# PLOS Pathogens

## Whole genome sequencing of Borrelia burgdorferi isolates reveals linked clusters of plasmid-borne accessory genome elements associated with virulence.
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | PPATHOGENS-D-23-00364 |
| Full Title: | Whole genome sequencing of Borrelia burgdorferi isolates reveals linked clusters of plasmid-borne accessory genome elements associated with virulence. |
| Short Title: | Borrelia burgdorferi genomes reveal linked clusters of accessory elements linked to virulence |
| Article Type: | Research Article |
| Section/Category: | Gram Negative Bacteria |
| Keywords: | Lyme disease, Borrelia burgdorferi, Whole Genome Sequencing;  Virulence; Pathogenicity |
| Abstract: | Lyme disease is the most common vector-borne disease in North America and Europe. The clinical manifestations of Lyme disease vary based on the genospecies of the infecting Borrelia burgdorferi spirochete, but the microbial genetic elements underlying these associations are not known. Here, we report the whole genome sequence (WGS) and analysis of 299 patient-derived B. burgdorferi sensu stricto(Bbss) isolates from patients in the Eastern and Midwestern US and Central Europe. We develop a WGS-based classification of Bbss isolates,confirm and extend the findings of previous single- and multi-locus typing systems, define the plasmid profiles of human-infectious Bbss isolates, annotate the core and strain-variable surface lipoproteome, and identify loci associated with disseminated infection. A core genome consisting of ~800 open reading frames and a core set of plasmids consisting of lp17, lp25, lp36, lp28-3, lp28-4, lp54, and cp26 are found in nearly all isolates. Strain-variable (accessory) plasmids and genes correlate strongly with phylogeny. Using genetic association study methods, we identify an accessory genome signature associated with dissemination and define the individual plasmids and genes that make up this signature. Strains within the RST1/WGS A subgroup, particularly a subset marked by the OspC type A genotype, are associated with increased rates of dissemination. OspC type A strains possess a unique constellation of strongly linked genetic changes including the presence of lp56 and lp28-1 plasmids and a cluster of genes that may contribute to their enhanced virulence compared to other genotypes. The patterns of OspC type A strains typify a broader paradigm across Bbss isolates, in which genetic structure is defined by correlated groups of strain-variable genes located predominantly on plasmids, particularly for expression of surface-exposed lipoproteins. These clusters of genes are inherited in blocks through strain-specific patterns of plasmid occupancy and are associated with the probability of invasive infection. |
| Additional Information: | |
| Question | Response |
| **Government Employee**<br><br>Are you or any of the contributing authors an employee of the United States government?<br><br>Manuscripts authored by one or more US Government employees are not copyrighted, but are licensed under a CC0 Public Domain Dedication, which allows unlimited distribution and reuse of the | No - No authors are employees of the U.S. government. |

| article for any lawful purpose. This is a legal requirement for US Government employees.<br><br>This will be typeset if the manuscript is accepted for publication. | |
|---|---|
| **Financial Disclosure**<br><br>Enter a financial disclosure statement that describes the sources of funding for the work included in this submission and the role the funder(s) played. This includes grants and any commercial funding of the work or authors.<br><br>This statement will be typeset if the manuscript is accepted for publication.<br><br>Review the submission guidelines and the instructions link below for detailed requirements and guidance. | This work was supported by a Doris Duke Charitable Foundation Physician Scientist Fellowship (to J.E.L), the National Institute of Allergy and Infectious Diseases (K99/R00148604 to J.E.L, U19AI110818 and U01 AI151812 to P.C.S.; R01AI045801 to I.S., and R21AI144916 to K.S.), the National Institute of Arthritis and Musculoskeletal and Skin Diseases (R01AR41511 to I.S. and K01AR062098 to K.S.), the Bay Area Lyme Foundation (to P.C.S. and J.E.L.), the Howard Hughes Medical Institute (P.C.S.), the Arthritis Foundation Fellowship (to K.S.), and the Slovenian Research Agency (P3-0296, J3-1744, and J3-8195 to F.S.). |
| **Competing Interests**<br><br>On behalf of all authors, disclose any competing interests that could be perceived to bias this work.<br><br>This statement will be typeset if the manuscript is accepted for publication.<br><br>Review the instructions link below and PLOS Pathogens' competing interests policy to determine what information must be disclosed at submission. | P.C.S. is a co-founder of, shareholder in, and consultant to Sherlock Biosciences and Delve Bio, as well as a board member of and shareholder in Danaher Corporation. K.S. served as a consultant for T2 Biosystems, Roche, BioMerieux, and NYS Biodefense Fund, for the development of a diagnostic assay in Lyme borreliosis. F.S. served on the scientific advisory board for Roche on Lyme disease serological diagnostics and on the scientific advisory board for Pfizer on Lyme disease vaccine, and is an unpaid member of the steering committee of the ESCMID Study Group on Lyme Borreliosis/ESGBOR. J.A.B. has received research funding from Analog Devices Inc., Zeus Scientific, Immunetics, Pfizer, DiaSorin and bioMerieux, and has been a paid consultant to T2 Biosystems, DiaSorin, and Roche Diagnostics.<br>G.P.W. reports receiving research grants from Institute for Systems Biology, Biopeptides, Corp., and Pfizer, Inc. He has been an expert witness in malpractice cases involving Lyme disease and babesiosis; and is an unpaid board member of the non-profit American Lyme Disease Foundation. |
| **Data Availability**<br><br>Provide a **Data Availability Statement** in the box below. This statement should detail where the data used in this submission can be accessed. This statement will be typeset if the manuscript is accepted for publication. | Genome sequences reported here have been deposited in Genbank under PRJNA923804. Code is available at https://github.com/JacobLemieux/borreliaseq. |

Before publication, authors are required to make all data underlying their findings fully available, without restriction. Review our PLOS Data Policy page for detailed information on this policy. Instructions for writing your Data Availability statement can be accessed via the Instructions link below.

1    **Title:** Whole genome sequencing of *Borrelia burgdorferi* isolates reveals linked clusters of

2    plasmid-borne accessory genome elements associated with virulence.

