

1 **Combining GWAS and F_{ST} -based approaches to identify targets of**
2 ***Borrelia*-mediated selection in natural rodent hosts**

3

4 Luca Cornetti^{1,2} & Barbara Tschirren^{3*}

5

6 ¹Department of Evolutionary Biology and Environmental Studies, University of

7 Zurich, Zurich, Switzerland

8 ²Zoological Institute, University of Basel, Basel, Switzerland

9 ³Centre for Ecology and Conservation, University of Exeter, Treliever Road,

10 Penryn, TR10 9FE, United Kingdom

11

12 *Correspondence:

13 Barbara Tschirren, Email: b.tschirren@exeter.ac.uk

14

15 Running title: Targets of *Borrelia*-mediated selection

16 **Abstract**

17 Recent advances in high-throughput sequencing technologies provide
18 opportunities to gain novel insights into the genetic basis of phenotypic trait
19 variation. Yet to date, progress in our understanding of genotype-phenotype
20 associations in non-model organisms in general and natural vertebrate
21 populations in particular has been hampered by small sample sizes typically
22 available for wildlife populations and a resulting lack of statistical power, as
23 well as a limited ability to control for false positive signals. Here we propose to
24 combine a genome-wide association (GWAS) and F_{ST} -based approach with
25 population-level replication to partly overcome these limitations. We present a
26 case study in which we used this approach in combination with Genotyping-
27 by-Sequencing (GBS) SNP data to identify genomic regions associated with
28 *Borrelia afzelii* resistance or susceptibility in the natural rodent host of this
29 Lyme disease-causing spirochete, the bank vole (*Myodes glareolus*). Using
30 this combined approach we identified four consensus SNPs located in exonic
31 regions of the genes *Slc26a4*, *Tns3*, *Wscd1* and *Espnl*, which were
32 significantly associated with the voles' *Borrelia* infectious status within and
33 across populations. Functional links between host responses to bacterial
34 infections and most of these genes have previously been demonstrated in
35 other rodent systems, making them promising new candidates for the study of
36 evolutionary host responses to *Borrelia* emergence. Our approach is
37 applicable to other systems and may facilitate the identification of genetic
38 variants underlying disease resistance or susceptibility, as well as other
39 ecologically relevant traits, in wildlife populations.

40

41 **Keywords:** host-parasite interactions, wild immunogenetics, pathogen-
42 mediated selection, evolutionary change, RAD sequencing (RAD-seq),
43 conservation genetics
44

45 **Introduction**

46 Testing evolutionary theories of host-parasite interactions and resistance
47 evolution is hampered by a lack of understanding of the genetic architecture
48 of host defence and susceptibility (Lazzaro & Little, 2009). This is particularly
49 the case for natural populations of non-model organisms, for which very little
50 functional genetic information is currently available (Spurgin & Richardson,
51 2010). Infectious diseases are a major cause of wildlife population declines
52 (Smith, Sax, & Lafferty, 2006) and pose a substantial threat to global
53 biodiversity (MacPhee & Greenwood, 2013). At the same time, wildlife
54 diseases can spillover into human populations and are thus of public health
55 concern (Daszak, Cunningham, & Hyatt, 2000; Jones et al., 2008;
56 Wiethoelter, Beltrán-Alcrudo, Kock, & Mor, 2015). A better understanding of
57 the genetic basis of wildlife disease resistance or susceptibility is thus crucial
58 for conservation efforts (Jones et al., 2007; Margres et al., 2018) but also to
59 understand and predict the dynamics of zoonotic diseases (Beldomenico &
60 Begon, 2010; Price, Spencer, & Donnelly, 2015).

61 Recent advances in high-throughput sequencing technologies have
62 made it possible to obtain extensive genomic information for non-model
63 organisms (Tagu, Colbourne, & Nègre, 2014). Yet challenges to associate
64 genetic variants with phenotypic traits of interest remain formidable (e.g. Hong
65 & Park, 2012). Two main approaches are typically used to identify genomic
66 regions of interest. The first approach, genome-wide association
67 (GWAS)(Amos, Driscoll, & Hoffman, 2011), tests for associations between
68 phenotypic traits of interest and genomic variants across the whole genome
69 (Petersen, Fredrich, Hoepfner, Ellinghaus, & Franke, 2017). A second

70 approach (F_{ST} outlier approach; Vitti, Grossman, & Sabeti, 2013) used to
71 identify genomic regions associated with phenotypic trait variation is based on
72 the premise that natural selection acting on a locus of interest will result in
73 differences in allele frequencies among populations subject to different
74 environmental conditions or showing different phenotypes. The F_{ST} outlier
75 approach identifies regions with unusually large (when compared to the
76 genome-wide F_{ST} distribution) genetic differentiation between populations
77 suggesting that they are under selection.

78 Although these two approaches have been successfully used to
79 identifying genomic regions of interest in humans (e.g. Visscher et al., 2017)
80 and model organisms (e.g. Flint & Eskin, 2012; Togninalli et al., 2018;
81 Wangler, Hu, & Shulman, 2017), so far outcomes were mixed for non-model
82 organisms in general, and natural vertebrate populations in particular (Santure
83 & Garant, 2018). A key limitation is that sample sizes available for GWAS in
84 wildlife populations are typically several magnitudes smaller than for humans
85 or model organisms (Amos et al., 2011; Hong & Park, 2012). Furthermore, it
86 remains notoriously difficult to control for false positive signals when testing
87 for associations between genetic variants and phenotypes (Amos et al., 2011;
88 McCarthy et al., 2008), and distinguishing between molecular signals of
89 natural selection and genetic drift with F_{ST} outlier approaches is non-trivial
90 (Hoban et al., 2016; Vitti et al., 2013).

91 In order to increase statistical power to identify genomic regions
92 associated with phenotypic traits of interest, and at the same time control for
93 false positive signals, it has been suggested to combine approaches and to
94 apply population level replication (Chanock et al., 2007; Santure & Garant,

95 2018; Schielzeth, Rios, & Burri, 2018). Yet, to our knowledge this approach
96 (i.e. the combination of GWAS and F_{ST} -based tests and application of
97 population-level replication) has not been used to identify targets of pathogen-
98 mediated selection in wildlife disease systems to date.

99 *Borrelia afzelii* is the most common *Borrelia* genospecies in Europe
100 and one of the causative agents of human Lyme disease (Steere, Coburn, &
101 Glickstein, 2004). It is transmitted by ticks (*Ixodes* sp.) and rodents, such as
102 the bank vole (*Myodes glareouls*), are its main natural hosts (Kurtenbach et
103 al., 2006; Mannelli, Bertolotti, Gern, & Gray, 2012). Recently, it has been
104 experimentally shown that *B. afzelii* has negative fitness consequences for
105 bank voles (Cayol et al., 2018). Defense mechanisms that prevent or control
106 *Borrelia* infections in natural hosts will thus be favored by natural selection.

107 Using a candidate gene approach, we have previously demonstrated
108 that naturally occurring genetic variants of Toll-like receptor 2 (*TLR2*) are
109 associated with the *Borrelia* infection status of bank voles (Cornetti et al.,
110 2018; Tschirren et al., 2013). Yet, *Borrelia* susceptibility is most likely
111 influenced by many genes and the variation explained by a single candidate
112 gene remains limited (Wilfert & Schmid-Hempel, 2008). In this study, we
113 performed genome-wide scans to identify potential targets of *Borrelia*-
114 mediated selection using a combination of complementary approaches and by
115 applying a population replication criterion. Specifically, we used (1) a GWAS
116 approach to identify genetic variants associated with *Borrelia* infection status,
117 and (2) a F_{ST} -based analysis between *Borrelia*-infected and -uninfected bank
118 voles within seven independent populations to identify outlier SNPs. The latter
119 was combined with (3) a population level replication criterion in which we

120 considered only SNPs that were identified as outliers in multiple populations.
121 Finally, we (4) overlaid the results of these complementary approaches to
122 identify consensus candidate SNPs, which were found to explain significant
123 amount of variation in bank vole *Borrelia* infection status.

124

125 **Materials and Methods**

126 *Field sampling*

127 Bank voles (*Myodes glareolus*) were captured during summer 2014 at seven
128 locations in the Kanton Graubünden, Switzerland (Table 1; Supporting
129 information Figure S1) using live-traps (Longworth Mammal Traps, Anglian
130 Lepidopterist Supplies). A high *Borrelia* infection prevalence has previously
131 been documented in bank voles at these sampling sites (Cornetti et al., 2018).
132 Caught bank voles (N = 177; Table 1) were weighed (to the nearest g), aged
133 following Gliwicz (1988)(adults (>20 g), subadults (15-20 g), and juveniles (<
134 15 g)), and a small ear biopsy was collected and stored in 95% ethanol. The
135 animals were then released at their capture site. Vole capture, handling and
136 tissue sampling complied with the current laws of Switzerland and were
137 performed under a license issued by the Department of Food Safety and
138 Animal Health of the Kanton Graubünden, Chur, Switzerland (permit number
139 2012_17).

140

141 *Borrelia* infection in bank voles

142 Genomic DNA was extracted from the ear biopsy using the QIAGEN DNeasy
143 Blood & Tissue Kit (Qiagen, Venlo, the Netherlands). To determine the *B.*
144 *afzelii* infection status of bank voles, we used a highly sensitive quantitative

145 real-time PCR (qPCR) assay using the *flaB* *B. afzelii*-specific primers Fla5F:
146 5'-CACCCAGCATCACTTTCAGGA-3' and Fla6R: 5'-
147 CTCCTCACCAGCAAAAAGA-3' (Råberg, 2012) on a StepOnePlus real-time
148 qPCR machine (Applied Biosystems, Foster City, CA, USA). We focused on
149 *B. afzelii* because a pilot study using reverse line blot (Herrmann et al., 2013)
150 had revealed that *B. afzelii* is the only *Borrelia* genospecies present in bank
151 voles at our study sites (unpublished data).

152 The amplification was carried out in a total volume of 20 µl, including
153 10 µl SYBR® Select Master Mix (2x, Applied Biosystems), 0.8 µl of each
154 primer (10 µM) and 4 µl extracted genomic DNA. The qPCR protocol
155 consisted of two initial holding steps first at 50 °C and then at 95 °C for 2 min
156 each, followed by 42 cycles of 95 °C for 15 sec, 59 °C for 30 sec, and 72 °C
157 for 30 sec (Råberg, 2012). Eight negative controls and eight serially diluted
158 positive controls were included on each plate. Samples with a cycle threshold
159 (Ct) value > 0 and a melting temperature between 76.4 °C and 77.8 °C were
160 considered to be *B. afzelii*-positive (Råberg, 2012). All samples were analyzed
161 in duplicate on two different plates (see Cornetti et al., 2018 for details).

162

163 *Genotyping-by-Sequencing and SNPs calling*

164 Samples of 118 adult bank voles were used for Genotyping-by-Sequencing
165 (GBS). Only adult bank voles were included because variation in *Borrelia*
166 infection status among juveniles and subadults is likely due to variation in
167 exposure rather than resistance (Tschirren et al., 2013). An equal number of
168 *Borrelia*-infected and uninfected individuals were randomly selected for each

169 site, whenever possible. Overall, 45% of the sequenced individuals were *B.*
170 *afzelii* infected (Table 1).

171 Extracted genomic DNA was sent to the GBS platform
172 (<http://www.biotech.cornell.edu>) of Cornell University, USA in July 2015. GBS
173 libraries were prepared using a double digest protocol with *SbfI* and *HpaII* as
174 restriction enzymes (Poland, Brown, Sorrells, & Jannink, 2012). Sequencing
175 (100-bp single-end reads) was performed on IlluminaHiSeq 2500.

176 In total 264,483,546 reads were obtained. Illumina adapters were
177 removed from raw sequences using Trimmomatic 0.33 (Bolger, Lohse, &
178 Usadel, 2014) . Sequences were aligned to the prairie vole (*Microtus*
179 *ochrogaster*) reference genome (MicOch1.0, (McGraw, Davis, Young, &
180 Thomas, 2011)) using Bowtie2 (Langmead & Salzberg, 2012). The prairie
181 vole and the bank vole are members of the same subfamily (Arvicolinae) and
182 their divergence time has been estimated to be 5.9 ± 0.8 Mya (95% confidence
183 interval: 4.6-7.6 Mya) based on nuclear genes (Abramson, Lebedev, Tesakov,
184 & Bannikova, 2009). To date, the prairie vole is the closest relative of the bank
185 vole with a high quality genome assembly (McGraw et al., 2011). The current
186 version of the prairie vole reference genome consists of 28 main scaffolds,
187 corresponding to 17 autosomes, the X chromosome and ten linkage groups
188 (Zerbino et al., 2018). The average mapping rate to the prairie vole genome
189 was about 80%. Samtools 1.3 was used to filter the BAM alignments for
190 quality (-q 20) before SNP calling was performed with GATK 3.7 (Van der
191 Auwera et al., 2013). The final set of SNPs was obtained with VCFtools 0.1.15
192 (Danecek et al., 2011) by filtering for quality (minimum genotype quality score
193 of 20), coverage (minimum genotype depth of 6 per individual) and rare

194 variants (minor allele frequency of 0.01), requiring that at least 70% of all
195 individuals passed the filters.

196

197 *Bank vole population structure*

198 Population structuring was assessed using a Multi-Dimensional Scaling
199 (MDS) approach implemented in Plink 1.90 (Chang et al., 2014), as well as
200 using the software Structure 2.3.4 (Pritchard, Stephens, & Donnelly, 2000). As
201 input for Structure we used a reduced dataset of 1555 SNPs, which included
202 variants with no missing data, and, to fulfill Structure model assumptions of
203 independence of loci (i.e. no linkage disequilibrium within populations), only
204 one SNP per read. Analyses were performed using an admixture model with
205 correlated allele frequencies for ten independent runs. We determined the
206 most likely number of genetic clusters (K) exploring K values between one
207 and seven (i.e. the number of sampling sites). Burn-in periods of 100,000
208 were used, followed by 500,000 iterations. The most plausible number of
209 genetically well-defined groups was determined by comparing the likelihood at
210 different K values (Pritchard et al., 2000) using Structure Harvester (Earl &
211 VonHoldt, 2012).

