or “T” at the same position. Copy number with respect to this allele is not always easy to
334 assess and can be very variable (10). Snp50173 and the corresponding restriction
335 endonuclease ApoI for PCR-REA SNP detection could provide an easier, alternative assay
336 for discriminating Bison-type strains from other Type C strains.

337 A study published by Castellanos and colleagues (17) developed a PCR-REA to detect a SNP
338 present on the *gyrB* gene at base position 1626 that allowed discrimination of Type III from
339 Type I and II strains. Additional SNPs in the *gyrA* gene were identified that could
340 differentiate Types I and III from Type II. In our study it is possible to use only a single SNP
341 to discriminate Type I, Type II and Type III based on the amplification and sequencing of a
342 fragment containing the snp3842359 where in the same position of the genome Type I strains
343 have a “G”, Type II have an “A” and Type III have a “C”. This is an improvement compared
344 with the system previously reported.

345 In conclusion, we developed a novel SNP-based genotyping assay based on the analysis of 14
346 SNPs that can be used to characterize *Map* isolates within 14 phylogenetic groups with a
347 higher discriminatory power compared to MIRU-VNTR assay and other typing methods. We
348 adopted an approach based on sequential detection of SNPs and a decision tree based on
349 PCR-restriction enzyme digestion to reduce the number of SNPs and required PCRs. This
350 novel assay can overcome some issues regarding the genotyping of isolates characterized as
351 Type I, Type III and Bison type. Continuous updating of genome sequences are needed in
352 order to better characterize new phylogenetic groups and SNP profiles. The novel SNP assay
353 is a discriminative, simple, reproducible method and requires only basic laboratory
354 equipment for the large-scale global typing of *Map* isolates.

355

356 **Acknowledgements**

357 This work was supported by the Scottish Government Rural and Environment Science and
358 Analytical Services Division. Célia Leão is a recipient of a PhD grant from “Fundação para a
359 Ciência e a Tecnologia” SFRH/BD/62469/2009.

360

361 **Conflict of interest statement**

362 The authors declare that there are no conflicts of interest.

363

364 **References**

- 365 1. **Naser SA, Sagrainsingh SR, Naser AS, Thanigachalam S.** 2014. *Mycobacterium*
366 *avium* subspecies *paratuberculosis* causes Crohn’s disease in some inflammatory
367 bowel disease patients. *World J Gastroenterol* **20**:7403-7415.
- 368 2. **Achtman M.** 2008. Evolution, population structure, and phylogeography of
369 genetically monomorphic bacterial pathogens. *Ann Rev Microbiol* **62**:53-70.
- 370 3. **Pavlik I, Horvathova A, Dvorska L, Bartl J, Svastova P, du Maine R, Rychlik I.**
371 1999. Standardisation of restriction fragment length polymorphism analysis for
372 *Mycobacterium avium* subspecies *paratuberculosis*. *J Microbiol Meth* **38**:155-167.
- 373 4. **Stevenson K, Hughes VM, de Juan L, Inglis NF, Wright F, Sharp JM.** 2002.
374 Molecular characterization of pigmented and non-pigmented isolates of
375 *Mycobacterium avium* subspecies *paratuberculosis*. *J Clin Microbiol* **40**:1798-1804.
- 376 5. **Motiwala AS, Strother M, Amonsin A, Byrum B, Naser SA, Stabel JR, Shulaw**
377 **WP, Bannantine JP, Kapur V, Sreevatsan S.** 2003. Molecular epidemiology of
378 *Mycobacterium avium* subsp. *paratuberculosis*: Evidence for limited strain diversity,
379 strain sharing, and identification of unique targets for diagnosis. *J Clin Microbiol*
380 **41**:2015-2016.