3

4    Jacob E. Lemieux[1,2], Weihua Huang[3,4], Nathan Hill[1,2], Tjasa Cerar[5], Lisa Freimark[2], Sergio

5    Hernandez[6], Matteo Luban[1,2], Vera Maraspin[7], Petra Bogovic[7], Katarina Ogrinc[7], Eva Ruzic-

6    Sabljic[5], Pascal Lapierre[6], Erica Lasek-Nesselquist[6], Navjot Singh[6], Radha Iyer[3], Dionysios

7    Liveris[3], Kurt D. Reed[8], John M. Leong[9], John A. Branda[1], Allen C. Steere[1], Gary P. Wormser[3],

8    Franc Strle[7], Pardis C. Sabeti[1,2,10,11*], Ira Schwartz[3,*], and Klemen Strle[1,6,*]

9

10   [1]Massachusetts General Hospital, Harvard Medical School, [2]Broad Institute of MIT and Harvard,

11   [3]New York Medical College, [4]East Carolina University, [5]University of Ljubljana, [6]Wadsworth

12   Center, [7]University Medical Center Ljubljana, [8]University of Wisconsin, [9]Tufts University,

13   Department of Molecular Biology and Microbiology [10]Harvard University, [11]Harvard T.H.Chan

14   School of Public Health.

15   *Contributed equally to this work

16   Correspondence to: lemieux@broadinstitute.org

17   Key words: Lyme disease, *Borrelia burgdorferi*, Whole Genome Sequencing; Virulence;

18   Pathogenicity

19

20

21

**Abstract:** Lyme disease is the most common vector-borne disease in North America and Europe. The clinical manifestations of Lyme disease vary based on the genospecies of the infecting *Borrelia burgdorferi* spirochete, but the microbial genetic elements underlying these associations are not known. Here, we report the whole genome sequence (WGS) and analysis of 299 patient-derived *B. burgdorferi* sensu stricto (*Bbss*) isolates from patients in the Eastern and Midwestern US and Central Europe. We develop a WGS-based classification of *Bbss* isolates, confirm and extend the findings of previous single- and multi-locus typing systems, define the plasmid profiles of human-infectious *Bbss* isolates, annotate the core and strain-variable surface lipoproteome, and identify loci associated with disseminated infection. A core genome consisting of ~800 open reading frames and a core set of plasmids consisting of lp17, lp25, lp36, lp28-3, lp28-4, lp54, and cp26 are found in nearly all isolates. Strain-variable (accessory) plasmids and genes correlate strongly with phylogeny. Using genetic association study methods, we identify an accessory genome signature associated with dissemination and define the individual plasmids and genes that make up this signature. Strains within the RST1/WGS A subgroup, particularly a subset marked by the OspC type A genotype, are associated with increased rates of dissemination. OspC type A strains possess a unique constellation of strongly linked genetic changes including the presence of lp56 and lp28-1 plasmids and a cluster of genes that may contribute to their enhanced virulence compared to other genotypes. The patterns of OspC type A strains typify a broader paradigm across *Bbss* isolates, in which genetic structure is defined by correlated groups of strain-variable genes located predominantly on plasmids, particularly for expression of surface-exposed lipoproteins. These clusters of genes are inherited in blocks through strain-specific patterns of plasmid occupancy and are associated with the probability of invasive infection.

1

## INTRODUCTION

Lyme disease is a heterogeneous illness caused by spirochetes of the *Borrelia burgdorferi* sensu lato (*Bbsl, sensu lato meaning 'in the broad sense'*) complex. *Bbsl* contains over 20 ~~subspecies (also termed~~ genospecies~~, genomic species~~), four of which cause the majority of disease in humans: *B. burgdorferi* sensu stricto (*Bbss, sensu stricto meaning in the strict sense*), *B. afzelii*, *B. garinii*, and *B. bavariensis* *[1]*. Nearly all Lyme disease in the US is caused by *Bbss*. In Europe, most infections are caused by *B. afzelii*, *B. garinii*, or *B. bavariensis,* whereas infection due to *Bbss* is rare. Infection with *Bbsl* usually presents as an expanding skin rash, erythema migrans (EM), at the site of the tick-bite. If untreated, spirochetes may disseminate to secondary sites ~~(a phenotype described as 'dissemination')~~, primarily other skin sites, the nervous system and joints [1,2]. In addition to clinical variation among *Bbsl* species, differences in virulence have also been noted between genotypes within *Bbss* [3–5]*,* and such phenotypes have been recapitulated in murine models [6–8]. These associations imply that microbial genetic loci likely influence the clinical manifestations of Lyme disease. Despite such evidence linking microbial genotype to clinical phenotype, the specific genes or loci responsible for the clinical manifestations of Lyme disease have not yet been identified.

*Bbss* genome analysis has been limited to date due to technical challenges of sequencing and assembly and difficulties of obtaining isolates from cases of human disease. The *Bbss* genome consists of a roughly one megabase of core genome (consisting of a ~900Kb chromosome and the plasmids cp26 and lp54), as well as numerous (>15) additional circular and linear ~~extrachromosomal DNA elements (colloquially termed~~ plasmids) [9,10]. Subsets of plasmids have high levels of homology (as exemplified by seven 32 kilobase circular plasmids (cp32) [11] and four 28-kilobase linear plasmids (lp28) [10] in the B31 reference isolate), which have diversified through duplication, recombination, and other ~~primordial~~ evolutionary events [12], The sheer number of plasmids and their extreme homology has made sequencing and

2

70     assembly of complete *Bbss* genomes a major challenge, particularly with widely-used short read

71     sequencing methods [13].

72         The technical challenges of sequencing and assembly are compounded by the difficulty

73     of obtaining isolates from human disease. It ~~has been~~ is possible to culture the organism from EM

74     lesions in the majority of cases, but this requires a skin biopsy and specialized culture

75     techniques, both of which are rarely used in routine clinical practice. The organism has

76     occasionally been cultured from CSF in patients with meningitis, but ~~extremely~~ rarely from

77     synovial fluid in patients with Lyme arthritis, the most common late disease manifestation in the

78     US. Thus, the ~~great~~ majority of available *Bbss* isolates are from patients with EM, an early

79     disease manifestation. As a result of these challenges, only a small number of human clinical

80     isolates have been sequenced and analyzed.  To our knowledge, no large WGS studies of

81     human isolates have been conducted. Fewer than 50 human isolates analyzed by WGS have

82     been publicly reported, either sporadically or included in cohorts consisting primarily of tick-

83     derived isolates [14–19].