212 Furthermore, we calculated the fixation index F_{ST} among populations
213 as a measure of population differentiation using the software Arlequin version
214 3.5 (Excoffier & Lischer, 2010) and tested for isolation-by-distance (IBD), by
215 correlating the pairwise genetic differentiation and geographic distance among
216 populations using Mantel test (Mantel, 1967). Linear geographic distances
217 between locations were calculated with the Geographic Distance Matrix
218 Generator version 1.2.3 (Ersts, 2017). The relationship between the linearized

219 F_{ST} ($F_{ST} / (1 - F_{ST})$) and the log-transformed linear geographic distance
220 (log(km)) was estimated using the Isolation-by-Distance Web Server (Jensen,
221 Bohonak, & Kelley, 2005).

222

223 *Identifying targets of Borrelia-mediated selection*

224 To identify putative targets of *Borrelia*-mediated selection we used a
225 combination of two approaches. First, we tested for an association between
226 bank vole SNPs and *Borrelia* infection status using the R package GenABEL
227 (Aulchenko, Ripke, Isaacs, & van Duijn, 2007). GenABEL allows performing
228 genome-wide association (GWAS) between SNPs and a phenotype while
229 correcting for population structure. We first computed a kinship matrix for our
230 samples using the function *ibs*. Then, using an Eigenstrat method, we
231 calculated the probability of each SNP to be associated with the phenotype
232 (i.e. *Borrelia* infection status), after correcting for kinship within the whole
233 dataset (Aulchenko et al., 2007) using the function *egscore*. This method uses
234 the genomic kinship matrix to derive axes of genetic variation (principal
235 components) and adjusts both the trait (i.e. *Borrelia* infection status) and the
236 genotypes onto these axes (Price et al., 2006). Corrected genotypes are
237 defined as residuals from regression of genotypes onto axes. Correlation
238 between corrected genotypes and the phenotype is computed, and test
239 statistics is defined as the square of this correlation times (N - K - 1), where N
240 is the number of genotyped subjects and K is the number of axes (Aulchenko
241 et al., 2007). The association analysis was performed using the combined
242 dataset (i.e. all individuals, N = 118, across all populations and all 21,811
243 SNPs).

244 Second, we complemented the GWAS approach with an independent
245 analysis based on F_{ST} to identify outlier SNPs between *Borrelia*-infected and
246 uninfected bank voles within populations. In each of the seven populations,
247 we calculated for each SNP the F_{ST} value between *Borrelia*-infected and
248 uninfected individuals using VCFtools 0.1.15 (Danecek et al., 2011). The
249 rationale of this approach is that no neutral population differentiation is
250 expected between infected and uninfected animals within a population, and
251 significantly differentiated SNPs suggest an association with *Borrelia*
252 resistance or susceptibility. Within each population, we selected the outlier
253 SNPs that lay within the top 10% of the population-specific F_{ST} distribution
254 (Bankers et al., 2017; Myles, Davison, Barrett, Stoneking, & Timpson, 2008;
255 Zueva et al., 2014). Given the significant population differentiation and distinct
256 genetic clustering of the sampled bank vole populations (see Results), these
257 populations represent independent replicates.

258 To control for false positive F_{ST} outliers within populations, we applied a
259 population replication criterion and only considered F_{ST} outliers that were
260 among the top 10% of the population-specific F_{ST} distribution in at least three
261 of the seven populations (40%) (see Simulations below). This approach is
262 conservative because it assumes that the same genetic mechanisms underlie
263 variation in *Borrelia* resistance or susceptibility in different populations.
264 Because genetic variants underlying variation in resistance or susceptibility to
265 infectious diseases are typically found to be exonic (Hill, 2012), we specifically
266 focused on SNPs located in exons in both approaches.

267 In a final step, we overlapped the results of the GWAS approach and of
268 the population-level replicated F_{ST} -based approach to identify consensus

269 candidate SNPs that were associated with *Borrelia* infection status in both
270 analyses. Candidate SNPs were annotated with SNPdat (Doran & Creevey,
271 2013), a tool specifically designed for non-model organisms. We then used a
272 generalized linear mixed model with a binomial error structure and site
273 included as a random effect to test for associations between these consensus
274 candidate SNPs and *Borrelia* infection status. Significance of fixed effects was
275 determined by comparing nested models (with and without the variable of
276 interest) using likelihood ratio tests. We considered both additive and
277 dominant modes of gene action of candidate SNPs. Analyses were performed
278 using the package lme4 (Bates, Maechler, & Bolker, 2011) in R version 3.6.1.
279 (R Core Team, 2014).

280

281 *Neutral simulations*

282 We performed neutral simulations to quantify the false-positive rates
283 when using a GWAS approach alone, a F_{ST} approach without population
284 replication, a F_{ST} approach with a two population replication criterion and a
285 F_{ST} approach with a three population replication criterion, as well as a
286 combined approach that focuses on consensus SNPs identified with both
287 approaches to identify putative targets of *Borrelia*-mediated selection. We
288 used the forward-time genetic simulator SLiM version 3.3 (Haller & Messer,
289 2019) to generate neutral polymorphisms based on the number of SNPs ($N =$
290 21'811) and individuals ($N = 118$) included in the empirical dataset, repeated
291 100 times. The commented SLiM script describing the simulation process is
292 presented in the Supporting information. In short, we generated sequences of
293 27'999 bp in length, with a mutation rate of 2×10^{-4} per base per generation

294 and a recombination rate of $r = 0.05$. At generation 1, seven subpopulations
295 appear and evolve independently for 4999 generations with some gene flow
296 among them. At the end of the simulations (i.e. after 4999 generations) a
297 subsample of each population is randomly selected according to the actual
298 sample size of the empirical data. The data is written in seven VCF files that
299 are used for further analyses.

300 The parameters used in the simulations were selected, after many pilot
301 runs, by taking into account computational and temporal constraints and in
302 order to obtain a total number of SNPs and population differentiation (in
303 particular in term of F_{ST} values) similar to the ones observed in the real
304 dataset. For each of the 100 simulated datasets, we randomly assigned
305 infected and uninfected individuals within each population, according to the
306 real data (for example, in Sagong, 10 infected and 9 uninfected individuals
307 were defined). Then, the seven VCF files were merged using GATK version 4
308 (McKenna et al., 2010); during this step, we also constrained the total number
309 of SNPs and the amount of missing data to be comparable to the observed
310 data. The resulting 100 VCF files, including 118 samples from seven
311 populations, were used to perform GWAS with GenABEL (Aulchenko, Ripke,
312 Isaacs, & van Duijn, 2007) based on the full dataset (118 samples), and F_{ST} -
313 based calculation between infected and uninfected individuals within each of
314 the populations.

315

316 **Results**

317 *Bank vole population structure*

318 A total of 21,811 SNPs were retained after quality filtering. Population
319 differentiation (F_{ST}) was comparably high considering the relatively small size
320 of the study area (maximum linear distance between sampling sites: 34 km),
321 and varied from 0.052 (Sagogn-Flims, the second closest locations) to 0.111
322 (Malans-Sagogn, the most distant locations, Supporting information Table
323 S1). All F_{ST} values were statistically significant (Supporting information Table
324 S1).

325 Within the study area, population structure was well defined. The MDS
326 analysis highlighted seven distinct groups corresponding to the seven
327 sampling locations (Figure 1). Similarly, using Structure we found that $K = 7$
328 was the most supported partition (Supporting information Figure S2), with the
329 seven well defined genetic clusters corresponding to the seven sampling
330 locations (Figure 1, Supporting information Figure S3). The relationship
331 between genetic and geographic distance was positive ($r = 0.87$, Supporting
332 information Figure S4) and statistically significant (Mantel test, $p < 0.001$),
333 suggesting pronounced isolation-by-distance across populations.

334

335 *Neutral simulations*

336 The neutral simulations demonstrated the weaknesses of using GWAS
337 and F_{ST} -based tests in isolation (Supporting information Figures S5 and S6)
338 and the significant reduction in false-positive rates when using a three-
339 population replication criterion for the F_{ST} -based test (Supporting information
340 Figure S5). The simulations furthermore showed that combining GWAS and
341 F_{ST} -based tests with population replication substantially reduces false-

342 positive rates, and thus increases the power to identify core candidate SNPs
343 (Figure 2).

344

345 *Identifying targets of Borrelia-mediated selection*

346 *a) GWAS approach*

347 Using a GWAS approach that corrects for population structure, we identified
348 1065 SNPs that were associated ($p < 0.05$) with the *Borrelia* infection status of
349 bank voles (Supporting information Figure S6). As expected given the
350 comparably small number of individuals included in this GWAS, none of these
351 associations reached statistical significance when accounting for multiple
352 testing using the Benjamini-Hochberg procedure (Supporting information
353 Figure S6). The identified SNPs were distributed across the whole genome
354 (Supporting information Figure S7). After correcting for chromosome size,
355 chromosome 15 and chromosome 19 showed the highest and lowest number
356 of putatively *Borrelia*-associated SNPs, respectively (chr 15: 1.02 SNPs/MB;
357 chr 19: 0.17 SNPs/MB; Supporting information Figure S7). Among the 1065
358 SNPs, 53 were located in exonic regions of 48 unique genes (Supporting
359 information Table S2).

360

361 *b) F_{ST} -based approach with population-level replication*

362 In a second step, we performed a F_{ST} -based analysis between *Borrelia*
363 infected and non-infected individuals within each of the seven populations and
364 applied a population-level replication criterion, by considering only F_{ST} outliers
365 that were among the top 10% of the population-specific F_{ST} distribution in at
366 least three of the seven population. The top 10% F_{ST} threshold ranged from

367 0.079 (Flims) to 0.145 (Rodels, Supporting information Figure S8). We
368 obtained 305 SNPs that showed consistent differentiation between *Borrelia*-
369 infected and uninfected bank voles in at least three independent population
370 replicates, 70 of which were also identified using the GWAS approach (Figure
371 2). Among the 305 SNPs, nine were located in exonic regions of nine unique
372 genes (Table 2).

373

374 *c) Consensus candidate loci*

375 When overlapping the results of the GWAS and the F_{ST} -based test, there
376 were four SNPs in coding regions that were associated with *Borrelia* infection
377 status in both analyses (Table 2). These four SNPs were located in the exonic
378 regions of the genes *Slc26a4* (Solute carrier family 26, member 4), *Tns3*
379 (Tensin 3), *Wscd1* (WSC domain containing 1) and *Espnl* (Espin-like). The
380 less frequent allele of *Slc26a4* (allele A; $\chi^2 = 4.414$, $DF = 1$, $P = 0.036$) and
381 *Tns3* (allele A; GLMM: $\chi^2 = 6.859$, $DF = 1$, $P = 0.009$) were associated with a
382 lower probability of *Borrelia* infection (Figure 3, Supporting information Table
383 S3), whereas the less frequent allele of *Wscd1* (allele G; $\chi^2 = 5.790$, $DF = 1$, P
384 = 0.016) and *Espnl* (allele C; $\chi^2 = 6.488$, $DF = 1$, $P = 0.011$) were associated
385 with a higher probability of *Borrelia* infection (Figure 3, Supporting information
386 Table S3). Models that treated the heterozygous and homozygous state of the
387 less frequent allele as separate genotypes are presented here. Models that
388 combined the heterozygous and homozygous state of the less frequent allele
389 are presented in Supporting information Table S4.

390

391 **Discussion**

392 For most non-model organisms, and wild-living vertebrates in particular, the
393 genetic basis underlying infectious disease resistance or susceptibility is
394 poorly understood (Spurgin & Richardson, 2010), which hampers progress in
395 our understanding of the eco-evolutionary dynamics of host-parasite
396 interactions and wildlife disease. Here we present a case study in a natural
397 rodent-*Borrelia* system in which we combined a GWAS and a F_{ST} -based
398 approach with population-level replication to identify genomic regions
399 associated with variation in *Borrelia* infection status.

400 *Borrelia* prevalence in bank voles was high at all seven study sites,
401 with 31-61% of adult bank voles being *Borrelia* infected across sites. Given
402 the negative fitness consequences a *Borrelia* infection has for the rodent host
403 (Cayol et al., 2018), natural selection is expected to favour mechanism that
404 control or prevent infection. Using a GWAS approach we identified a large
405 number (1065) of SNPs that were associated ($p < 0.05$) with the voles'
406 *Borrelia* infection status while controlling for population structure. Given the
407 relatively relaxed criteria that were applied to identify genotype-phenotype
408 associations in this GWAS approach, many of these SNPs were likely false
409 positives. Indeed, the risk of false positive signals in GWAS increases when a
410 large number of SNPs but a small number of individuals are included (Hong &
411 Park, 2012), which was the case in our study, and is typical for most studies
412 on wildlife populations.

413 In order to increase statistical power to detect true signals while
414 controlling for false positives, we complemented the GWAS approach with an
415 F_{ST} -based analysis, replicated across populations. We observed pronounced
416 (given the relatively small geographical distances among study sites) genetic

417 differentiation of bank voles across study sites and strong isolation-by-
418 distance. Furthermore, Bayesian assignment analysis identified seven well
419 defined genetic clusters, corresponding to the seven sampling sites. These
420 results indicate that the seven bank vole populations are sufficiently isolated
421 to represent independent replicates for the identification of targets of *Borrelia*-
422 mediated selection. Within each population we identified SNPs lying within the
423 top 10% of the population-specific F_{ST} distribution and considered SNPs that
424 were among these extreme 10% in multiple populations to be putatively
425 associated with *Borrelia* infection status. This F_{ST} -based approach assumes
426 that the same mechanisms underlie *Borrelia* resistance or susceptibility in
427 multiple populations, which does not necessarily need to be the case. Indeed,
428 quantitative trait loci (QTL) are often found to be population-specific
429 (Schielzeth et al., 2018; Tschirren & Bensch, 2010), and previous work in
430 other systems has demonstrated variation in host responses to pathogen
431 infection across populations (Bankers et al., 2017; Kurtz et al., 2014).
432 Applying a population-level replication criterion is conservative and will
433 prevent identifying such population-specific SNPs associated with *Borrelia*
434 infection status, yet it also allows us to control for false positive signals.
435 Indeed, the neutral simulations demonstrated the power of using a population
436 replication criterion to reduce false positive signals.