- 381 6. **Pillai SR, Jayarao BM, Gummo JD, Hue Jr EC, Tiwari D, Stabel JR, Whitlock**
382 **RH.** 2001. Identification and sub-typing of *Mycobacterium avium* subsp. *avium* by
383 randomly amplified polymorphic DNA. *Vet Microbiol* **79**:275-284.
- 384 7. **Thibault VC, Grayon M, Boschioli ML, Hubbans C, Overduin P, Stevenson K,**
385 **Gutierrez MC, Supply P, Biet F.** 2007. New variable-number tandem-repeats
386 markers for typing *Mycobacterium avium* subsp. *paratuberculosis* and *M. avium*
387 strains: comparison with IS900 and IS1245 restriction fragment length polymorphism
388 typing. *J Clin Microbiol* **45**:2404-2410.
- 389 8. **Sevilla I, Li L, Amonsin A, Garrido J, Geijo MV, Kapur V, Juste RA.** 2008.
390 Comparative analysis of *Mycobacterium avium* subsp. *paratuberculosis* isolates from
391 cattle, sheep and goats by short sequence repeat and pulsed-field gel electrophoresis
392 typing. *BMC Microbiol* **8**:204.
- 393 9. **Stevenson K, Álvarez J, Bakker D, Biet F, de Juan L, Denham S, Dimareli Z,**
394 **Dohmann K, Gerlach G-F, Heron I, Kopecna M, May L, Pavlik I, Sharp JM,**
395 **Thibault VC, Willemsen P, Zadoks R, Greig A.** 2009. Occurrence of
396 *Mycobacterium avium* subspecies *paratuberculosis* across host species and European
397 countries with evidence for transmission between wildlife and domestic ruminants.
398 *BMC Microbiol* **9**:212.
- 399 10. **Bryant JM, Thibault VC, Smith DGE, McLuckie J, Heron I, Sevilla IA, Biet F,**
400 **Harris SR, Maskell DJ, Bentley SD, Parkhill J, Stevenson K.** 2015. Phylogenomic
401 exploration of the relationships between strains of *Mycobacterium avium* subspecies
402 *paratuberculosis*. *BMC Genomics* In press
- 403 11. **Bryant JM, Schürch AC, van Deutekom H, Harris SR, de Beer JL, de Jager V,**
404 **Kremer K, van Hijum SA, Siezen RJ, Borgdorff M, Bentley SD, Parkhill J, van**
405 **Soolingen D.** 2013. Inferring patient to patient transmission of *Mycobacterium*
406 *tuberculosis* from whole genome sequencing data. *BMC Infect Dis* **13**:100.

- 407 12. **Ahlstrom C, Barkema HW, Stevenson K, Zadoks RN, Biek R, Kao R, Trewby H,**
408 **Hauptstein D, Kelton DF, Fecteau G, Labrecque O, Keefe GP, McKenna SLB,**
409 **Buck JD.** 2015. Limitations of variable number of tandem repeat typing identified
410 through whole genome sequencing of *Mycobacterium avium* subsp. *paratuberculosis*
411 on a national and herd level. *BMC Genomics* **16**:161.
- 412 13. **Stucki D, Malla B, Hostettler S, Huna T, Feldmann J, Yeboah-Manu D, Borrell**
413 **S, Fenner L, Comas I, Coscollà M, Gagneux S.** 2012. Two new rapid SNP-Typing
414 methods for classifying *Mycobacterium tuberculosis* Complex into the main
415 phylogenetic lineages. *PLOS ONE* **7**:e41253.
- 416 14. **Joshi D, Harris NB, Waters R, Thacker T, Mathema B, Krieswirth B, Sreevatsan**
417 **S.** 2012. Single Nucleotide Polymorphisms in the *Mycobacterium bovis* genome
418 resolve phylogenetic relationships. *J Clin Microbiol* **50**:3853-3861.
- 419 15. **Octavia S, Lan R.** 2007. Single-Nucleotide-Polymorphism typing and genetic
420 relationships of *Salmonella enterica* Serovar Typhi Isolates. *J Clin Microbiol*
421 **45**:3795–3801.
- 422 16. **Morelli G, Song Y, Mazzoni CJ, Eppinger M, Roumagnac P, Wagner DM,**
423 **Feldkamp M, Kusecek B, Vogler AJ, Li Y, Cui Y, Thomson NR, Jombart T,**
424 **Leblois R, Lichtner P, Rahalison L, Petersen JM, Balloux F, Keim P, Wirth T,**
425 **Ravel J, Yang R, Carniel E, Achtman M.** 2010. *Yersinia pestis* genome sequencing
426 identifies patterns of global phylogenetic diversity. *Nat Genet* **42**:1140-1143.
- 427 17. **Castellanos E, Aranaz A, Romero B, de Juan L, Álvarez J, Bezos J, Rodríguez S,**
428 **Stevenson K, Mateos A, Domínguez L.** 2007. Polymorphisms in *gyrA* and *gyrB*
429 genes among *Mycobacterium avium* subsp. *paratuberculosis* Type I, II, and III
430 isolates. *J Clin Microbiol* **45**:3439-3442.