84         Genotyping systems have been developed to subclassify *Bbss* strains using single or

85     multiple genomic regions (reviewed in [20]). Two of the most commonly used typing methods

86     are based on a restriction-fragment length polymorphisms in the 16S-23S ribosomal RNA

87     spacer region [21,22], termed ribosomal spacer type (RST), and on sequence variation of outer

88     surface protein C (OspC), one of the most variable *Bbss* proteins [23,24]. RST typing subdivides

89     *Bbss* into 3 types, referred to as RST1, RST2 and RST3 [6], whereas OspC typing subdivides

90     *Bbss* into ~30 OspC genotypes of which >24 cause infection in humans [25–27]. RST and OspC

91     are in linkage disequilibrium on the core genome, and each RST genotype is generally

92     associated with particular OspC types (e.g., RST1 mostly corresponds to OspC types A and B

93     and RST2 corresponds primarily to OspC types F, H, K and N) [27]), whereas RST3 is the most

94     variable and correlates with the remaining OspC types. In addition to these genotyping

95   methods, multilocus sequence typing (MLST), which is based on eight chromosomal

96   housekeeping genes, has been used to further sub-stratify the strains [27,28]. According to the

97   *Borrelia* MLST database (https://pubmlst.org/borrelia/), >900 MLST sequence types have been

98   identified.

99        Application of targeted genotyping methods has previously established a link between

100  *Bbss* microbial genotype and several phenotypic properties including dissemination, disease

101  severity, immunogenicity, and distinct clinical presentation [1,3,5,6,8,26,27,29–32]. For

102  example, using RST and OspC genotyping we previously showed that RST1 OspC type A

103  strains have greater proclivity to disseminate, are more immunogenic, are associated with more

104  symptomatic early infection, and with a greater frequency of post-infectious Lyme arthritis.

105  However, these approaches lack the resolution to reconstruct a detailed evolutionary history or

106  to define individual genes or loci underlying phenotypic variability. The limitations of previous

107  studies have been further compounded by the absence of large cohorts of patient-derived

108  isolates accompanied by detailed clinical information. Here, we used whole genome sequencing

109  to characterize in detail the genomes – including the core genome and associated plasmids – of

110  299 patient-derived *Bbss* strains. The isolates were collected primarily from patients with EM,

111  ~~the initial skin lesion of the infection,~~ over three decades across Northeastern and Midwestern

112  US and Central Europe. We carried out phylogenetic and phylogeographic analysis, and

113  identified particular *Bbss* genomic groups, plasmids, and individual open reading frames (ORFs)

114  associated with ~~tissue invasive (~~disseminated~~)~~ human disease.

115

**MATERIALS and METHODS**

117  **Selection of *B. burgdorferi* isolates (see Supplemental Table 1)**. In total, 299 *Bbss* isolates

118  collected from 299 patients over a 30-year period (1992-2021) were included in this study: 202

119　from the Northeastern US, 61 from the Midwestern US and 36 from Slovenia (Central Europe).

120　The majority (97%) of isolates were derived from skin (n = 287) or blood (n = 2) of patients (9

121　were derived from ~~cerebrospinal fluid~~ [CSF]) by culturing in BSK or MKP medium [33,34]. All

Is this acromyn used again?

122　patients met the US Centers for Disease Control and Prevention (CDC) criteria for Lyme

123　disease [35]. Only low passage isolates (passage <5) were used for WGS.

124　*Northeastern United States:* The 201 isolates from the Northeastern US were collected at two

125　geographic locations: 113 from New England (primarily from contiguous regions of

126　Massachusetts, Rhode Island, and Connecticut) and 88 from New York State. The New York

127　strains belong to a larger collection of more than 400 clinical isolates, collected between 1992-

128　2005, that had been previously typed at the *rrs-rrlA* IGS and *ospC* loci [4,31]. To account for the

129　full diversity of *Bbss* genotypes found in the collection, isolates with the best sequence quality

130　from each OspC major group were selected for this study in accordance with their prevalence in

131　the entire collection. All of the latter isolates were cultured from skin biopsies of infected

132　patients, rather than from blood or CSF (Supplemental Tables 1 and 2).

133　*Midwestern United States:* The 62 isolates from the Midwestern US were derived from

134　specimens submitted to the Marshfield Laboratories (Marshfield, WI) for *Borrelia* culture from

135　1993 to 2003 (Supplemental Tables 1 and 2).

136　*Central Europe (Slovenia)*: The 36 isolates from Slovenia represent all *Bbss* isolates that were

137　cultured from patients over a 27-year period (1994-2021), who were evaluated at the Lyme

138　borreliosis outpatient clinic at the University Medical Center Ljubljana (UMCL).

139

140　**Selection of patients.** This study involves secondary use of deidentified archival clinical

141　isolates and patient data collected in previous studies and was approved by the Massachusetts

142    General Hospital Institutional Review Board (IRB) under protocol 2019P001864. Patients

143    included in this study were diagnosed with early Lyme disease and were classified as having

144    either localized or disseminated infection. Early Lyme disease was defined by the presence of at

145    least one EM skin lesion or symptoms consistent with Lyme neuroborreliosis along with a

146    positive CSF culture. Localized infection was defined by a single culture positive EM skin lesion

147    in the absence of clinical and/or microbiological evidence of dissemination to a secondary site.