437 The simulations furthermore showed that false-positive rates are
438 further reduced when GWAS and F_{ST} -based approaches are combined. Using
439 such a combined approach, we identified four consensus polymorphisms
440 located in exonic regions representing promising targets of *Borrelia*-mediated
441 selection in bank voles. Interestingly, for most of the genes in which the

442 consensus SNPs were located, a functional association with response to
443 bacterial infection has previously been demonstrated in rodents: *Slc26a4*
444 encodes the anion exchanger Pendrin and is expressed in membranes of ion-
445 transporting epithelia where it regulates luminal pH and fluid transport
446 (Royaux et al., 2000). *Slc26a4* has been found to be significantly upregulated
447 in mouse macrophages experimentally stimulated with live *B. burgdorferi*
448 (Gautam et al., 2011). Similarly, mice infected with the bacterial pathogen
449 *Bordetella pertussis*, the etiological agent of whooping cough, exhibited
450 significant *Slc26a4* upregulation (Scanlon et al., 2014). *Tns3* encodes a
451 phosphoprotein thought to act as a link between the extracellular matrix and
452 the cytoskeleton (Lo, 2004). This gene was found to be significantly
453 upregulated in mice experimentally infected with the bacterium
454 *Mycobacterium bovis*, the main etiological agent of bovine tuberculosis
455 (Aranday-Cortes et al., 2012). Another member of the tensin gene family,
456 *Tensin 1*, was downregulated and upregulated 4 h and 24 h respectively, after
457 stimulation of mouse macrophages with live *B. burgdorferi* (Gautam et al.,
458 2011), suggesting a general role of the tensin gene family in the response to
459 bacterial infections. *Wscd1* encodes a protein with sulfotransferase activity
460 (Smith, Blake, Kadin, Richardson, & Bult, 2018). *Wscd1* has been found to be
461 significantly downregulated in mice five days after infection with the bacterium
462 *Yersinia pseudotuberculosis* (Heine et al., 2018). Taken together, these
463 previous findings suggest that several genes identified as candidates in our
464 study play a role in the response to bacterial infections in rodents. However,
465 the exact mechanisms by which these genes may confer resistance or

466 susceptibility to *Borrelia* in bank voles are currently unknown and will be the
467 focus of future work.

468 Previously, we have demonstrated that naturally occurring
469 polymorphisms at the innate immune receptor Toll-like receptor 2 (*TLR2*) are
470 significantly associated with the *Borrelia* infection status of bank voles in
471 multiple populations, including the current study population (Cornetti et al.,
472 2018; Tschirren et al., 2013; Tschirren, 2015). These findings were in line with
473 biomedical research that has identified *TLR2* as a candidate gene for *Borrelia*
474 resistance in laboratory mice (Alexopoulou et al., 2002; Dennis et al., 2009;
475 Singh & Girschick, 2006; Wooten et al., 2002). Interestingly, however, *TLR2*
476 was not identified as a potential candidate gene in this study, also not when
477 considering the results of GWAS and F_{ST} -based approaches separately. It
478 demonstrates the limitations of reduced-representation sequencing
479 approaches, which do not cover the whole genome (Davey & Blaxter, 2010).
480 Depending on the density and distribution of SNPs, as well as recombination
481 rates in regions of interest, the power to detect signals might be low. In our
482 study, no GBS SNP was close enough to possibly pick up signals of selection
483 acting on *TLR2*. In fact the closest SNP was about 46kb away from *TLR2*, a
484 physical distance larger than the linkage disequilibrium decay estimated for
485 natural rodent populations ($r^2 < 0.2$ ~20 kb; Staubach et al, 2012). This
486 suggests that some other putative candidate regions were likely missed with
487 our approach because of insufficient SNP coverage. At the same time, we can
488 exclude the possibility that the SNPs identified in our study were physically
489 linked to *TLR2*. *TLR2* is located on chromosome 1 in the prairie vole, whereas
490 *Tns3*, *Wscd1* and *Espnl* are located on chromosome 7, linkage group 1 and

491 linkage group 4, respectively. *Slc26a4* is located on chromosome 1 as well,
492 but separated by more than 40 Mb from *TLR2*. One possible way to obtain a
493 more conclusive list of candidate genes associated with *Borrelia* infectious
494 status would be to increase SNP density or sequence the whole genome.
495 However, costs associated with especially the latter approach are still
496 excruciatingly high, in particular when large numbers of individuals are
497 included and the study species has a large genome size, as is the case for
498 mammals (Catchen et al., 2017).

499 In conclusion, by combining GWAS and F_{ST} -based approach with
500 population-level replication we identified consensus SNPs in exonic regions of
501 genes for which a functional association with host responses to bacterial
502 infections has previously been demonstrated. These loci thus represent
503 promising new candidate genes that may allow tracking evolutionary changes
504 in host populations in response to *Borrelia* emergence. More generally, the
505 combined approach used in this case study can be applied to other systems
506 and may contribute to a better understanding of genotype-phenotype
507 associations in wildlife populations.

508

509 **Acknowledgements**

510 This study was supported by the University of Zurich Research Priority
511 Program “Evolution in Action: from Genomes to Ecosystems”, the Faculty of
512 Science of the University of Zurich, the Baugarten Stiftung and the Stiftung für
513 wissenschaftliche Forschung an der Universität Zürich. We thank the
514 numerous people who contributed to sample collection in the field, Peter

515 Fields for suggestions on genetic simulations, and Jacek Radwan and two
516 anonymous reviewers for their constructive comments on the manuscript.

517

518

519 **Author Contributions**

520 LC and BT designed the research, LC performed laboratory work and

521 analysed the data, LC and BT wrote the paper.

522

523 **Data accessibility**

524 The Genotyping-by-Sequencing data are deposited in NCBI BioProject

525 PRJNA306409, SRA experiment SRR3031372. The quality filtered SNP file

526 used for population genomic analyses and information on population and

527 infection status of individual bank voles are deposited in the Dryad repository

528 doi:10.5061/dryad.c866t1g3t.

529 **References**

- 530 Abramson, N. I., Lebedev, V. S., Tesakov, A. S., & Bannikova, A. A. (2009).
531 Supraspecies relationships in the subfamily Arvicolinae (Rodentia,
532 Cricetidae): An unexpected result of nuclear gene analysis. *Molecular*
533 *Biology*, 43(5), 834–846. doi:10.1134/S0026893309050148
- 534 Alexopoulou, L., Thomas, V., Schnare, M., Lobet, Y., Anguita, J., Schoen, R.
535 T., ... Flavell, R. A. (2002). Hyporesponsiveness to vaccination with
536 *Borrelia burgdorferi* OspA in humans and in *TLR1*- and *TLR2*-deficient
537 mice. *Nature Medicine*, 8(8), 878–84. doi:10.1038/nm732
- 538 Amos, W., Driscoll, E., & Hoffman, J. I. (2011). Candidate genes versus
539 genome-wide associations: Which are better for detecting genetic
540 susceptibility to infectious disease? *Proceedings of the Royal Society B:*
541 *Biological Sciences*, 278(1709), 1183–1188. doi:10.1098/rspb.2010.1920
- 542 Aranday-Cortes, E., Hogarth, P. J., Kaveh, D. A., Whelan, A. O., Villarreal-
543 Ramos, B., Lalvani, A., & Vordermeier, H. M. (2012). Transcriptional
544 profiling of disease-induced host responses in bovine tuberculosis and
545 the identification of potential diagnostic biomarkers. *PLoS ONE*, 7(2).
546 doi:10.1371/journal.pone.0030626
- 547 Aulchenko, Y. S., Ripke, S., Isaacs, A., & van Duijn, C. M. (2007). GenABEL:
548 An R library for genome-wide association analysis. *Bioinformatics*,
549 23(10), 1294–1296. doi:10.1093/bioinformatics/btm108
- 550 Bankers, L., Fields, P., McElroy, K. E., Boore, J. L., Logsdon, J. M., &
551 Neiman, M. (2017). Genomic evidence for population-specific responses
552 to co-evolving parasites in a New Zealand freshwater snail. *Molecular*
553 *Ecology*, 26(14), 3663–3675. doi:10.1111/mec.14146

554 Bates, D., Maechler, M., Bolker, B. & Walker, S. (2011). Fitting linear mixed-
555 effects models using lme4. *Journal of Statistical Software*, 67(1), 1-48.

556 Beldomenico, P. M., & Begon, M. (2010). Disease spread, susceptibility and
557 infection intensity: vicious circles? *Trends in Ecology and Evolution*,
558 25(1), 21–27. doi:10.1016/j.tree.2009.06.015

559 Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible
560 trimmer for Illumina sequence data. *Bioinformatics*, 30(15), 2114–2120.
561 doi:10.1093/bioinformatics/btu170

562 Catchen, J. M., Hohenlohe, P. A., Bernatchez, L., Funk, W. C., Andrews, K.
563 R., & Allendorf, F. W. (2017). Unbroken : RADseq remains a powerful tool
564 for understanding the genetics of adaptation in natural populations.
565 *Molecular Ecology Resources*, 17, 362–365. doi:10.1111/1755-
566 0998.12669

567 Cayol, C., Giermek, A., Gomez-Chamorro, A., Hytönen, J., Kallio, E. R.,
568 Mappes, T., ... Koskela, E. (2018). *Borrelia afzelii* alters reproductive
569 success in a rodent host. *Proceedings of the Royal Society B: Biological*
570 *Sciences*, 285(1884), 20181056. doi:10.1098/rspb.2018.1056

571 Chang, C. C., Chow, C. C., Tellier, L. C. A. M., Vattikuti, S., Purcell, S. M., &
572 Lee, J. J. (2014). Second-generation PLINK: rising to the challenge of
573 larger and richer datasets, 1–16. doi:10.1186/s13742-015-0047-8

574 Chanock, S. J., Manolio, T., Boehnke, M., Boerwinkle, E., Hunter, D. J.,
575 Thomas, G., ... Collins, F. S. (2007). Replicating genotype–phenotype
576 associations. *Nature*, 447, 655-660. doi: 10.1038/447655a

577 Cornetti, L., Hilfiker, D., Lemoine, M., & Tschirren, B. (2018). Small-scale
578 spatial variation in infection risk shapes the evolution of a *Borrelia*

579 resistance gene in wild rodents. *Molecular Ecology*, 27(17), 3515–3524.
580 doi:10.1111/mec.14812

581 Cornetti, L. & Tschirren, B. (2020) Data from: Combining GWAS and F_{ST} -
582 based approaches to identify targets of *Borrelia*-mediated selection in
583 natural rodent hosts; Dryad; doi: 10.5061/dryad.c866t1g3t.

584 Cornetti, L. & Tschirren, B. (2020). NCBI Sequence Read Archive (BioProject
585 ID: PRJNA306409, SRA experiment SRR3031372).
586 <https://www.ncbi.nlm.nih.gov/bioproject/PRJNA306409>.

587 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M.
588 A., ... Durbin, R. (2011). The variant call format and VCFtools.
589 *Bioinformatics*, 27(15), 2156–2158. doi:10.1093/bioinformatics/btr330

590 Daszak, P., Cunningham, A. A., & Hyatt, A. D. (2000). Emerging infectious
591 diseases of wildlife - threats to biodiversity and human health. *Science*,
592 287, 443–449. doi:10.1126/science.287.5452.443

593 Davey, J. L., & Blaxter, M. W. (2010). RADseq: Next-generation population
594 genetics. *Briefings in Functional Genomics*, 9(5–6), 416–423.
595 doi:10.1093/bfgp/elq031

596 Dennis, V. A., Dixit, S., O'Brien, S. M., Alvarez, X., Pahar, B., & Philipp,
597 M. T. (2009). Live *Borrelia burgdorferi* spirochetes elicit inflammatory
598 mediators from human monocytes via the toll-like receptor signaling
599 pathway. *Infection and Immunity*, 77(3), 1238–1245.
600 doi:10.1128/IAI.01078-08

601 Doran, A. G., & Creevey, C. J. (2013). Snpdat: easy and rapid annotation of
602 results from de novo snp discovery projects for model and non-model
603 organisms. *BMC Bioinformatics*, 14(1), 45. doi:10.1186/1471-2105-14-45

604 Earl, D. A., & VonHoldt, B. M. (2012). STRUCTURE HARVESTER: A website
605 and program for visualizing STRUCTURE output and implementing the
606 Evanno method. *Conservation Genetics Resources*, 4(2), 359–361.
607 doi:10.1007/s12686-011-9548-7

608 Ersts, P. J. (2017). Geographic Distance Matrix Generator (version 1.2.3).
609 Retrieved July 18, 2017, from
610 http://biodiversityinformatics.amnh.org/open_source/gdmg

611 Excoffier, L., & Lischer, H. E. L. (2010). Arlequin suite ver 3.5: A new series of
612 programs to perform population genetics analyses under Linux and
613 Windows. *Molecular Ecology Resources*, 10(3), 564–567.
614 doi:10.1111/j.1755-0998.2010.02847.x

615 Flint, J., & Eskin, E. (2012). Genome-wide association studies in mice. *Nature*
616 *Reviews. Genetics*, 13(11), 807–817. doi:10.1038/nrg3335

617 Gautam, A., Dixit, S., Philipp, M. T., Singh, S. R., Morici, L. A., Kaushal, D., &
618 Dennis, V. A. (2011). Interleukin-10 alters effector functions of multiple
619 genes induced by *Borrelia burgdorferi* in macrophages to regulate Lyme
620 disease inflammation. *Infection and Immunity*, 79(12), 4876–4892.
621 doi:10.1128/IAI.05451-11

622 Gliwicz, J. (1988). Seasonal dispersal in non-cyclic populations of
623 *Clethrionomys glareolus* and *Apodemus flavicollis*. *Acta Theriologica*, 33,
624 263–272. doi:10.4098/AT.arch.88-20

625 Haller, B. C., & Messer, P. W. (2019). SLiM 3: Forward genetic simulations
626 beyond the Wright-Fisher model. *Molecular Biology and Evolution*, 36(3),
627 632–637. doi:10.1093/molbev/msy228

628 Heine, W., Beckstette, M., Heroven, A. K., Thiemann, S., Heise, U., Nuss, A.

629 M., ... Dersch, P. (2018). Loss of CNFYtoxin-induced inflammation drives
630 *Yersinia pseudotuberculosis* into persistency. *PLoS Pathogens* 14(2),
631 e1006858. doi:10.1371/journal.ppat.1006858

632 Hill, A. V. S. (2012). Evolution, revolution and heresy in the genetics of
633 infectious disease susceptibility. *Philosophical Transactions of the Royal*
634 *Society B: Biological Sciences*, 367(1590), 840–849.
635 doi:10.1098/rstb.2011.0275

636 Hoban, S., Kelley, J. L., Lotterhos, K. E., Antolin, M. F., Bradburd, G., Lowry,
637 D. B., ... Whitlock, M. C. (2016). Finding the genomic basis of local
638 adaptation: pitfalls, practical solutions, and future directions. *American*
639 *Naturalist*, 188(4), 379–397. doi:10.1086/688018