- 431 18. **Gastaldelli M, Stefani E, Lettini AA, Pozzato N.** 2011. Multiplexed typing of
432 *Mycobacterium avium* subsp. *paratuberculosis* Type I, II and III by Luminex xMAP
433 suspension array. J Clin Microbiol **49**:389-391.
- 434 19. **Wynne JW, Bellera C, Boyda V, Francisc B, Gwoźdźd J, Carajiasd M, Heinea**
435 **HG, Wagnere J, Kirkwoode CD, Michalskia WP.** 2014. SNP genotyping of animal
436 and human derived isolates of *Mycobacterium avium* subsp. *paratuberculosis*. Vet
437 Microbiol **172**:479-485.
- 438 20. **Stevenson K.** 2010. Comparative Differences between Strains of *Mycobacterium*
439 *avium* subsp. *paratuberculosis* p126-137. In Behr MA, Collins DM (eds),
440 Paratuberculosis Organism, Disease, Control. CAB International, Oxfordshire, UK.
- 441 21. **Stevenson K.** 2015. Genetic diversity of *Mycobacterium avium* subsp.
442 *paratuberculosis* and the influence of strain type on infection and pathogenesis: A
443 review. Vet Res **46**:64.
- 444 22. **Li L, Bannantine JP, Zhang Q, Amonsin A, May BJ, Alt D, Banerji N, Kanjilal**
445 **S, Kapur V.** 2005. The complete genome sequence of *Mycobacterium avium*
446 subspecies *paratuberculosis*. Proc Natl Acad Sci USA **102**:12344-12349.
- 447 23. **Wynne JW, Seeman T, Bulach D, Coutts SA, Talaat AM, Michalski WP.** 2010.
448 Re-sequencing the *Mycobacterium avium* subspecies *paratuberculosis* K10 genome:
449 improved annotation and revised genome sequence. J Bacteriol **192**:6319-6320.
- 450 24. **Hunter PR, Gaston MA.** 1988. Numerical index of the discriminatory ability of
451 typing systems: an application of Simpson's index of diversity. J Clin Microbiol
452 **26**:2465-2466.
- 453 25. **Castellanos E, de Juan L, Domínguez L, Aranaz A.** 2012. Progress in molecular
454 typing of *Mycobacterium avium* subspecies *paratuberculosis*. Res Vet Sci **92**:169-
455 179.

- 456 26. **Sevilla I, Garrido JM, Geijo M, Juste RA.** 2007. Pulsed-field gel electrophoresis
457 profile homogeneity of *Mycobacterium avium* subsp. *paratuberculosis* isolates from
458 cattle and heterogeneity of those from sheep and goats. *BMC Microbiol* 7:18.
- 459 27. **Gagneux S, Small PM.** 2007. Global phylogeography of *Mycobacterium tuberculosis*
460 and implications for tuberculosis product development. *Lancet Infect Dis* 7:328-337.
- 461 28. **Pearson T, Busch JD, Ravel J, Read TD, Rhoton SD, U'Ren JM, Simonson TS,**
462 **Kachur SM, Leadem RR, Cardon ML, Van Ert MN, Huynh LY, Fraser CM,**
463 **Keim P.** 2004. Phylogenetic discovery bias in *Bacillus anthracis* using single-
464 nucleotide polymorphisms from whole-genome sequencing. *Proc Natl Acad Sci USA*
465 **101:**13536-13541.
- 466 29. **Whittington RJ, Marsh IB, Choy E, Cousins D.** 1998. Polymorphisms in *IS1311*,
467 an insertion sequence common to *Mycobacterium avium* and *M. avium* subsp.
468 *paratuberculosis*, can be used to distinguish between and within these species. *Mol*
469 *Cell Probe* **12:**349-358.
- 470 30. **Whittington RJ, Marsh IB, Whitlock RH.** 2001. Typing of *IS1311* polymorphisms
471 confirms that bison (*Bison bison*) with paratuberculosis in Montana are infected with
472 a strain of *Mycobacterium avium* subsp *paratuberculosis* distinct from that occurring
473 in cattle and other domesticated livestock. *Mol Cell Probes* **15:**139-145.
- 474