148    Disseminated infection was defined by a positive blood or CSF culture or PCR, multiple EM

149    lesions, and/or signs of neurological involvement. We were able to classify 291 or the 299

150    (97.3%) isolates as Disseminated or Localized by these criteria. Clinical records were not

151    available to classify 8/299 (2.7%), and these isolates were excluded from analyses of

152    dissemination. A measure of bloodstream dissemination was available for 212/299 (70.9%) of

153    isolates, with blood PCR available for 106/299 (35.4%) and blood culture available for a disjoint

154    set of 106/299 (35.4%) of all isolates. Multiple EM was present in 57 / 290 (19.7%); among

155    patients with a single EM, 23/88 (26.1%) had a positive blood culture and 28/86 (32.6%) had a

156    positive PCR. Complications such as Lyme neuroborreliosis were defined by clinical criteria and

157    based on assessment by the treating clinician. In Europe, central nervous system (CNS)

158    pleocytosis and intrathecal production of *Borrelia* antibodies were required for diagnostic

159    determination of Lyme neuroborreliosis, following the EFNS guidelines [36]. Summary statistics

160    of isolates by group is provided in Supplemental Table 1. The list of isolates and associated

161    metadata is provided in Supplemental Table 2.

162    **Whole-Genome Sequencing**. *Bbss* DNA was isolated from the cultured isolates with either the

163    IsoQuick kit (Orca Research, Bothell, WA), the Gentra PureGene DNA Isolation Kit (Qiagen

164    Inc., Valencia, CA), or the DNEasy kit (Qiagen Inc, Valencia, CA). Short-read next-generation

165    sequencing (NGS) library construction was performed using the Nextera XT Library Prep Kit

166    (Illumina, San Diego, CA). DNA quantification was performed ~~in a 96-well microplate~~ using the

167    SpectraMax Quant dsDNA Assay Kit and the Gemini XPS Fluorometer (Molecular Devices, San

168    Jose, CA), or ~~in a single tube~~ using the Qubit 2.0 fluorometer (Thermo Fisher Scientific,

169    Springfield Township, NJ). Library quality was examined using the 4200 TapeStation and D1000

170    ScreenTape (Agilent, Santa Clara, CA). Paired-end sequencing (2 × 150 or 250 cycles) was

171    performed using the NextSeq 550 or MiSeq system (Illumina).

172    **Bioinformatics Data Analysis.** Trimmomatic v0.39 [37] was used for trimming and cleaning of

173    raw sequence reads; SPAdes v3.14.1 [38] for *de novo* genome assembly; QUAST [39] for

174    quality assessment and assembly visualization; Kraken2 [40] v2.1.1 for digital cleaning of

175    assembled genomic sequence by using taxonomy classification; mlst v2.19.0

176    (https://github.com/tseemann) for MLST [41] identification from assembled sequences; k-mer

177    weighted inner product (kWIP) [42] v0.2.0 for alignment-free, k-mer-based relatedness analysis;

178    prokka v1.14.6 [43] for genome sequence annotation; Roary [44] for core- and pan-genome

179    analysis; FastTree v2.1.11 [45] for phylogeny tree generation. Bioconductor [46] packages in R

180    [47] v4.1.1 and/or RStudio v2021.09.0+351, such as ggplot2 [48], ggtree [49], ggtreeExtra, and

181    ggstar, were also used for phylogeny tree generation. MLST definitions were downloaded from

182    pubMLST. Multidimensional scaling (MDS) was calculated on the kWIP distances using the

183    command mdscale() in R. Fisher's exact test was used for pairwise comparison of categorical

184    variables using the fisher.test() function in R. The MiniKraken2 database was constructed for

185    Kraken2 from complete bacterial, archaeal, and viral genomes in RefSeq as of March 12, 2020.

186    To characterize the plasmid content of individual isolates, we took two approaches. We first

187    aligned the contigs to the B31 reference and quantified a plasmid as present or absent if greater

188     than 50% of the reference genome plasmid was covered by contigs. As a complementary

189     approach, we built a hidden Markov model (HMM) of PFam32 genes using HMMer [50] and

190     searched the resulting profile against the assemblies to identify PFam32 genes. We then

191     aligned the resulting putative PFam32 genes against a set of canonical PFam32 genes,

192     provided by Dr. Sherwood Casjens, that have been used to determine plasmid types in

193     published reports [51]. For each putative PFam32 gene, if a match with <5% identity was

194     present in the list of annotated PFam32 genes, we marked the isolate as having a copy of the

195     closest-matching PFam32 based on sequence identity. If no PFam32 within these thresholds

196     could be identified, the closest PFam32 family member was considered unknown and not

197     assigned in this analysis.

198

199     **RESULTS**

200     *Whole-genome sequencing of human Borrelia burgdorferi* sensu stricto *isolates*

201     To gain insight into the evolution, population structure, and pathogenesis of *Bbss* in human

202     infection, we sequenced the complete genomes of 299 *Bbss* from human cases of early Lyme

203     disease. We sequenced their whole genomes at a median coverage of 57.6x (interquartile

204     range [IQR] 27.6x - 130.8x). The *de novo* assemblies produced high-quality, nearly-complete

205     genomic assemblies with a median total length of 1.34 megabases (Mb) (IQR 1.30 - 1.37 Mb).

206     Final assemblies contained a median of 107 contigs per isolate (IQR 88.0 - 137.5) and had a

207     median N50 of 213,476 bases (IQR 80,809 - 221,506 bases). We were unable to finish

208     assembly of plasmids due to repetitive plasmid sequences. Assembly statistics are given in

209     Supplemental Table 3.

210         As an initial characterization of divergence between strains without any reference or

211     annotation, we applied alignment-free, kmer-based analysis (kWIP) to the WGS data and

212    identified three major clusters based on their genetic distances (Figure 1C and D, Figure S1).

213    This unbiased distance analysis revealed that a single lineage (WGS A) was divergent from all

214    other isolates (Figure 1C and D). The remaining isolates are grouped into two stable clusters

215    (WGS groups B and C). RST type 1 was divergent from the other two WGS groups, but RST 2

216    and 3 were mixed between WGS groups B and C (Figures 1C and 1D).