640 Hong, E. P., & Park, J. W. (2012). Sample size and statistical power
641 calculation in genetic association studies. *Genomics & Informatics*, 10(2),
642 117–122. doi:10.5808/GI.2012.10.2.117

643 Jensen, J. L., Bohonak, A. J., & Kelley, S. T. (2005). Isolation by distance,
644 web service. *BMC Genetics*, 6, 13. doi:10.1186/1471-2156-6-13

645 Jones, K. E., Patel, N. G., Levy, M. A., Storeygard, A., Balk, D., Gittleman, J.
646 L., & Daszak, P. (2008). Global trends in emerging infectious diseases.
647 *Nature*, 451(7181), 990–993. doi:10.1038/nature06536

648 Jones, M. E., Jarman, P. J., Lees, C. M., Hesterman, H., Hamede, R. K.,
649 Mooney, N. J., ... McCallum, H. (2007). Conservation management of
650 Tasmanian devils in the context of an emerging, extinction-threatening
651 disease: Devil facial tumor disease. *EcoHealth*, 4(3), 326–337.
652 doi:10.1007/s10393-007-0120-6

653 Kurtenbach, K., Hanincová, K., Tsao, J. I., Margos, G., Fish, D., & Ogden, N.

654 H. (2006). Fundamental processes in the evolutionary ecology of Lyme
655 borreliosis. *Nature Reviews Microbiology*, 4(9), 660–669.
656 doi:10.1038/nrmicro1475

657 Kurtz, J., Behrens, S., Schulenburg, H., Bornberg-Bauer, E., Peuß, R.,
658 Milutinović, B., ... Esser, D. (2014). Infection routes matter in population-
659 specific responses of the red flour beetle to the entomopathogen *Bacillus*
660 *thuringiensis*. *BMC Genomics*, 15(1), 445. doi:10.1186/1471-2164-15-445

661 Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with
662 Bowtie 2. *Nature Methods*, 9(4), 357–9. doi:10.1038/nmeth.1923

663 Lazzaro, B. P., & Little, T. J. (2009). Immunity in a variable world.
664 *Philosophical Transactions of the Royal Society B: Biological Sciences*,
665 364(1513), 15–26. doi:10.1098/rstb.2008.0141

666 Lo, S. H. (2004). Tensin. *International Journal of Biochemistry and Cell*
667 *Biology*, 36(1), 31–34. doi:10.1016/S1357-2725(03)00171-7

668 MacPhee, R. D. E., & Greenwood, A. D. (2013). Infectious disease,
669 endangerment, and extinction. *International Journal of Evolutionary*
670 *Biology*, 2013, 1–9. doi:10.1155/2013/571939

671 Mannelli, A., Bertolotti, L., Gern, L., & Gray, J. (2012). Ecology of *Borrelia*
672 *burgdorferi* sensu lato in Europe: transmission dynamics in multi-host
673 systems, influence of molecular processes and effects of climate change.
674 *FEMS Microbiology Reviews*, 36(4), 837–61. doi:10.1111/j.1574-
675 6976.2011.00312.x

676 Mantel, N. (1967). The detection of disease clustering and a generalized
677 regression approach. *Cancer Research*, 27, 209-220.

678 Margres, M. J., Jones, M., Epstein, B., Kerlin, D. H., Comte, S., Fox, S., ...

679 Storfer, A. (2018). Large-effect loci affect survival in Tasmanian devils (
680 *Sarcophilus harrisii*) infected with a transmissible cancer. *Molecular*
681 *Ecology*, 27(21), 4189–4199. doi:10.1111/mec.14853

682 McCarthy, M. I., Abecasis, G. R., Cardon, L. R., Goldstein, D. B., Little, J.,
683 Ioannidis, J. P. a, & Hirschhorn, J. N. (2008). Genome-wide association
684 studies for complex traits: consensus, uncertainty and challenges. *Nature*
685 *Reviews. Genetics*, 9(5), 356–369. doi:10.1038/nrg2344

686 McGraw, L. A., Davis, J. K., Young, L. J., & Thomas, J. W. (2011). A genetic
687 linkage map and comparative mapping of the prairie vole (*Microtus*
688 *ochrogaster*) genome. *BMC Genetics*, 12(1), 60. doi:10.1186/1471-2156-
689 12-60

690 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky,
691 A., ... DePristo, M. A. (2010). The genome analysis toolkit: A MapReduce
692 framework for analyzing next-generation DNA sequencing data. *Genome*
693 *Research*, 20(9), 1297–1303. doi:10.1101/gr.107524.110

694 Myles, S., Davison, D., Barrett, J., Stoneking, M., & Timpson, N. (2008).
695 Worldwide population differentiation at disease-associated SNPs. *BMC*
696 *Medical Genomics*, 1(1), 22. doi:10.1186/1755-8794-1-22

697 Petersen, B. S., Fredrich, B., Hoepfner, M. P., Ellinghaus, D., & Franke, A.
698 (2017). Opportunities and challenges of whole-genome and -exome
699 sequencing. *BMC Genetics*, 18(1), 1–13. doi:10.1186/s12863-017-0479-5

700 Poland, J. A., Brown, P. J., Sorrells, M. E., & Jannink, J. L. (2012).
701 Development of high-density genetic maps for barley and wheat using a
702 novel two-enzyme genotyping-by-sequencing approach. *PLoS ONE*, 7(2).
703 doi:10.1371/journal.pone.0032253

704 Price, A. L., Patterson, N. J., Plenge, R. M., Weinblatt, M. E., Shadick, N. A.,
705 & Reich, D. (2006). Principal components analysis corrects for
706 stratification in genome-wide association studies. *Nature Genetics*, 38(8),
707 904–909. doi:10.1038/ng1847

708 Price, A. L., Spencer, C. C. A., & Donnelly, P. (2015). Progress and promise
709 in understanding the genetic basis of common diseases. *Proceedings of*
710 *the Royal Society B: Biological Sciences*, 282(1821), 20151684.
711 doi:10.1098/rspb.2015.1684

712 Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of population
713 structure using multilocus genotype data. *Genetics*, 155(2), 945–959.
714 doi:10.1111/j.1471-8286.2007.01758.x

715 R Core Team (2014). R: A language and environment for statistical
716 computing. R Foundation for Statistical Computing, Vienna, Austria.

717 Råberg, L. (2012). Infection intensity and infectivity of the tick-borne pathogen
718 *Borrelia afzelii*. *Journal of Evolutionary Biology*, 25(7), 1448–53.
719 doi:10.1111/j.1420-9101.2012.02515.x

720 Royaux, I. E., Suzuki, K., Mori, A., Katoh, R., Everett, L. A., Kohn, L. D., &
721 Green, E. D. (2000). Pendrin, the protein encoded by the pendred
722 syndrome gene (PDS), is an apical porter of iodide in the thyroid and is
723 regulated by thyroglobulin in FRTL-5 cells. *Endocrinology*, 141(2), 839–
724 845. doi:10.1210/endo.141.2.7303

725 Santure, A. W., & Garant, D. (2018). Wild GWAS - association mapping in
726 natural populations. *Molecular Ecology Resources*, 18(4), 729–738.
727 doi:10.1111/1755-0998.12901

728 Scanlon, K. M., Gau, Y., Zhu, J., Skerry, C., Wall, S. M., Soleimani, M., &

729 Carbonetti, N. H. (2014). Epithelial anion transporter Pendrin contributes
730 to inflammatory lung pathology in mouse models of *Bordetella pertussis*
731 infection. *Infection and Immunity*, 82(10), 4212–4221.
732 doi:10.1128/iai.02222-14

733 Schielzeth, H., Rios, A., & Burri, R. (2018). Success and failure in replication
734 of genotype-phenotype associations: How does replication help in
735 understanding the genetic basis of phenotypic variation in outbred
736 populations? *Molecular Ecology Resources*, 18(4), 739–754.
737 doi:10.1111/1755-0998.12780

738 Singh, S. K., & Girschick, H. J. (2006). Toll-like receptors in *Borrelia*
739 *burgdorferi*-induced inflammation. *Clinical Microbiology and Infection*,
740 12(8), 705–717. doi:10.1111/j.1469-0691.2006.01440.x

741 Smith, C. L., Blake, J. A., Kadin, J. A., Richardson, J. E., & Bult, C. J. (2018).
742 Mouse Genome Database (MGD)-2018: Knowledgebase for the
743 laboratory mouse. *Nucleic Acids Research*, 46(D1), D836–D842.
744 doi:10.1093/nar/gkx1006

745 Smith, K. F., Sax, D. F., & Lafferty, K. D. (2006). Evidence for the role of
746 infectious disease in species extinction and endangerment. *Conservation*
747 *Biology*, 20(5), 1349–1357. doi:10.1111/j.1523-1739.2006.00524.x

748 Spurgin, L. G., & Richardson, D. S. (2010). How pathogens drive genetic
749 diversity: MHC, mechanisms and misunderstandings. *Proceedings of the*
750 *Royal Society B: Biological Sciences*, 277(1684), 979–88.
751 doi:10.1098/rspb.2009.2084

752 Staubach, F., Lorenc, A., Messer, P. W., Tang, K., Petrov, D. A., & Tautz, D.
753 (2012). Genome patterns of selection and introgression of haplotypes in

754 natural populations of the house mouse (*Mus musculus*). *PLoS Genetics*,
755 8(8), e1002891. doi:10.1371/journal.pgen.1002891

756 Steere, A. C., Coburn, J., & Glickstein, L. (2004). The emergence of Lyme
757 disease. *Journal of Clinical Investigation*, 113(8), 1093–1101.
758 doi:10.1172/JCI200421681

759 Tagu, D., Colbourne, J. K., & Nègre, N. (2014). Genomic data integration for
760 ecological and evolutionary traits in non-model organisms. *BMC*
761 *Genomics*, 15(1), 490. doi:10.1186/1471-2164-15-490

762 Togninalli, M., Seren, Ü., Meng, D., Fitz, J., Nordborg, M., Weigel, D., ...
763 Grimm, D. G. (2018). The AraGWAS Catalog: A curated and
764 standardized *Arabidopsis thaliana* GWAS catalog. *Nucleic Acids*
765 *Research*, 46(D1), D1150–D1156. doi:10.1093/nar/gkx954

766 Tschirren, B., Andersson, M., Scherman, K., Westerdahl, H., Mittl, P. R., &
767 Råberg, L. (2013). Polymorphisms at the innate immune receptor *TLR2*
768 are associated with *Borrelia* infection in a wild rodent population.
769 *Proceedings of the Royal Society B: Biological Sciences*, 280, 20130364.
770 doi:10.1098/rspb.2013.0364

771 Tschirren, B. (2015). *Borrelia burgdorferi* sensu lato infection pressure shapes
772 innate immune gene evolution in natural rodent populations across
773 Europe. *Biology Letters*, 11, 20150263. doi:10.1098/rsbl.2015.0263

774 Tschirren, B. & Bensch, S. (2010). Genetics of personalities: no simple
775 answers for complex traits. *Molecular Ecology*, 19(4), 624–626. doi:
776 10.1111/j.1365-294X.2009.04519.x

777 Van der Auwera, G. A., Carneiro, M. O., Hartl, C., Poplin, R., del Angel, G.,
778 Levy-Moonshine, A., ... DePristo, M. A. (2013). From FastQ data to high-

779 confidence variant calls: the genome analysis toolkit best practices
780 pipeline. In *Current Protocols in Bioinformatics*. John Wiley & Sons, Inc.
781 doi:10.1002/0471250953.bi1110s43

782 Visscher, P. M., Wray, N. R., Zhang, Q., Sklar, P., McCarthy, M. I., Brown, M.
783 A., & Yang, J. (2017). 10 years of GWAS discovery: biology, function,
784 and translation. *American Journal of Human Genetics*, *101*(1), 5–22.
785 doi:10.1016/j.ajhg.2017.06.005

786 Vitti, J. J., Grossman, S. R., & Sabeti, P. C. (2013). Detecting natural
787 selection in genomic data. *Annual Review of Genetics*, *47*, 97–120.
788 doi:10.1146/annurev-genet-111212-133526

789 Wangler, M. F., Hu, Y., & Shulman, J. M. (2017). *Drosophila* and genome-
790 wide association studies: a review and resource for the functional
791 dissection of human complex traits. *Disease Models & Mechanisms*,
792 *10*(2), 77–88. doi:10.1242/dmm.027680

793 Wiethoelter, A. K., Beltrán-Alcrudo, D., Kock, R., & Mor, S. M. (2015). Global
794 trends in infectious diseases at the wildlife–livestock interface.
795 *Proceedings of the National Academy of Sciences USA*, *112*(31), 9662–
796 9667. doi:10.1073/pnas.1422741112

797 Wilfert, L., & Schmid-Hempel, P. (2008). The genetic architecture of
798 susceptibility to parasites. *BMC Evolutionary Biology*, *8*(1), 1–8.
799 doi:10.1186/1471-2148-8-187

800 Wooten, R. M., Ma, Y., Yoder, R. A., Brown, J. P., Weis, J. H., Zachary, J. F.,
801 ... Weis, J. J. (2002). Toll-like receptor 2 is required for innate, but not
802 acquired, host defense to *Borrelia burgdorferi*. *Journal of Immunology*,
803 *168*(1), 348–355.

804 Zerbino, D. R., Achuthan, P., Akanni, W., Amode, M. R., Barrell, D., Bhai, J.,
805 ... Flicek, P. (2018). Ensembl 2018. *Nucleic Acids Research*, 46(D1),
806 D754–D761. doi:10.1093/nar/gkx1098

807 Zueva, K. J., Lumme, J., Veselov, A. E., Kent, M. P., Lien, S., & Primmer, C.
808 R. (2014). Footprints of directional selection in wild atlantic salmon
809 populations: Evidence for parasite-driven evolution? *PLoS ONE*, 9(3).
810 doi:10.1371/journal.pone.0091672

811

812 **Tables**

813

814 **Table 1. Sampling locations and number of analysed bank voles**

815 Elevation and study site coordinates, the number of genotyped adult bank voles (N), the number of genotyped *Borrelia*-free bank
 816 voles (N uninf) and the number of genotyped *Borrelia*-infected bank voles (N inf) and *Borrelia* prevalence in adult bank voles at the
 817 study sites are reported.