475 **Figure 1.** Whole genome SNP-based phylogenetic tree of 133 *Map* isolates included in this
476 study (strain sequence reference MAPMRI numbers are indicated). Previously described
477 lineages and sub-groups A and B described in this study are highlighted in grey.

478

479 **Figure 2.** Whole genome SNP-based phylogenetic tree of 115 Type C *Map* isolates (strain
480 sequence reference MAPMRI numbers are indicated). Ten clades within the phylogenetic
481 sub-group A can be distinguished by PCR and sequencing following the analysis of ten SNPs
482 (grey boxes). Black circles indicate the 30 strains used for the validation of the method.

483

484 **Figure 3.** Work flow with a schematic representation of the decision tree with the sequential
485 numbered steps and restriction enzymes, SNPs positions, expected bases and SNP profiles
486 obtained based on the SNPs analysis. * SNP that can be detected only by sequencing of the
487 amplified product. # Bison type strains from India. § Bison type strains from US.

488 |

Table 1. Additional *Map* field isolates used for validation of SNP assay

Isolate	Host	Geographic location	Multiplex PFGE profile	INMV profile*	SNP profile (this study)
F043531	cattle	Northern Ireland (UK)	2.1	1	1
F012398	cattle	Cumbria (UK)	2.1	1	11
F005713	cattle	Wiltshire (UK)	2.1	1	11
C217551	cattle	Selkirkshire (UK)	2.1.	1	9
C216785/2	cattle	Mid Glamorgan (UK)	2.1	1	9
C219376	cattle	Shropshire (UK)	2.1	1	9
C221325	cattle	Gwynedd (UK)	2.1	1	9
C524656	cattle	Aberdeenshire (UK)	2.1	1	1
C326442	cattle	Ayr (UK)	2.1	1	9
C216962/6	cattle	Leicestershire (UK)	2.1	1	9
C1	goat	São Miguel (Azores, PT)	ND	ND	3
C2	goat	São Miguel (Azores, PT)	ND	ND	3
C4	goat	São Miguel (Azores, PT)	ND	ND	3

C7	goat	São Miguel (Azores, PT)	ND	ND	3
C9	goat	São Miguel (Azores, PT)	ND	ND	3
C11	goat	São Miguel (Azores, PT)	ND	ND	3
C13	goat	São Miguel (Azores, PT)	ND	ND	3
C14	goat	São Miguel (Azores, PT)	ND	ND	3
C16	goat	São Miguel (PT)	ND	ND	3
C4A4	goat	São Miguel (Azores, PT)	ND	ND	3
B1	cattle	Vila do Conde (PT)	ND	2	9
B3	cattle	Vila do Conde (PT)	ND	2	9
B13	cattle	Póvoa do Varzim (PT)	ND	2	11
B18	cattle	Póvoa do Varzim (PT)	ND	2	11
B21	cattle	Póvoa do Varzim (PT)	ND	2	11
B22	cattle	Póvoa do Varzim (PT)	ND	2	11

ND - not determined; UK – United Kingdom; PT – Portugal. *Profile number assigned by INRA and published in an online database <http://mac->

inmv.tours.inra.fr

Table 2. PCR primers and restriction endonucleases used in the SNP assay for this study