217         We next constructed both maximum-likelihood (ML) and maximum clade credibility

218    (MCC) phylogenetic trees using core genome elements (as defined by Roary[44], see methods)

219    from WGS (Figure 2). WGS groups defined by k-mer distance corresponded to the ML clade

220    structure on the core-genome tree and the associated OspC types (Figure 2A and E). However,

221    they revealed substructure within these groups, particularly WGS group B, which we split into

222    subclusters B.1 and B.2 (Figure 2B and S3B). We also inferred MCC trees using Bayesian

223    methods as implemented in BEAST. A MCC tree is shown in Figure S2; ML and MCC trees

224    were in broad agreement, and the posterior probability of all nodes separating WGS groups was

225    > 0.99, indicating that the distance-based clustering was phylogenetically well-supported.

226

227    *Comparison of Bbss isolates using classical genotyping approaches*

228    We typed these isolates using the RST, OspC, and MLST typing schemes and compared WGS

229    type to these existing methods (Figures 1A and 1B). Among the 299 strains, 98 were RST1

230    (32.7%), 112 were RST2 (37.4%), and 89 (29.8%) were RST3; 52 (17%) were OspC type K, 44

231    (15%) were OspC type A, 46 (15%) were OspC type B, and 21 (7%) were OspC type H.  As

232    demonstrated previously [4,30], there was a strong linkage between RST and OspC type

233    (Fisher's exact test, $p < 1 \times 10^{-6}$).

234         In Slovenia in Europe, the most common isolates were RST1 (75%), >60% of which

235    were OspC type B. In contrast, the most common *Bbss* isolates in the US were RST2 (41% in

9

236    Northeastern US and 49% in Midwestern US), whereas RST1 strains comprised 32% of the

237    strains in Northeastern US and only 10% in the Midwest. Further, certain OspC types have

238    distinct geographic distributions. For example, OspC type L is found only in the Midwestern US

239    and Slovenia and OspC types Q, R and S have only been isolated from European patients

240    [26,27]. These findings are consistent with previous reports that found genetic differences in

241    *Bbss* populations based on geography [26,27]. WGS groups were strongly associated with RST

242    (Figures 1A-B, Fisher's exact test, p < 1 x 10$^{-6}$) and OspC type (Figure 1A-B, S1; Fisher's exact

243    test, p < 1 x 10$^{-6}$). RST1 / Osp C type A/B sequences consistently clustered as a single clade in

244    the core genome phylogenetic tree and MDS of k-mer distances (Figures 1C and S1),

245    demonstrating agreement between typing methods. In contrast, RST2 and RST3 were both

246    polyphyletic in the WGS data and contained within separate WGS groups (Figures 1C and 1D).

247    Trees inferred from core genome sequences (Figure 2D, left panel) differed in the relatedness

248    of major clades from those inferred from accessory genome sequences (as defined by Roary

249    [44], see methods) (Figure 2D, right panel), but agreed on the substructure and sample

250    membership of individual clades. This pattern, which affects major clades as a whole, indicates

251    the occurrence of recombination events deep in the evolutionary history between core and

252    accessory genome sequences.

253        Similarly, OspC types were monophyletic on the WGS tree (Figure 2E) and on a tree

254    built from OspC sequences (Figure 2F), but WGS type was polyphyletic on the OspC tree

255    (Figure 2G). Consistent with this polyphyly, face-to-face comparison of core genome and OspC

256    trees demonstrated that in many cases, closely related OspC sequences were part of distinct

257    WGS groups (Figure 2H). For example, the OspC type L isolates from the Midwestern US and

258    Slovenia are on different branches of the core genome phylogenetic tree (Figure S2H). Thus,

259    RST and OspC typing methods identify substructure in *Bbss* genomes, and largely agree on the

260     divergent RST1 / OspC A/B clade. In contrast, RST does not capture fine-grain genetic

261     structure, and OspC sequence distance does not correlate with genome-wide distance between

262     isolates.

263

264     *Population geographic structure:*

265     We next explored the relationship between genetic markers and geography. WGS group was

266     strongly associated with broad geographic region (US Northeast, US Midwest, EU Slovenia)

267     (Fisher's exact test, $p < 1 \times 10^{-6}$), similar to the findings with previously evaluated genetic

268     markers including RST (Fisher's exact test, $p < 1 \times 10^{-6}$) and OspC type  (Fisher's exact test, $p <$

269     $1 \times 10^{-6}$) (counts by geographic region are shown in Figures 1A-B).

270          Using finer-grained geographic clustering among subregions in the Northeastern US

271     (New York, Massachusetts, Connecticut, and Rhode Island), geographic region was significantly

272     associated with WGS group (Fisher's exact test, $p = 0.009$), suggesting that geographic

273     structuring of genotypes also occur on a regional scale (Figure S2E). The number of ORFs in

274     the genome differed significantly by region within a given WGS group (Figure 3A). In the US

275     Northeast and in Slovenia, WGS groups differed significantly by the number of ORFs (Figure

276     3B). As core genome size is relatively constant among strains regardless of geographic

277     location, the differences in accessory genome size across different populations, even within a

278     given genomic group with a single common ancestor, suggests that the diversification of

279     accessory genome size may be one mechanism by which strains adapt to distinct ecological

280     factors in each geographic region. Slovenian isolates are clustered in two well-defined

281     monophyletic groups (Figure 2C), suggesting at least two inter-continental exchanges (Figure

282     S2C), consistent with a previous report [15]. There were numerous (>10) exchanges between

283     samples in the US midwest and northeast (Figure S2D).

284      We attempted to define the timing of these exchanges by inferring a time-stamped

285    phylogeny using BEAST (Supplemental Note 1). Together, these models demonstrate a remote

286    (hundreds of thousands to tens of millions of years) TMRCA for human-infectious strains of

287    *Bbss*, consistent with previous estimates [52]. Precise timing requires more accurate knowledge

288    of the mutation rate in *Bbss.*

289

290    *Associations between genotype and Bbss dissemination in patients:*

291    Dissemination is a crucial clinical event that enables the progression of disease from an EM skin

292    lesion to more severe Lyme disease complications such as meningitis, carditis, and arthritis.

293    Given the previously-reported associations between single-locus genetic markers and

294    dissemination[4,5,8,30], we investigated the relationship between genotype and dissemination.