818

| Location | Label | Elevation (masl) | North | East | N | N uninf | N inf | <i>Borrelia</i> prevalence (%) |
|-----------------|--------------|-------------------------|--------------|-------------|----------|----------------|--------------|---------------------------------------|
| Bonaduz | BON | 944 | 46.799 | 9.352 | 16 | 8 | 8 | 50.0 |
| Rodels | ROD | 630 | 46.760 | 9.425 | 17 | 13 | 4 | 31.2 |
| Sagogn | SAG | 693 | 46.783 | 9.233 | 19 | 10 | 9 | 48.4 |
| Flims | FLI | 1138 | 46.827 | 9.280 | 15 | 9 | 6 | 54.5 |
| Malans | MAL | 560 | 46.992 | 9.557 | 19 | 10 | 9 | 44.8 |
| Passugg | PAS | 732 | 46.840 | 9.538 | 13 | 6 | 7 | 61.5 |

| | | | | | | | | |
|---------|-----|-----|--------|-------|----|---|----|------|
| Trimmis | TRI | 762 | 46.882 | 9.559 | 19 | 9 | 10 | 44.8 |
|---------|-----|-----|--------|-------|----|---|----|------|

819

820

821 **Table 2. F_{ST} outlier SNPs**

822 Exonic SNPs that were identified as outliers in multiple populations when comparing *Borrelia*-infected and *Borrelia*-free bank voles
 823 using a F_{ST} -based approach. SNPs in bold were also identified to be associated with *Borrelia*-infection status using a GWAS
 824 approach. The SNP position refers to the prairie vole genome version MicOch1.0.

825

| Chromosome | SNP Position | Start of exon | End of exon | Protein ID | Gene description | Number of populations in which the SNP was an outlier |
|------------|-----------------|-----------------|-----------------|---------------------------|---|---|
| 1 | 82812170 | 82812122 | 82812285 | ENSMUSP00000001253 | Solute carrier family 26, member 4 | 3 |
| 5 | 90628057 | 90626456 | 90630654 | ENSMUSP0000000081880 | Golgin subfamily A member 4 | 3 |
| 7 | 81147668 | 81147563 | 81147724 | ENSMUSP0000000081864 | Dynein, axonemal, heavy chain 17 | 3 |

| | | | | | | |
|-----|---------|---------|----------|----------------|-----------------------|---|
| 7 | 2713691 | 2713673 | 27136920 | ENSMUSP0000010 | WSC domain containing | 4 |
| | 7 | 6 | | 4150 | 1 | |
| 7 | 2581122 | 2581081 | 25812290 | ENSMUSP0000005 | Haspin | 3 |
| | 6 | 2 | | 5806 | | |
| 8 | 7265577 | 7265562 | 72655797 | ENSMUSP0000005 | Sideroflexin 3 | 3 |
| | 7 | 7 | | 9419 | | |
| LG1 | 8299897 | 8299839 | 8299951 | ENSMUSP0000002 | Tensin 3 | 4 |
| | | | | 0695 | | |
| LG4 | 6044582 | 6044511 | 60446729 | ENSMUSP0000008 | Espin-like | 3 |
| | 2 | 6 | | 6294 | | |
| LG5 | 3878388 | 3878385 | 38783985 | ENSMUSP0000003 | Talin-1 | 3 |
| | 1 | 7 | | 0187 | | |

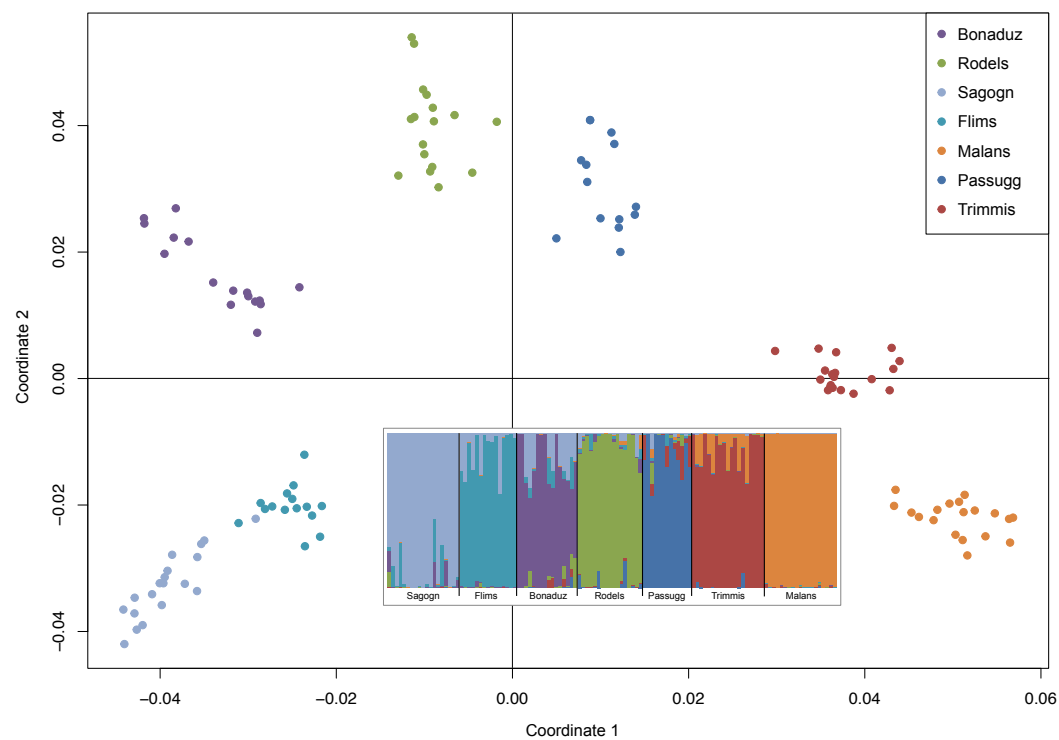
826

827

828 **Figures**

829 **Figure 1. Multi-dimensional scaling of bank vole genetic diversity.**

830 Different colours represent different sampling sites. The inset shows the proportion of ancestry for each sampled bank vole (N =
831 118) for seven genetic clusters inferred with STRUCTURE (see Supplementary Figures S2 and S3 for additional information).



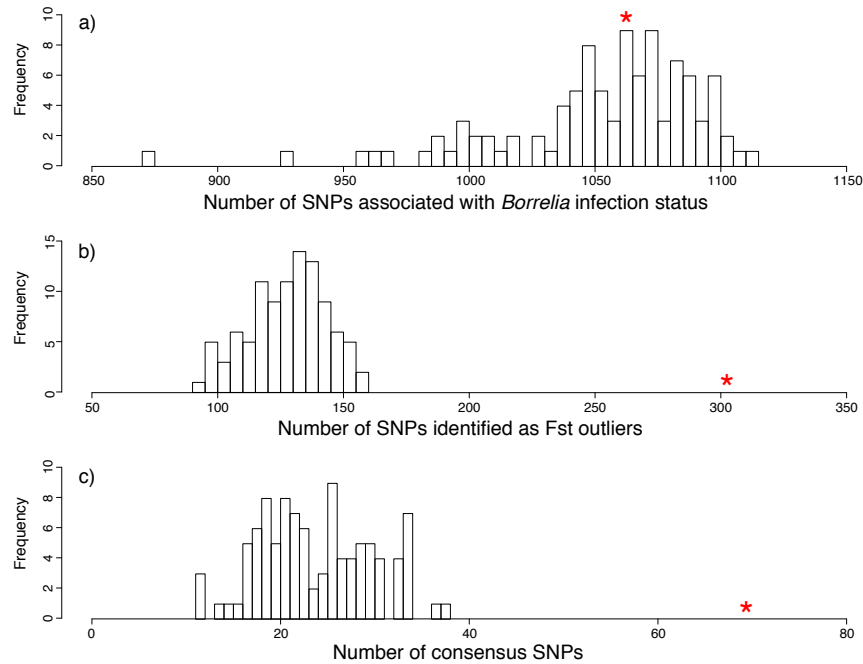
832

833

834 **Figure 2. Simulations of false-positive rates.**

835 We used a simulation approach to quantify the false positive rate of the GWAS approach (a), the F_{ST} -based approach with a three
836 population replication criterion (b), and (c) the combined approach (i.e. consensus SNPs identified in (a) and (b)). The red asterisk
837 indicates the number of identified SNPs observed in the real data using the respective approach.

838

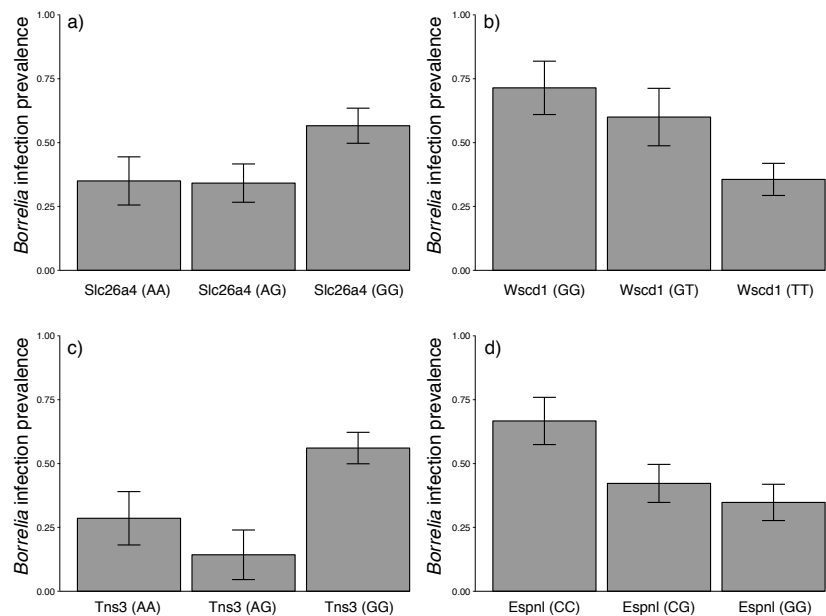


839

840 **Figure 3. Genetic polymorphisms at the four consensus candidate loci are associated with *Borrelia* infection status in**
841 **bank voles.**

842 Animals that carried the rarer allele of *Slc26a4* (a, allele A) and *Tns3* (c, allele A) were less likely to be *Borrelia*-infected, whereas
843 animals that carried the rarer allele of *Wscd1* (b, allele G) and *Espnl* (d, allele C) were more likely to be *Borrelia*-infected. Error bars
844 represent standard errors.

845



846

Supporting information

Combining GWAS and F_{ST} -based approaches to identify targets of *Borrelia*-mediated selection in natural rodent hosts

L. Cornetti & B. Tschirren

1. Supporting Methods

Neutral simulations

Script used for the simulations

```
initialize() {
  initializeMutationRate(2e-4); ## mutation rate
  initializeMutationType("m1", 0.5, "f", 0.0); ## mutation type description: non-
coding or synonymous
  initializeGenomicElementType("g1", c(m1), c(100)); ## mutation occurrence
  initializeGenomicElement(g1, 0, 27999); ## size of the simulated
chromosome
  initializeRecombinationRate(0.05); ## recombination rate
}
1 { ## at generation 1 seven subpopulations appear
  sim.addSubpop("p1", 100); ## population size of p1
  sim.addSubpop("p2", 100); ## population size of p2
  sim.addSubpop("p3", 100); ## population size of p3
  sim.addSubpop("p4", 100); ## population size of p4
  sim.addSubpop("p5", 100); ## population size of p5
  sim.addSubpop("p6", 100); ## population size of p6
  sim.addSubpop("p7", 100); ## population size of p7
  p1.setMigrationRates(c(p2,p3,p4,p5,p6,p7),
c(0.05,0.05,0.05,0.01,0.03,0.01)); ## migration rates into population p1 from
the others
  p2.setMigrationRates(c(p1,p3,p4,p5,p6,p7),
c(0.05,0.03,0.03,0.01,0.05,0.03)); ## migration rates into population p2 from
the others
  p3.setMigrationRates(c(p1,p2,p4,p5,p6,p7),
c(0.05,0.03,0.05,0.01,0.01,0.01)); ## migration rates into population p3 from
the others
  p4.setMigrationRates(c(p1,p2,p3,p5,p6,p7),
c(0.05,0.03,0.05,0.01,0.03,0.03)); ## migration rates into population p4 from
the others
  p5.setMigrationRates(c(p1,p2,p3,p4,p6,p7),
c(0.01,0.01,0.01,0.01,0.03,0.03)); ## migration rates into population p5 from
the others
```

```

    p6.setMigrationRates(c(p1,p2,p3,p4,p5,p7),
c(0.03,0.05,0.01,0.03,0.01,0.05)); ## migration rates into population p6 from
the others
    p7.setMigrationRates(c(p1,p2,p3,p4,p5,p6),
c(0.01,0.03,0.01,0.03,0.05,0.03)); ## migration rates into population p7 from
the others
}
4999 late() { ## number of generation simulated
    bonaduz = p1.sampleIndividuals(16).genomes; ## number of samples
selected from p1 according to the sample size of Bonaduz
    bonaduz.outputVCF(filePath="/home/p1.vcf"); ## the SNPs are written in a
VCF file
    rodels = p2.sampleIndividuals(17).genomes; ## number of samples
selected from p2 according to the sample size of Rodels
    rodels.outputVCF(filePath="/home/p2.vcf"); ## the SNPs are written in a
VCF file
    sagogn = p3.sampleIndividuals(19).genomes; ## number of samples
selected from p3 according to the sample size of Sagogn
    sagogn.outputVCF(filePath="/home/p3.vcf"); ## the SNPs are written in a
VCF file
    flims = p4.sampleIndividuals(15).genomes; ## number of samples selected
from p4 according to the sample size of Flims
    flims.outputVCF(filePath="/home/p4.vcf"); ## the SNPs are written in a VCF
file
    malans = p5.sampleIndividuals(19).genomes; ## number of samples
selected from p5 according to the sample size of Malans
    malans.outputVCF(filePath="/home/p5.vcf"); ## the SNPs are written in a
VCF file
    passugg = p6.sampleIndividuals(13).genomes; ## number of samples
selected from p6 according to the sample size of Passugg
    passugg.outputVCF(filePath="/home/p6.vcf"); ## the SNPs are written in a
VCF file
    trimmis = p7.sampleIndividuals(19).genomes; ## number of samples
selected from p7 according to the sample size of Trimmis
    trimmis.outputVCF(filePath="/home/p7.vcf"); ## the SNPs are written in a
VCF file
}

```

Supporting Figures

Figure S1. Map of the sampling sites in the Kanton Graubünden, Switzerland.