SNP name	Primer name	Primer sequence (5' to 3')	Annealing temperature (°C)	Product size (bp)	Enzyme	Restriction results§
snp3842359*	MAP_F1	CACCTGGCCAAGTACTACCA	63	528	BsmBI	(A) Type C - 528 bp
	MAP_R1	GCGATGTCATGATGCTGCTG				(G/C) Type S - 312, 216 bp
snp343677	TypeS_F	AACACCAGGATCGCGTTCTT	65	511	Avall	(G) Type I - 297, 152, 62 bp
	TypeS_R	CAATTAGCGGTCGAGTCGTC				(A) Type III - 293, 218 bp
snp50173	BisonF	GGACGATTACTCGGTTCCAG	63	469	ApoI	(T) Bison group - 226, 192, 51 bp
	BisonR	ACCCGTGTTCCGGCTACCT				(G) Cattle group - 277, 192 bp
snp4111202	SNP4_F	GTCAGAAACATCCCGCCTTC	65	461	FatI	(G) Type C sub-group A – 284, 177 bp
	SNP4_R	GTATTGAGTGAGGCAAGCGG				(C) Type C sub-group B - 461 bp
snp3879247	SNP5_F	GTTGATCGACAGCGAGTGC	64	465	BlpI/DdeI	(C) – 227, 238 bp
	SNP5_R	GTGGTGTCCGAGGTGAACTT				(T) – 465 bp
snp2939977	SNP6_F	TATCTCCAAGGACGCATTCC	64	516	-	-
	SNP6_R	CTGCCATGTCCGTCCTTAAT				
snp1932058	SNP7_F	GGCTTGAAACTCCAAGTCT	63	452	-	-
	SNP7_R	CGTCGTACATCCTCGTGGT				
snp1327872	SNP8_F	GCGCTTGTTGACAGGTTGA	64	528	AvaI/BsoBI	(G) – 292, 171, 65 bp
	SNP8_R	TACGACGAAGACCCCGACTA				(T) – 463, 65 bp
snp3844632	SNP9_F	GATCGATGCGGAGCTCGT	64	457	FatI	(G) – 457 bp
	SNP9_R	TGACAGGAAGGTCCATAGCC				(C) – 241, 217 bp
snp1966028	SNP10_F	GTCGAGGGCTTCCAGGTT	67	427	SapI/EarI	(A) – 427 bp
	SNP10_R	GTCTGAGGCCAGCGACAC				(C) – 246, 182 bp
snp305277 [†]	SNP11_F	CCATCCCGAGTTCAACAAGT	64	461	BspMI/BfuAI	(G) – 310, 151 bp
	SNP11_R	ACTTGTCGGGGTTGTAGCTG				(A) – 461 bp
snp4339946	SNP12_F	AACCGCTCAAGGCGAAAG	64	464	BstEII	(T) – 292, 172 bp
	SNP12_R	TCCCTTATCTGCGAAGTGCT				(A) – 464 bp
snp2087274	SNP13_F	CAGACCGAGCACCTCCTG	65	453	HpyAV	(C) – 453 bp
	SNP13_R	CCGCGTTGAAGGATCTCAAG				(A) – 227, 226 bp
snp1686154	SNP14_F	GAATCCCCGGAAGTGGTG	65	525	MscI	(G) – 525 bp

SNP14_R GCAGTCCAGATAACGGAACG

(A) – 284, 241 bp

*SNP can also distinguish between Type I and III (Type I has a "G" and Type III has a "C" at base position 3842359 detectable by sequencing the PCR product); †SNP can also distinguish between two phylogenetic subgroups of Bison group isolates (Figure 3); § in parenthesis is the expected base at the correspondent SNP position; - not applicable.