295    We scored isolates as either disseminated or localized based on certain clinical characteristics

296    of the patients from whom they were obtained, particularly having multiple vs 1 EM skin lesion

297    and having neurologic Lyme disease as well as having positive culture or PCR results for *Bbss*

298    in blood.

299      WGS groups differed from each other in their propensity to disseminate (p = 0.059 for 3

300    groups; p = 0.012 for 4 groups, Fisher's exact test) (Figure 3C, Figure S3C). Slovenian isolates

301    disseminated at a lower rate (25%) than US isolates (42.7%) (p = 0.045, Fisher's exact test),

302    and the relationship between WGS groups and dissemination was slightly stronger when testing

303    US isolates only (p = 0.02 for 3 groups; p = 0.004 for 4 groups, Fisher's exact test). WGS group

304    A isolates from the US, which correlate with OspC type A and RST1 strains, showed the highest

305    rate of dissemination (51.4%) whereas US WGS group B isolates had the lowest rate of

306    dissemination (32.4%).  Within WGS group B, there was evidence of substructure (Figure S3).

307    US B.1 isolates disseminated at a higher rate (40.0%) than B.2 isolates (18.4%) (Figure S3C).

308    Consistent with previous observations [3,4] and with the general alignment of WGS,

309    RST, and OspC type, RST type was also associated with dissemination (p = 0.010, Fisher's

310    exact test), with RST1 having the greatest propensity to disseminate and RST3 the lowest [4,5]

311    (Figure S4B). OspC type A was also associated with dissemination (p = 0.008, Fisher's exact

312    test, Figure S4A). A significant association with dissemination could not be detected when OspC

313    type was tested as a categorical variable with 23 categories (p = 0.3, Fisher's exact test, Figure

314    S4), but power is reduced by many categories.

315    The propensity to disseminate varied greatly among the US and Slovenian isolates,

316    which is likely due to the major genetic differences in isolates between the two regions (Figure

317    3C). In Slovenia, the predominant WGS group A isolates are OspC type B and all the WGS-B.2

318    isolates are ospC type L (Figure S4). This correlation was particularly notable for WGSA strains,

319    which were recovered from patients with disseminated Lyme disease at a rate of 51.4% in the

320    US vs 23.1% in Slovenia. WGS-B.2 isolates in the US possess the lowest dissemination rate

321    (18.4%), whereas those from Slovenia showed a higher dissemination rate of 30% (Figure 3D

322    and S4A). Taken together, these data confirm that rates of dissemination vary by genotype and

323    demonstrate that WGS A/RST1, particularly a subset distinguished by OspC type A strains, is a

324    genetically distinct lineage with higher rates of dissemination.

325

326    *Plasmid associations with WGS profiles:*

327    As most of the genetic variation in *Bbss* occurs on plasmids [51,53,54], we investigated the

328    variation in plasmid content across genotypes. Assembly and analysis of plasmid sequences is

329    challenging because the length of repeated sequences in plasmids is greater than the read

330    length generated by the short-read Illumina sequencing technology used in this study [13]. To

331    circumvent this, we exploited the relationship between plasmid partition genes (plasmid family

13

332    32; PFam32) and plasmid types [12,51], putatively identifying the presence or absence of a

333    plasmid by the presence/absence of unique PFam32 sequences (Figure 4). After annotating all

334    PFam32 genes in the assemblies using an HMM, we linked each putative PFam32 to a plasmid

335    by finding the closest match by sequence homology from a curated list of PFam32 protein

336    sequences (see methods).

337         Applying this method to each strain, we created a comprehensive map of plasmids

338    across *Bbss* strains (Figure 4A-B). While a few plasmids are found more broadly, distinct

339    genotypes and WGS groups contain unique constellations of plasmids. Several plasmids,

340    including cp26, lp54, lp36, lp25, lp28-4, lp28-3 are found in nearly all isolates (Figure 4A-B) and

341    others such as cp32-7, cp32-5, cp32-6, cp32-9, and cp32-3 are found in most strains. Other

342    plasmids were more variable and only found in certain genotypes. OspC type A strains

343    possessed a distinct plasmid profile, containing lp56 and a unique version of lp28-1 (marked by

344    the lp28-1 PFam32 as well as a previously-annotated "orphan" PFam32 sequence, BB_F13.

345    When found in isolation, BB_F13 defines an lp28-11 plasmid [51], so is annotated as such,

346    although in many cases it may signify a subtype of lp28-1 rather than an entirely new plasmid

347    (especially OspC type A isolates whose reference is likely similar to the B31 reference[9,10]).

348    Based on PFam32 sequences, WGS A strains also contained lp28-2 and most also contained

349    lp38. OspC type K strains also contained a relatively homogenous subset of plasmids including

350    lp21, lp28-5, lp28-6, cp32-12. WGS-A/ RST1 genotypes were the least heterogeneous with

351    respect to plasmid diversity and OspC type, whereas WGS-B and WGS-C groups (RST2 and

352    RST3) were more diverse, although the subset of RST2 strains consisting of OspC type K

353    isolates was also relatively homogenous. Curiously, lp28-9 was found only in Slovenian RST1

354    isolates (Figure 4), the majority of which were OspC type B (Figure 1); cp32-12, cp32-9, and

355    cp32-1 were also found more commonly in Slovenian isolates.

356     Many plasmids (e.g. lp28-1, lp28-2, lp38 and numerous others ) were found in multiple

357     distinct branches of the phylogenetic tree suggesting a complex inheritance pattern of

358     polyphyletic loss and/or recombination. This is consistent with the previously observed

359     reassortment between core genome elements and accessory genome elements (Figure 2D)

360     and genetic markers such as OspC (Figure 2H). For example, OspC types B and N both

361     contained lp28-8, whereas OspC type K genotype is most closely correlated with the lp21, lp28-

362     5 and cp32-12 pattern. lp56 is associated with OspC type A and OspC type I.