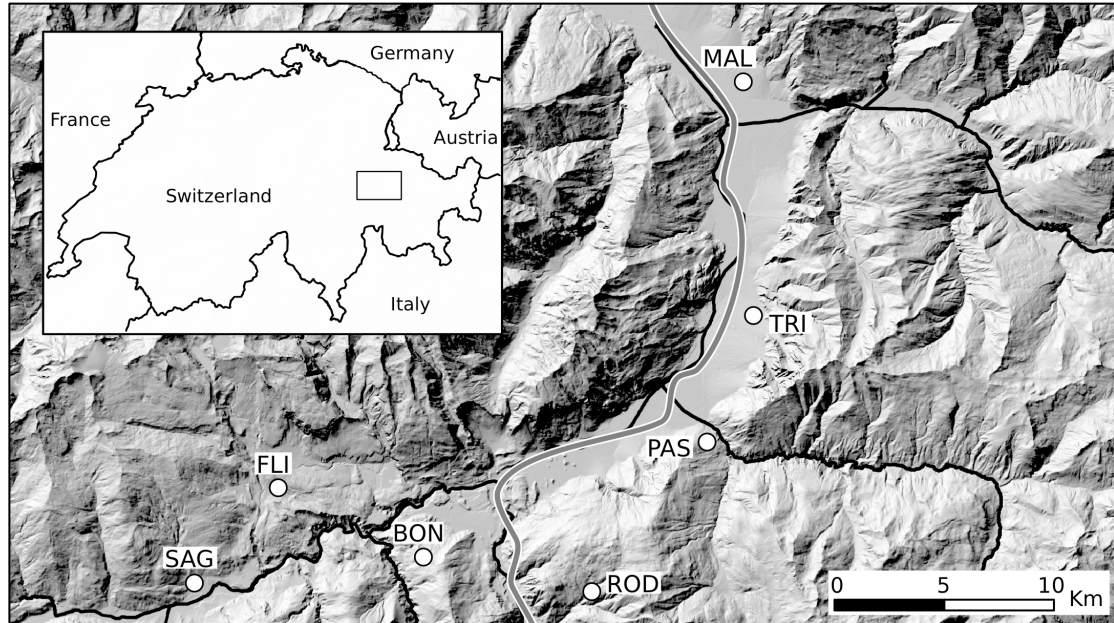


Figure S2. Estimate of the number of genetically well defined groups (K) based on mean likelihood (Pritchard, Stephens, & Donnelly, 2000).

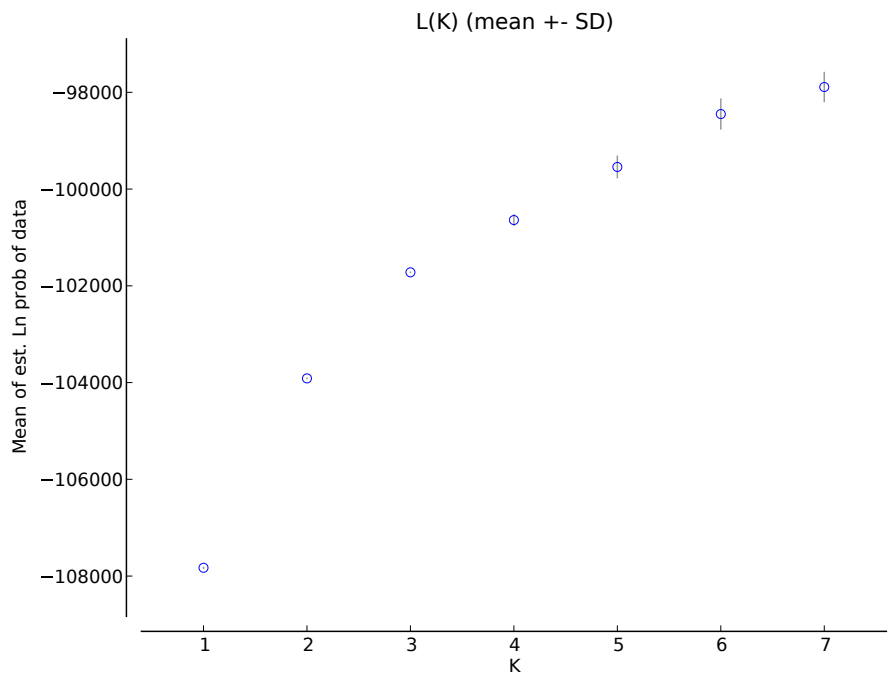


Figure S3. STRUCTURE plot describing the bank vole population structure in the study area using K=2 to K=7 as most probable number of genetic groups.

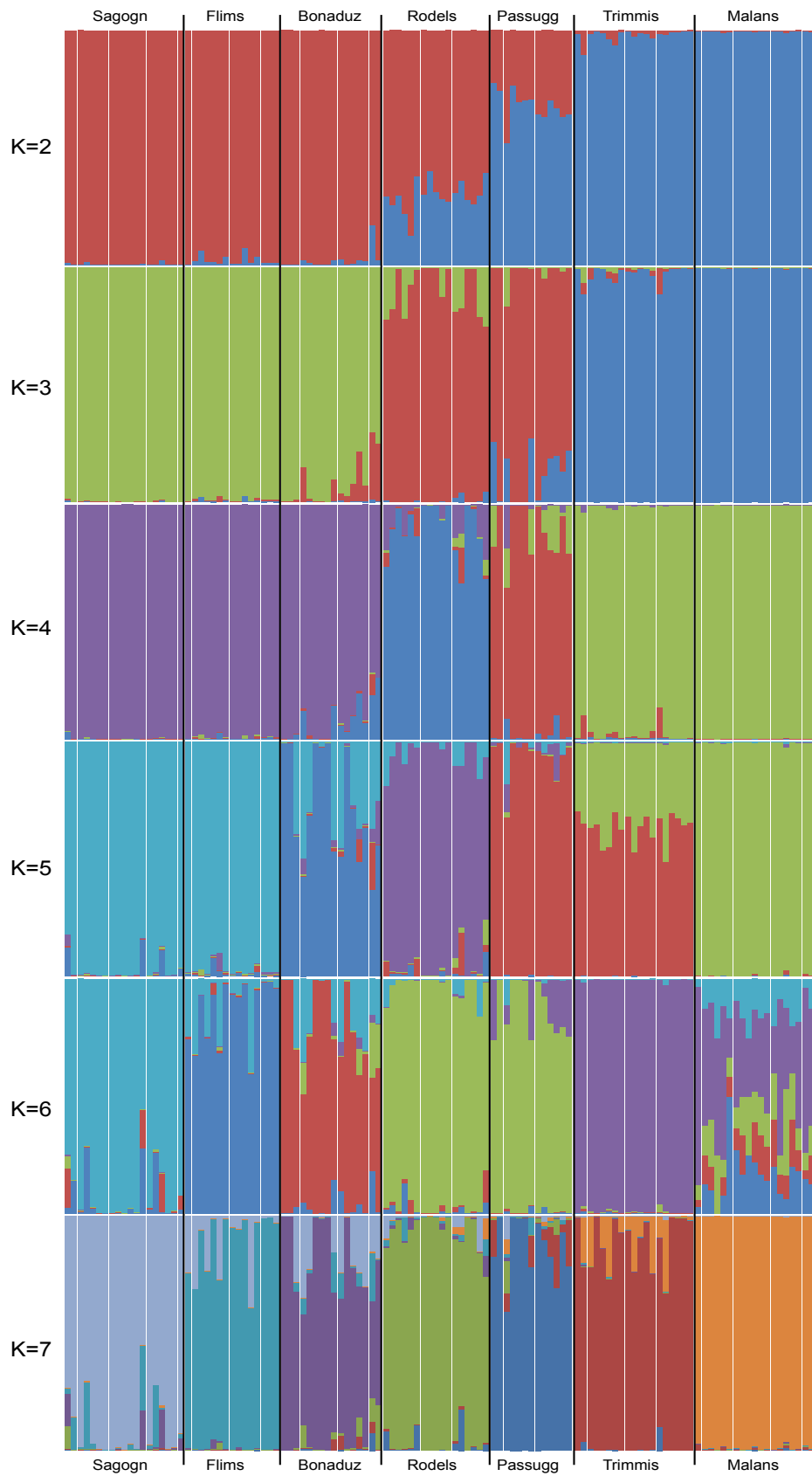


Figure S4. Isolation-by-distance of bank vole populations across our study sites.

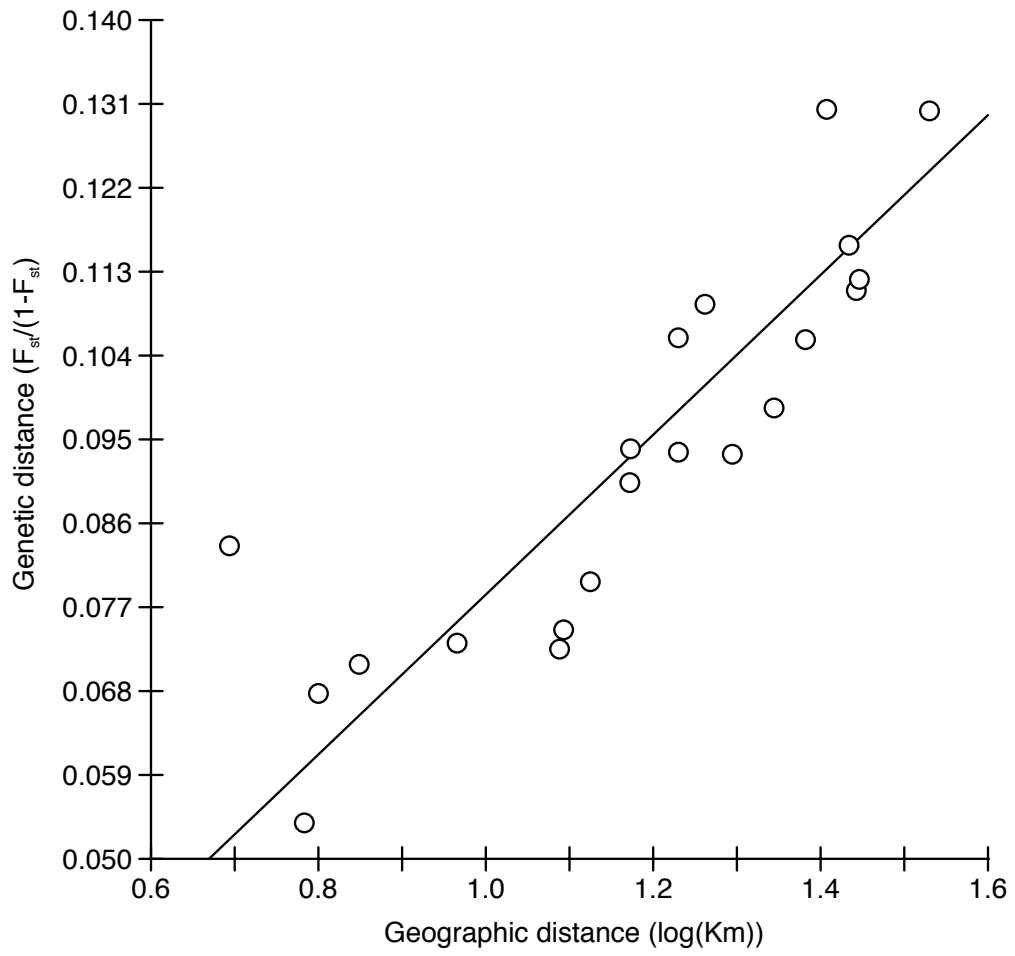


Figure S5 False positive rates of F_{ST} -based approach.

We used a simulation approach to quantify the false positive rate of the F_{ST} -based approach when using no population replication (a), a two population replication criterion (b), and a three population replication criterion (as used in the main study) (c). The red asterisk indicates the number of outliers observed in the real data using the respective approach.