Table 3. SNP profiles of Type C *Map* isolates in phylogenetic sub-group A used in this study

Phylogenetic group	SNP profile	No. isolates verified	Base at SNP position*													
			3842359	343677	50173	4111202	3879247	2939977	1932058	1327872	3844632	1966028	305277	4339946	2087274	1686154
Reference base (K10)			A	A	G	G	C	G	G	G	G	A	G	T	C	A
Clade 1	1	7	A	A	G	G	C	G	A	G	G	A	G	T	C	G
Clade 2	2	1	A	A	G	G	C	A	G	G	G	A	G	T	C	G
Clade 3	3	11	A	A	G	G	T	G	G	G	G	A	G	T	C	G
Clade 4	4	2	A	A	G	G	C	G	G	T	G	A	G	T	C	G
Clade 5	5	1	A	A	G	G	C	G	G	G	C	A	G	T	C	G
Clade 6	6	1	A	A	G	G	C	G	G	G	G	A	A	T	C	G
Clade 7	7	1	A	A	G	G	C	G	G	G	G	C	G	T	C	G
Clade 8	8	1	A	A	G	G	C	G	G	G	G	A	G	A	C	G
Clade 9	9	16	A	A	G	G	C	G	G	G	G	A	G	T	A	G
Clade 10	10	9	A	A	G	G	C	G	G	G	G	A	G	T	C	A
new clade	11	6	A	A	G	G	C	G	G	G	G	A	G	T	C	G

Total no. 56

* SNP position in the revised *Map* K10 sequence (Accession number: AE016958.1).