363     Specific plasmids showed significant associations with dissemination. The presence of

364     lp28-1 was associated with dissemination (OR 1.9, p = 0.01, Fisher's exact test), as was cp32-

365     11 (OR 1.9, p = 0.01) and cp32-4 (OR 2.0, p = 0.01) (Figure 4C-D, Supplemental Table 3). The

366     lp38 plasmid is present in roughly half of US isolates but absent in all Slovenian isolates and

367     demonstrated a trend toward being associated with dissemination (OR 1.6, p = 0.05).

368     To confirm the accuracy of these plasmid differences across genotype, we also

369     constructed a map of plasmid occupancy across strains by an alternate approach. We aligned

370     contigs from assembled genomes to the B31 reference sequence and annotated a plasmid as

371     "present" if the assembled contigs covered a majority of the reference plasmid sequence (Figure

372     S5A-C). Only plasmids present in the B31 reference genome were considered in this analysis.

373     These results were qualitatively similar to those obtained using the PFam32 sequences (Figure

374     S5, Supplemental Table 4) confirming that cp26, lp54, lp17, lp28-3, lp28-4 and lp36 were

375     present in nearly all strains whereas other plasmids were more variable.

376     Together, these analyses reveal a core set of plasmids present across *Bbss* strains as

377     well as strain-variable plasmids that are associated with distinct geographic and clinical features

378     (i.e., propensity to disseminate) of *Bbss*, suggesting that they contain individual genetic

379     elements that may underlie distinct disease phenotypes.

380

381    *Strain variation in core, accessory, and surface lipoproteome*

382    In an effort to implicate individual genetic elements in dissemination, the core and accessory

383    genome elements were identified in each of the sequenced isolates and all ORFs in the *de novo*

384    assemblies were annotated and clustered using BLAST, splitting clusters whose BLAST

385    homology was < 80% (Figure 5). Plotting the presence or absence of a given core or accessory

386    genome element adjacent to each isolate in the phylogeny reveals consistent patterns of ORF

387    presence/absence across closely related groups of isolates. Each of the genomic groups

388    contained unique clusters of ORFs in the accessory genome (Figure 5). The accessory genome

389    phylogenetic tree (Figure 2D, right) provided an alternative and more natural clustering of

390    accessory genome elements and PFam32 sequences (Figure S6A-B).

391        The most invasive genotype (WGS A) was associated with the largest pan-genome,

392    whereas the less invasive groups (WGS Group B and C) were associated with smaller genomes

393    (Figure 3A,B). Although many genes do not have a known function, we prioritized surface-

394    expressed lipoproteins (Figure 6) for further analysis because of their important roles in Lyme

395    disease pathogenesis and immunity (reviewed in [1,55]). We focused on the subset of all

396    lipoprotein ORFs demonstrated to be located on the surface of the spirochete [56] and divided

397    them into core (Figure 6A) and strain-variable (Figure 6B). The *Bbss* core lipoproteome (Figure

398    6A) consists of approximately 45 surface lipoprotein groups that are present in almost every

399    isolate. These include OspA and B, complement regulator acquiring surface proteins

400    (CRASPS), as well as several other lipoproteins whose functions are less well-understood. The

401    accessory lipoproteome (Figure 6B) consists of approximately 100 lipoprotein groups that are

402    strain-variable. These include lipoproteins found in only subsets of isolates, such as BB_A69

403    and BB_E31, and others, such as Decorin binding protein A (BB_A24) and OspC (BB_B19),

16

404     which were found in almost every isolate but broken into separate ortholog groups because of

405     extensive allelic diversity. Strain-specific clusters were also present in major gene families of

406     Erps[57,58] (Figure S7A) and Mlps[59,60] (Figure S7B). Larger numbers of these multi-gene

407     family members were found in more invasive WGS groups (A and C) (Figure 6C). The number

408     of lipoproteins in a given isolate was associated with the probability of dissemination ($\beta_1 = 0.037$

409     +/- 0.017, p = 0.03, logistic regression, Figure 7D). A stronger effect was seen for Erps ($\beta_1 =$

410     0.087 +/- 0.053, logistic regression, Figure 7D) with a trend toward significance (p = 0.1). In

411     contrast, the total number of ORFs and the number of Mlp alleles were not significant in logistic

412     regression models (p = 0.45 and p=0.38, respectively, Figure 7D). Aggregating mean effects by

413     OspC types (Figure S7E) showed similar trends.

414         Several lipoprotein groups, such as BBK32, BBK07, and BBK52 were found in almost all

415     strains, but were not found in a subset of closely related genotypes. Notably, CspZ (BBH_06)

416     and two other lipoproteins encoded on lp28-3, BB_H37 and BB_H32, were lost in two divergent

417     subsets of Slovenian isolates (Figure 6A), suggesting multiple independent loss events in

418     evolutionary history. Interestingly, these two subsets were either WGS-A or WGS-B.2, strains

419     with the greatest and least probability of dissemination (Figure S3). The increased frequency of

420     loss of lp28-3 in Slovenian isolates implies that this plasmid is likely non-essential for human

421     infection. Moreover, this finding suggests that the selective forces acting on lp28-3 may differ in

422     Europe and the US.

423         Many genes had evidence of recurrent loss or gain. For example, one cluster that shows

424     this pattern in Figure 5B contains the lipoproteins BB_J45, BB_J34, and BB_J36 along with 12

425     other genes annotated on the lp38 in B31, suggesting that these lipoproteins had been lost or

426     gained multiple times in the evolutionary tree as a part of a pattern that involved most or all of

427     lp38.

428

429 *Associations between Accessory Genome Elements, Genotype, and Dissemination*

430 The genetic basis of the phenotypic differences between these strains most likely includes

431 nucleotide-level variation in chromosomal and plasmid DNA as well as variation in gene

432 presence or absence in the accessory genome (which is primarily plasmid-borne). While it is not

433 feasible to resolve these associations definitively in this study, we attempted to identify

434 preliminary ORF-level associations by clustering ORFs according to homology using Roary [44].

435 We then applied linear mixed models genome-wide study approaches to identify ortholog

436 groups associated with disseminated infection (Figure 7A-B). We used the approach of Earle et.

437 al [61] to distinguish "locus" and "lineage" effects by identifying lineages that were associated

438 with a phenotype.