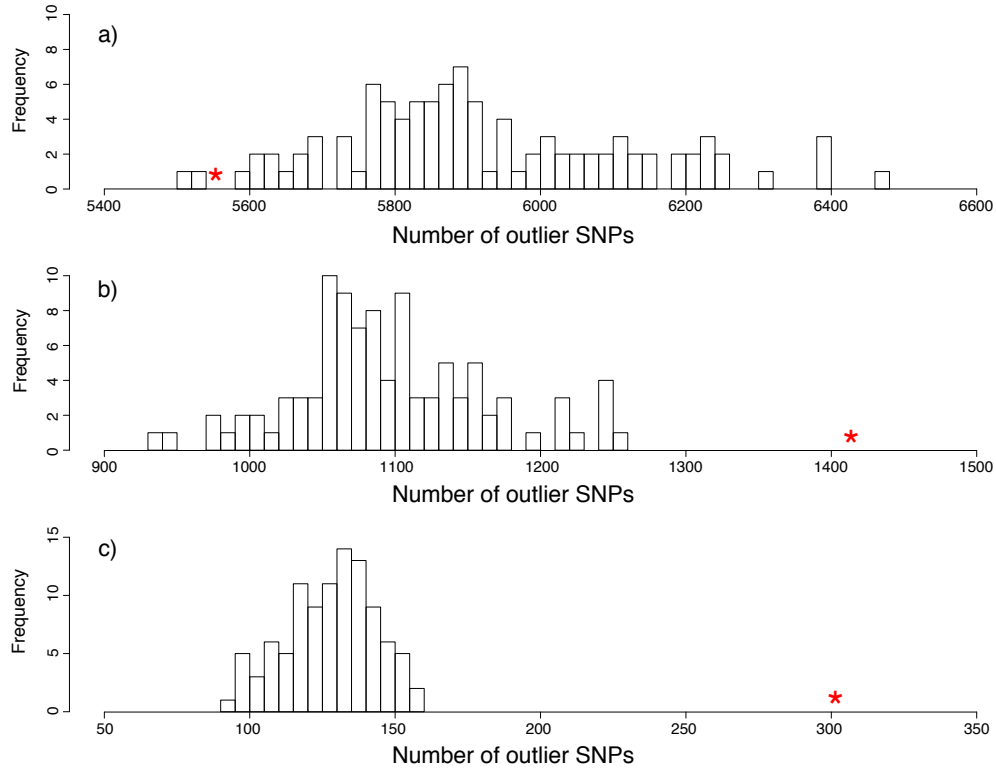
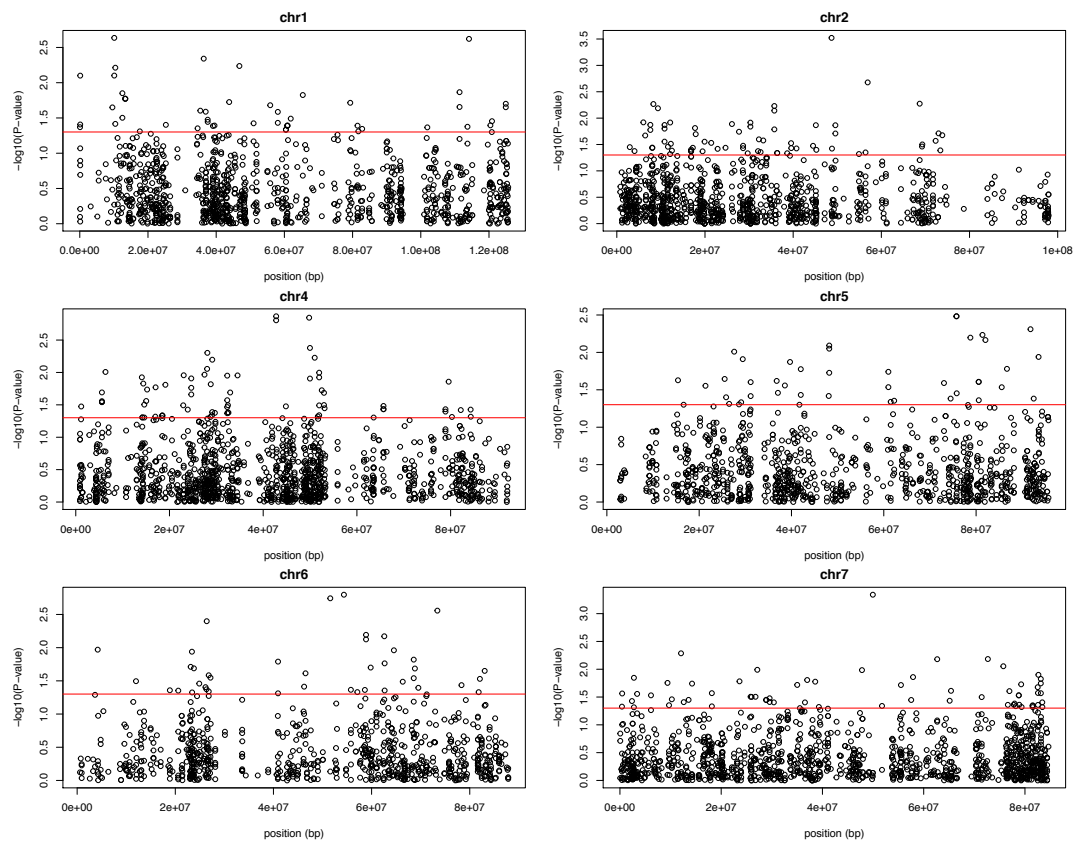
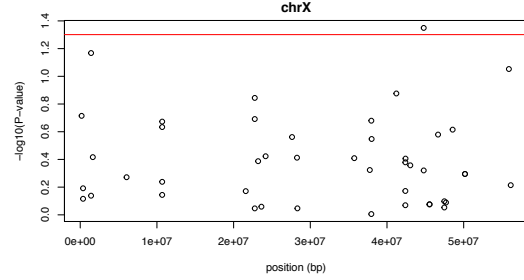
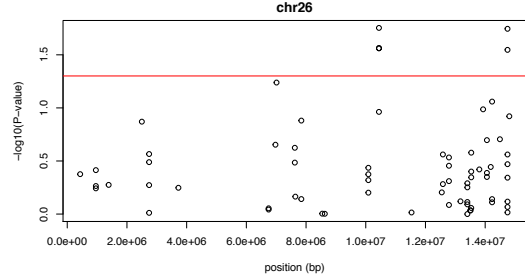
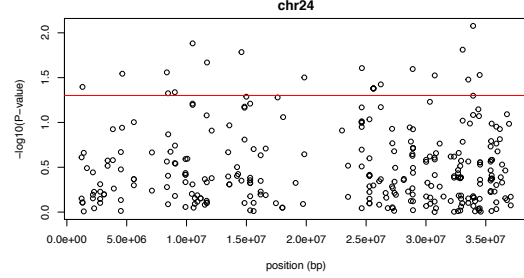
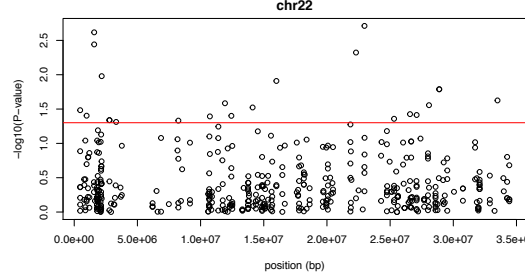
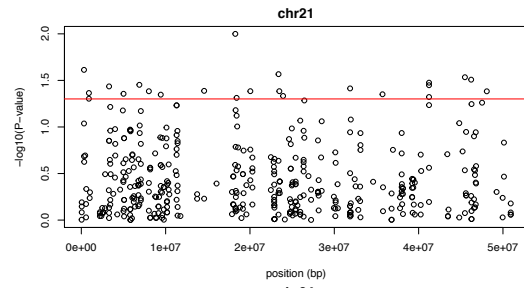
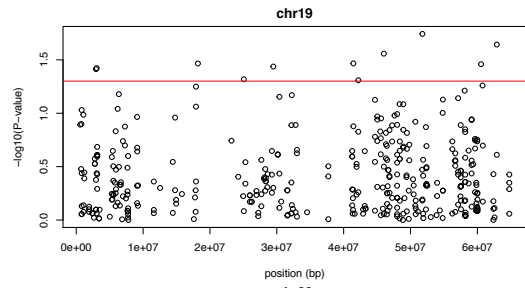
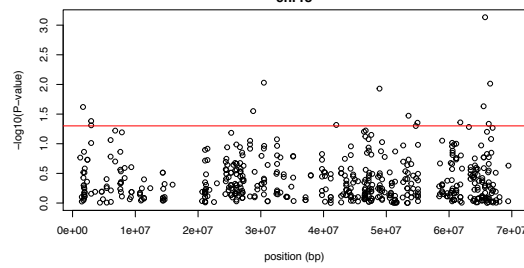
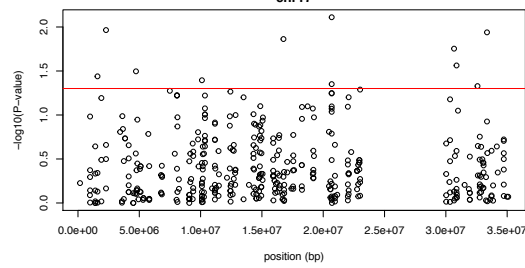
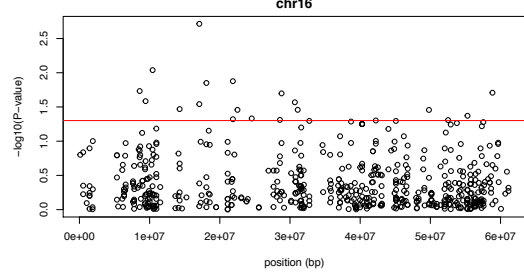
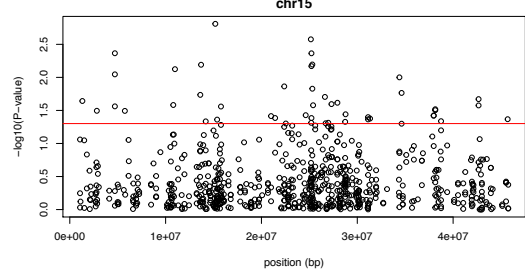
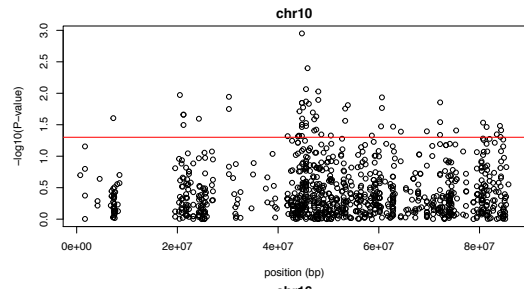
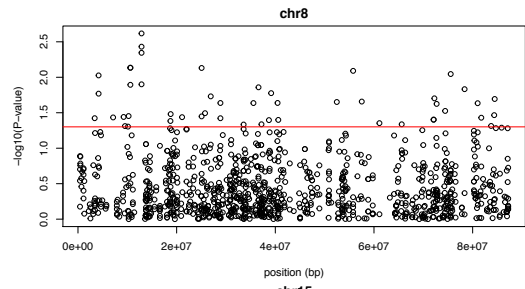


Figure S6. Association scan for *Borrelia* infection status in bank voles using a GWAS approach. Associations between SNPs and bank vole *Borrelia* infection status after correcting for population structure ($-\log_{10}(P\text{-value})$) are plotted for each chromosome and linkage group. The red line represents the $P = 0.05$ threshold. Dots above this line represent SNPs possibly associated with *Borrelia* infection status.





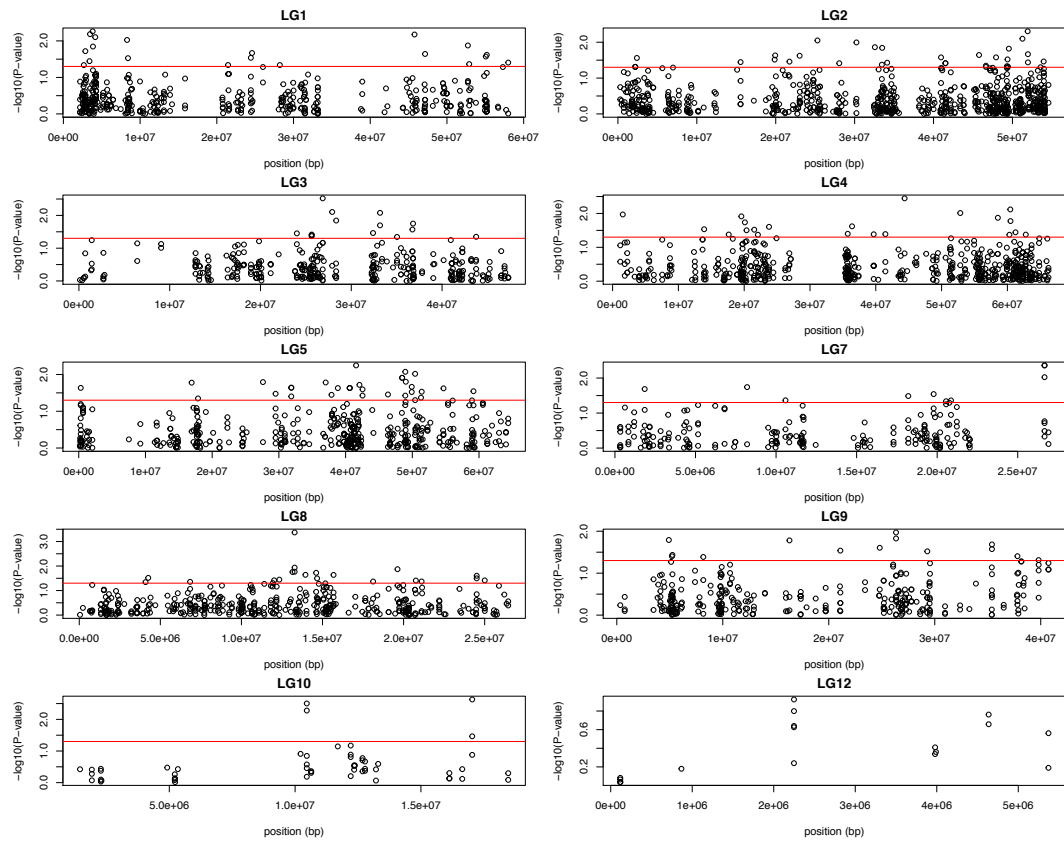


Figure S7. Number of putatively *Borrelia* infection status associated SNPs per MB across chromosomes and linkage groups. Candidate SNPs were identified using a GWAS approach (see Figure S6). The main 28 scaffolds of the prairie vole (*Microtus ochrogaster*) reference genome, version MicOch1.0, are shown.