Defining SNP base marked in bold type

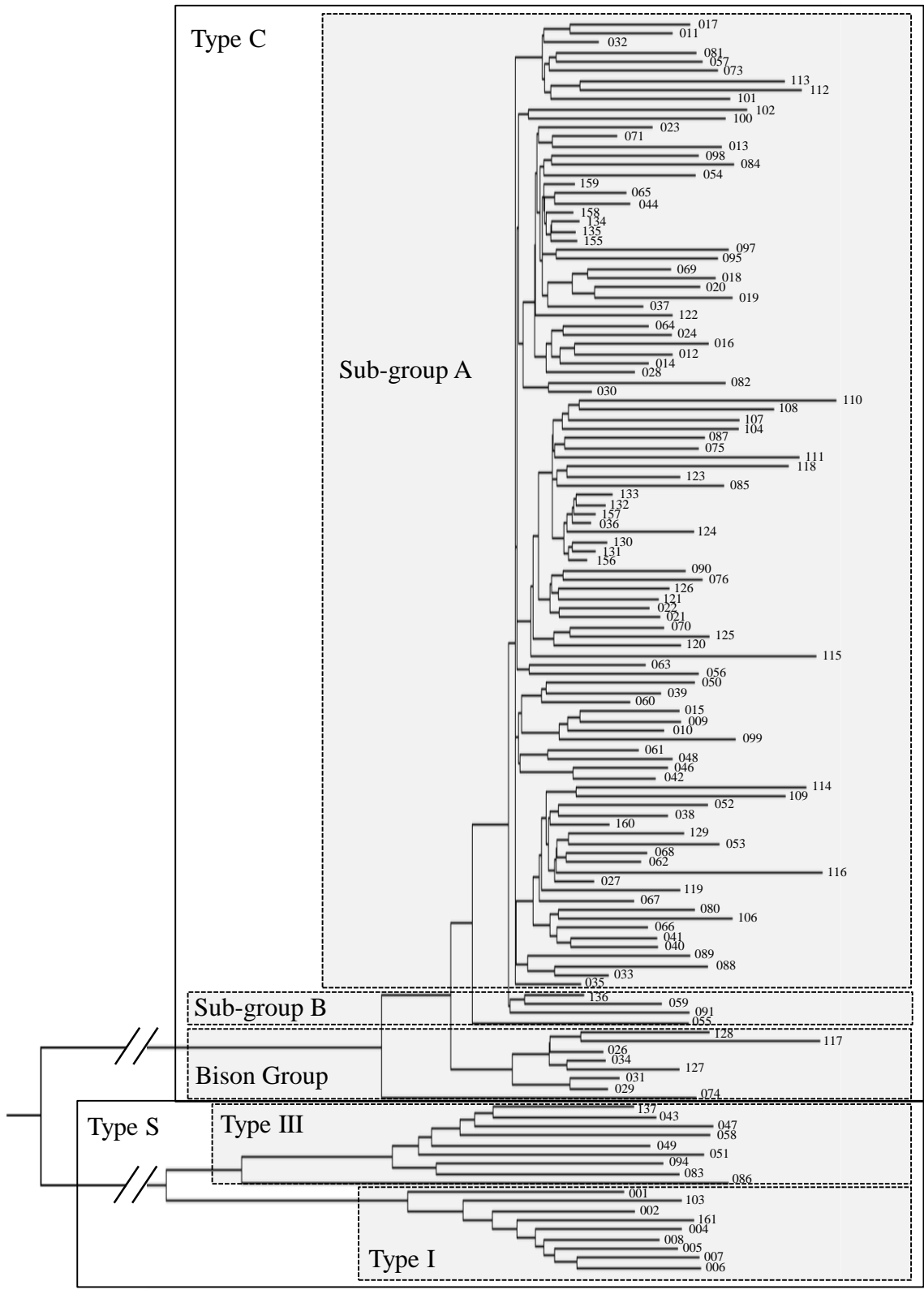


Fig. 1

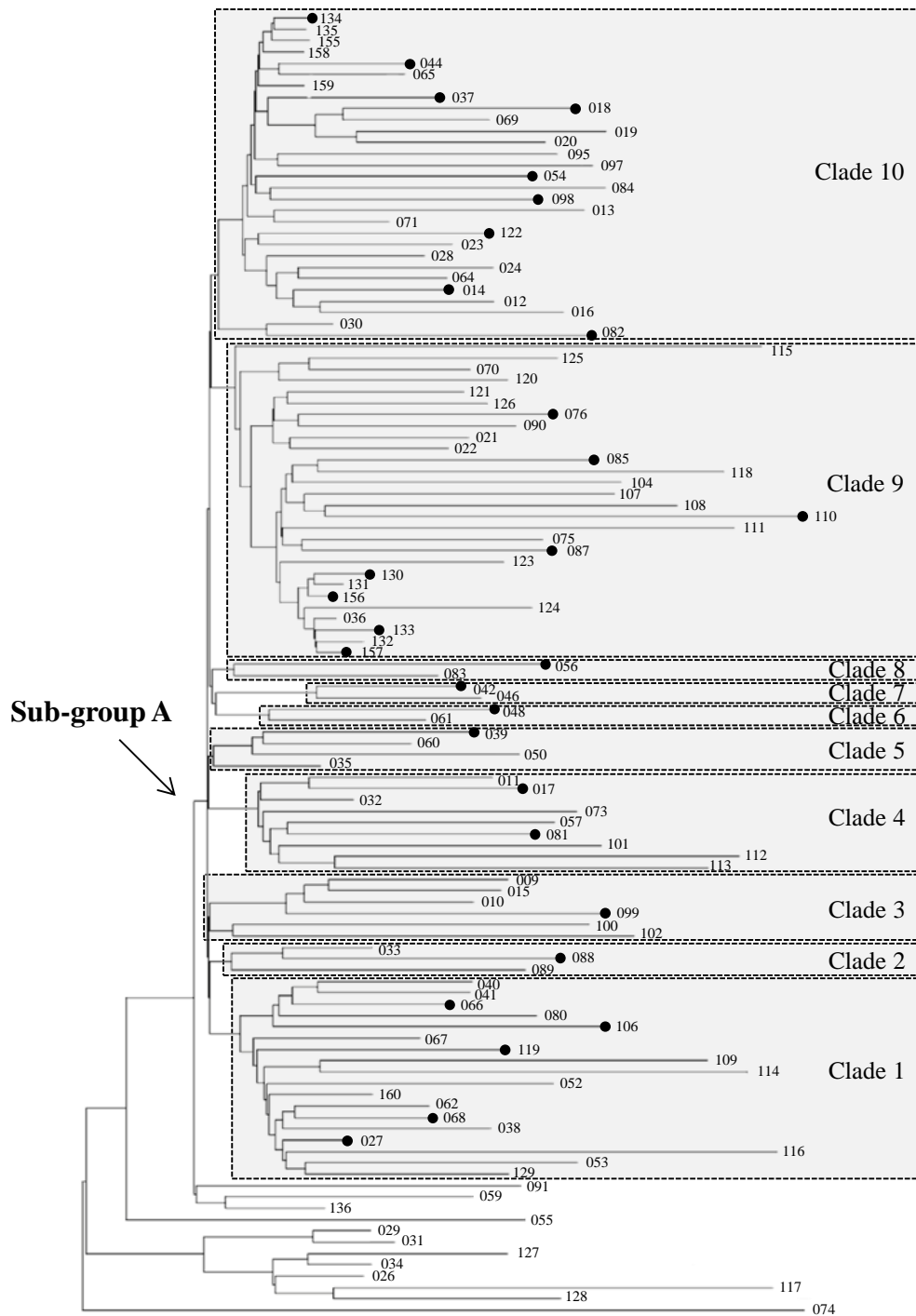


Fig. 2

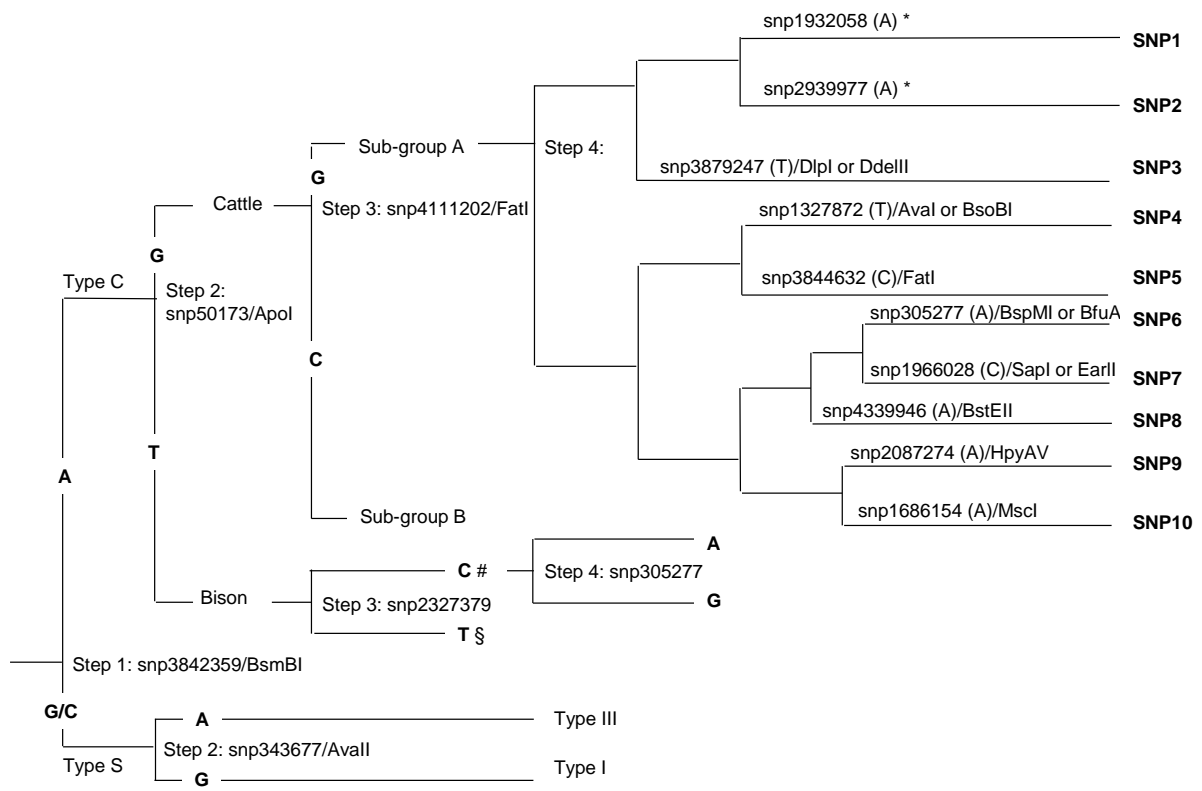


Fig. 3