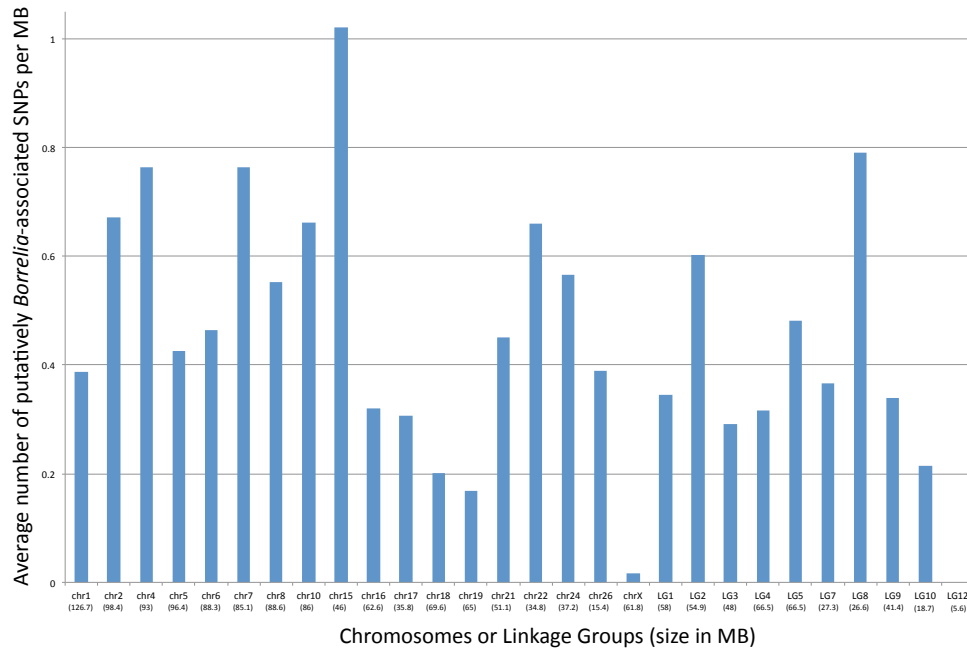
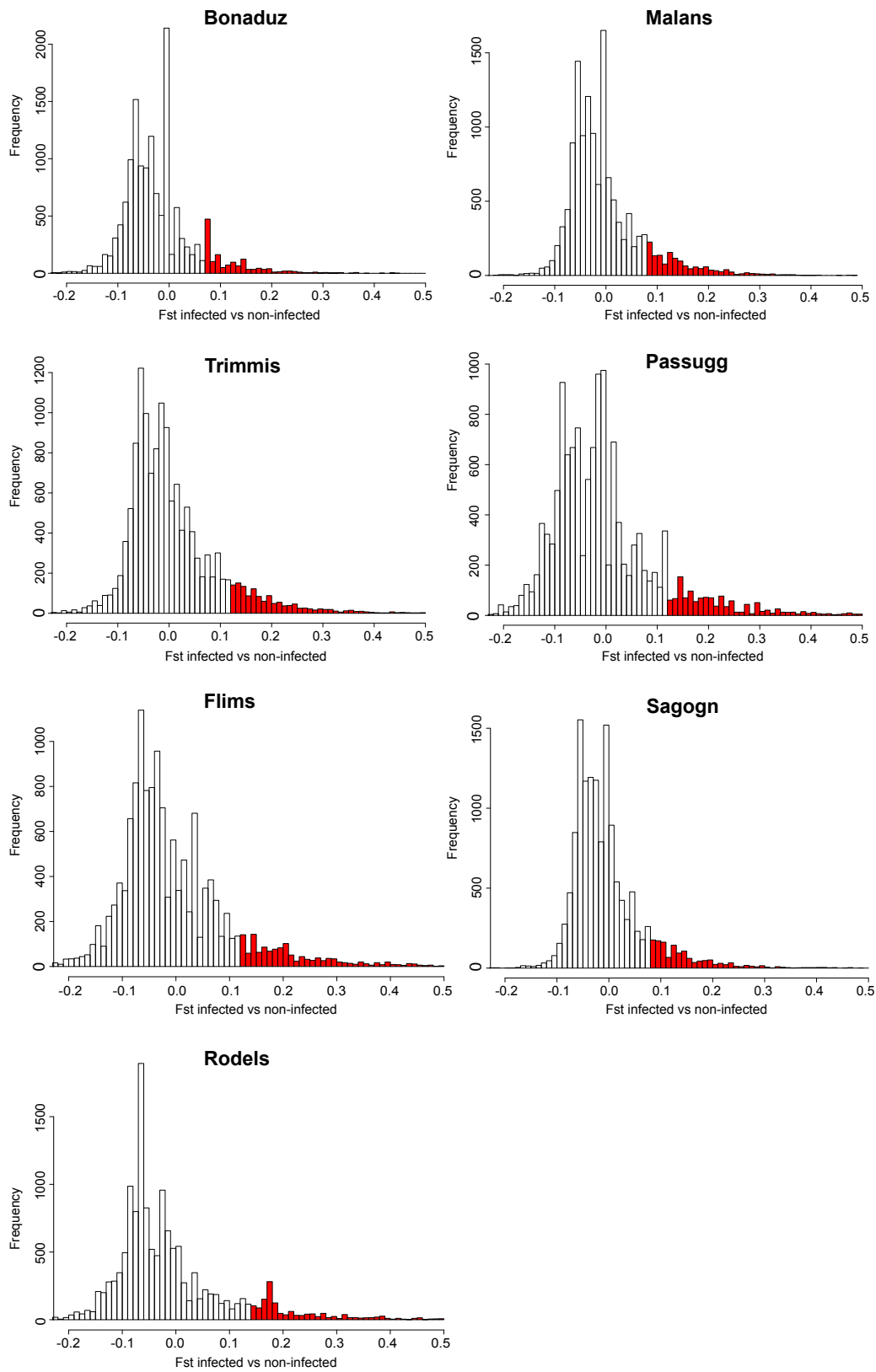


Figure S8. F_{ST} distribution of SNPs across the seven study populations. The red colour represents the SNPs falling into the 10% highest F_{ST} values when comparing *Borrelia*-infected vs uninfected bank voles within each population.



Supporting Tables

Table S1. Pairwise F_{ST} based on 21,811 genome-wide SNPs (above the diagonal) and linear distance (km) between locations (below the diagonal). All F_{ST} values are statistically significant ($p < 0.05$).

| | Bonaduz | Rodels | Sagogn | Flims | Malans | Passugg | Trimmis |
|---------|---------|--------|--------|-------|--------|---------|---------|
| Bonaduz | - | 0.063 | 0.068 | 0.062 | 0.110 | 0.083 | 0.092 |
| Rodels | 7.06 | - | 0.084 | 0.073 | 0.101 | 0.069 | 0.083 |
| Sagogn | 9.24 | 14.86 | - | 0.052 | 0.111 | 0.095 | 0.099 |
| Flims | 6.31 | 13.33 | 6.07 | - | 0.094 | 0.086 | 0.087 |
| Malans | 25.55 | 27.71 | 33.90 | 27.95 | - | 0.096 | 0.066 |
| Passugg | 14.89 | 12.39 | 24.09 | 19.70 | 16.98 | - | 0.076 |
| Trimmis | 18.27 | 16.99 | 27.16 | 22.10 | 12.25 | 4.94 | - |

Table S2. List of 53 exonic SNPs putatively associated with *Borrelia* infection status in bank voles, identified using a GWAS approach. Genes with more than one SNP per GBS read are underlined, and genes for which SNPs were observed in two different exons and two different GBS reads are highlighted with an asterisk (*). Candidate SNPs identified with both the GWAS and the F_{ST} -based approach with population replication are highlighted in bold. The SNP position refers to the prairie vole genome version MicOch1.0.

| Chromosome Number | SNP Position | Start of exon | End of exon | Protein ID | Gene description |
|-------------------|-----------------|-----------------|-----------------|---------------------------|--|
| LG2 | 46643179 | 46643140 | 46643278 | ENSMUSP00000001183 | formiminotransferase cyclodeaminase |
| 1 | 82812170 | 82812122 | 82812285 | ENSMUSP00000001253 | solute carrier family 26, member 4 |
| <u>21</u> | <u>900079</u> | <u>899816</u> | <u>900277</u> | <u>ENSMUSP00000005015</u> | <u>papillary renal cell carcinoma (translocation-associated)</u> |
| <u>21</u> | <u>900160</u> | <u>899816</u> | <u>900277</u> | <u>ENSMUSP00000005015</u> | <u>papillary renal cell carcinoma (translocation-associated)</u> |
| 10 | 43932666 | 43932619 | 43932690 | ENSMUSP00000010007 | succinate dehydrogenase complex, subunit B, iron sulfur (Ip) |
| <u>15</u> | <u>25218452</u> | <u>25218178</u> | <u>25218816</u> | <u>ENSMUSP00000016172</u> | <u>cadherin, EGF LAG seven-pass G-type receptor 1</u> |
| <u>15</u> | <u>25218467</u> | <u>25218178</u> | <u>25218816</u> | <u>ENSMUSP00000016172</u> | <u>cadherin, EGF LAG seven-pass G-type receptor 1</u> |
| <u>15</u> | <u>25218528</u> | <u>25218178</u> | <u>25218816</u> | <u>ENSMUSP00000016172</u> | <u>cadherin, EGF LAG seven-pass G-type receptor 1</u> |
| LG8 | 4254355 | 4253368 | 4254449 | ENSMUSP00000016399 | tubulin, beta 1 class VI |
| LG1 | 8299897 | 8299839 | 8299951 | ENSMUSP00000020695 | tensin 3 |
| 7 | 29565141 | 29565139 | 29565229 | ENSMUSP00000021259 | guanylate cyclase 2e |
| 15 | 2823221 | 2823152 | 2823685 | ENSMUSP00000023720 | keratin 84 |
| 24 | 8453004 | 8453003 | 8453147 | ENSMUSP00000026500 | advillin |
| 22 | 11931974 | 11931956 | 11932112 | ENSMUSP00000026922 | homer scaffolding protein 2 |
| LG5 | 54679380 | 54679319 | 54680259 | ENSMUSP00000027035 | SRY (sex determining region Y)-box 17 |
| LG2 | 23030453 | 23030389 | 23030539 | ENSMUSP00000027257 | MIT, microtubule interacting and transport, domain containing 1 |
| LG4 | 60476015 | 60476012 | 60476048 | ENSMUSP00000027532 | selenocysteine lyase |
| 6 | 23429687 | 23428593 | 23430689 | ENSMUSP00000027706 | leucine rich repeat protein 2, neuronal |
| 5 | 74480869 | 74480842 | 74481033 | ENSMUSP00000034961 | immunoglobulin superfamily, DCC subclass, member 3 |

| | | | | | |
|------------|-----------------|-----------------|-----------------|---------------------------|--|
| 5 | 81380322 | 81380142 | 81380335 | ENSMUSP00000035232 | kelch domain containing 8B |
| LG4 | 19571145 | 19570722 | 19571251 | ENSMUSP00000036699 | immunity-related GTPase family, Q |
| 4 | 17036091 | 17035966 | 17036136 | ENSMUSP00000039271 | fucokinase |
| 7 | 83424112* | 83423459 | 83425210 | ENSMUSP00000043643 | BAH domain and coiled-coil containing 1 |
| 7 | 83439688* | 83439647 | 83439722 | ENSMUSP00000043643 | BAH domain and coiled-coil containing 1 |
| 2 | 2946490 | 2946456 | 2946531 | ENSMUSP00000046544 | small G protein signaling modulator 1 |
| 24 | 33941679 | 33941663 | 33941821 | ENSMUSP00000048309 | stabilin 2 |
| 21 | 14567230 | 14566282 | 14567610 | ENSMUSP00000052306 | phospholipid phosphatase related 4 |
| 7 | 25811175 | 25810812 | 25812290 | ENSMUSP00000055806 | germ cell associated 2, haspin |
| 22 | 25313647 | 25313056 | 25313865 | ENSMUSP00000073855 | synemin, intermediate filament protein |
| 2 | 6446782 | 6446705 | 6447009 | ENSMUSP00000075488 | serine/arginine repetitive matrix 4 |
| LG4 | 13915780 | 13915721 | 13915902 | ENSMUSP00000076093 | utrophin |
| <u>2</u> | <u>16862787</u> | <u>16862779</u> | <u>16862865</u> | <u>ENSMUSP00000078198</u> | <u>rabphilin 3A</u> |
| <u>2</u> | <u>16862810</u> | <u>16862779</u> | <u>16862865</u> | <u>ENSMUSP00000078198</u> | <u>rabphilin 3A</u> |
| 15 | 26699638 | 26699630 | 26699696 | ENSMUSP00000079575 | RIKEN cDNA 1810041L15 gene |
| 4 | 79583497 | 79583457 | 79583535 | ENSMUSP00000079752 | low density lipoprotein receptor-related protein 2 |
| 7 | 77798465 | 77798372 | 77798499 | ENSMUSP00000081398 | kinesin family member 19A |
| LG4 | 60445822 | 60445116 | 60446729 | ENSMUSP00000086294 | espin-like |
| 15 | 28737242 | 28737225 | 28737356 | ENSMUSP00000086582 | meiotic double-stranded break formation protein 1 |
| 7 | 78870838 | 78870772 | 78870933 | ENSMUSP00000091439 | myosin XVB |
| 6 | 23134593 | 23134517 | 23134713 | ENSMUSP00000092148 | neurofascin |
| LG9 | 16293329 | 16293222 | 16293422 | ENSMUSP00000093587 | spectrin repeat containing, nuclear envelope 1 |
| LG8 | 25004183 | 25003387 | 25004577 | ENSMUSP00000096813 | zinc finger, CCHC domain containing 3 |
| 8 | 74431200 | 74431102 | 74432722 | ENSMUSP00000096972 | cyclin M2 |
| 2 | 4016724 | 4016690 | 4016844 | ENSMUSP00000099642 | acetyl-Coenzyme A carboxylase beta |
| 10 | 75362360 | 75362298 | 75362460 | ENSMUSP00000102320 | podocan |
| 22 | 2174486 | 2174390 | 2174566 | ENSMUSP00000102745 | myosin VIIA |

| | | | | | |
|----|-----------------|-----------------|-----------------|---------------------------|---|
| 7 | 27136917 | 27136736 | 27136920 | ENSMUSP00000104150 | WSC domain containing 1 |
| 7 | 30401083 | 30400232 | 30401249 | ENSMUSP00000104314 | netrin 1 |
| 5 | 15411871 | 15411823 | 15411894 | ENSMUSP00000111093 | adaptor protein complex AP-1, mu 2 subunit leucine-rich repeat, immunoglobulin-like and transmembrane domains 1 |
| 6 | 55777259 | 55777095 | 55777561 | ENSMUSP00000113964 | |
| 8 | 38709839 | 38709811 | 38709966 | ENSMUSP00000115356 | aldehyde dehydrogenase 3 family, member B2 |
| 6 | 40936018 | 40935948 | 40936067 | ENSMUSP00000126464 | acyl-Coenzyme A oxidase 2, branched chain |
| 15 | 10977462 | 10977321 | 10977467 | ENSMUSP00000129100 | 5-oxoprolinase (ATP-hydrolysing) |

Table S3. Allele frequencies at the four consensus candidate genes identified with both the GWAS and the F_{ST} -based approach. Frequencies of the reference allele are shown.

| | | | Gene | | | |
|------------------|-----|-----------------------------------|----------|----------|---------|------------|
| | N | <i>Borrelia</i> -infection status | Slc26a4 | Tns3 | WSC1 | Espin-like |
| Chromosome | | | chr1 | chr7 | LG1 | LG4 |
| Position | | | 82812170 | 27136917 | 8299897 | 60445822 |
| Alleles | | | A / G | G / A | G / T | G / C |
| Reference allele | | | A | G | G | G |
| | 118 | all | 35.5 | 83.9 | 19.8 | 58.0 |
| | 53 | infected | 27.5 | 92.7 | 28.9 | 48.1 |
| | 65 | uninfected | 42 | 76.1 | 12.5 | 66.1 |

Table S4. Analysis of consensus candidate loci

We used generalized linear mixed models with a binomial error structure and location included as a random effect to test for associations between consensus SNPs and *Borrelia* infection status of bank voles. Models that combined the heterozygous and homozygous state of the less frequent allele are presented here. Models that treated the heterozygous and homozygous state of the less frequent allele as separate genotypes are presented in the main text (see Figure 3 in main text for a visualization of the effects).

| <i>SNP</i> | <i>Test statistics</i> |
|-------------------|-------------------------------------|
| <i>Slc26a4</i> | $\chi^2 = 5.679, DF = 1, P = 0.017$ |
| <i>Tns3</i> | $\chi^2 = 9.346, DF = 1, P = 0.002$ |
| <i>Wscd1</i> | $\chi^2 = 5.639, DF = 1, P = 0.018$ |
| <i>Espin-like</i> | $\chi^2 = 3.162, DF = 1, P = 0.075$ |

References

Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of population structure using multilocus genotype data. *Genetics*, *155*(2), 945–959. doi:10.1111/j.1471-8286.2007.01758.x