439 Two lineages, defined by principal components of the distance matrix between isolates,

440 were significantly associated with the phenotype of dissemination (MDS10, p = 0.02, Wald's

441 test), and a second component was borderline associated (MDS8, p = 0.08, Wald's test). The

442 results of all analyses are reported in Supplemental Table 5 and lipoprotein-specific analyses in

443 Supplemental Table 6. In ancestry-adjusted association logistic regression analysis in which

444 principal components were included as covariates [62], only a handful of loci were associated

445 with phenotype, and their genomic position was distributed throughout the genome with no

446 strong spatial pattern (Figure 7B). The uncorrected association statistics showed somewhat

447 stronger correlations that were concentrated in the plasmids (Figure 7A).

448 We also used the pan-genome association approach to identify associations between

449 ortholog groups and single-locus genetic markers. Single-locus genetic markers were strongly

450 linked to genetic variation in ORFs, particularly among plasmids (Figure 8; Supplemental Table

451 7 for OspC Type A; Supplemental Table 8 for OspC Type K; Supplemental Table 9 for RST1).

452    The strongest effects were seen among surface-exposed lipoproteins [56] (Figure S8).

453    Together, these results, along with those of Figure 6, demonstrate that individual *Bbss*

454    genotypes represent a tightly-linked set of genetic variation that confers a distinct surface

455    lipoproteome.

456          Due to the structural patterns of genetic diversity in *Bbss*, ORFs associated with

457    phenotype without ancestry correction (Figures 6D and Figure 7A) should not be ignored. Due

458    to the near-complete linkage (e.g. Figure 8) between genetic elements in the accessory

459    genome, individual loci with strong, causal effects on a given phenotype may not be separable

460    from their set of linked variants, i.e. their background lineage. OspC type A strains, which are

461    included among the strains with the highest rates of dissemination in this study (Figure S4) and

462    as reported previously [3,4], and which have been linked to more severe symptoms of Lyme

463    disease [3] (Figure S4C), are strongly associated with a set of approximately 75 loci (OR > 50)

464    including a DbpA ortholog group (OR 4964, $p = 1.9 \times 10^{-48}$, likelihood ratio test), an OspC

465    ortholog group (OR 2951, $p = 1.9 \times 10^{-48}$, likelihood ratio test), and BB_H26 (OR 2186, $p = 4.9 \times$

466    $10^{-38}$, likelihood ratio test). These and other linked alleles were strongly correlated with one

467    another ($r = 0.94$, $p < 2.2 \times 10^{-16}$ for DbpA/group1807 and OspC/group1021; $r = 0.85$, $p < 2.2 \times$

468    $10^{-16}$ for DbpA/group and BB_H26). In many cases this linkage is physical due to presence on

469    the same replicon (e.g. the BB_J alleles on lp38), strongly linked allelic groups may also be

470    present on distinct replicons (e.g. DbpA on lp54 and OspC on cp26). While the strong

471    correlations between individual alleles make it difficult to separate the statistical effects of

472    individual alleles, such correlations are also the characteristic and defining feature of *Bbss*

473    lineages.

474

475    **Discussion:**

476   The sequencing and analysis of 299 human clinical isolates of *Bbss* that we report here

477   provides a previously unavailable level of resolution into the *Bbss* genetic and geographic

478   diversity of *Bbss* strains causing Lyme disease. Our collection of WGS assemblies from these

479   isolates—which were collected across distinct geographic regions, and which were linked to

480   certain clinical manifestations, and systematically typed with RST, OspC, and MLST—lays a

481   foundation for further research and advances our understanding of Lyme disease in several

482   ways.

483   First, our results confirm and extend previous findings on the microbial genetic basis of

484   disease manifestations in humans. Prior studies have identified genetic markers and correlated

485   their presence with specific clinical findings [1,3,5,6,8,26,27,29–32], but the relationships among

486   these markers and specific *Bbss* genes that cause phenotypic differences had not yet been

487   studied due to limitations of existing typing systems and a lack of human isolates. Along with the

488   novel genetic diversity uncovered by sequencing additional clinical isolates, the statistical

489   evidence linking genetic elements to dissemination and geography that was observed in this

490   study will be useful in prioritizing candidate genes and/or loci for further experimental evaluation.

491   For example, we confirm here previous findings that WGS A / RST1 — particularly the subtype

492   defined by OspC type A — is genetically distinct [27,63–65], and we identify certain genetic

493   alterations associated with this lineage, including having a larger number of ORFs than other

494   lineages. These ORFs are found on a strain-specific constellation of plasmids, including lp28-1

495   and lp56. This is consistent with previous findings that have linked the presence of lp28-1 to

496   infectivity in mouse models [66–69]. Importantly, these results extend previous findings which

497   showed that RST1 OspC type A strains are associated with more severe Lyme disease [3], by

498   identifying candidate plasmids lp28-1 and lp56 as potential genetic factors associated with

499   greater virulence of these *Bb* genotypes in patients. We show that this association, derived from

500   mouse models, extends to humans.

501         Second, the microbial genetic association studies presented here begin to resolve the

502   individual genetic elements underlying certain human phenotypes of Lyme disease. Using two

503   different methods to infer the presence or absence of plasmids, we provide the first plasmid

504   presence / absence maps of a large collection of human clinical isolates. Integrating this

505   information with associations at the level of individual ORFs provides a clearer view of the

506   potential determinants of distinct phenotypes. While we cannot yet resolve the causative loci on

507   lp28-1 or lp56 that enhance the pathogenicity of OspC type A strains, we highlight candidate loci

508   and quantify the statistical evidence for each locus considered. ORFs in these plasmids such as

509   BB_Q67 (which encodes a restriction enzyme modification system [70,71]), BB_Q09, BB_Q05,

510   BB_Q06, BB_Q07, BB_J31, BB_J41 and others (Supplemental Table 8) are tightly linked to the

511   OspC type A genotype and are candidates for further experimental study.

512         In addition, our sub-analysis of surface-exposed lipoprotein sequences (Figures 6A and

513   6B) may also be useful for experimental follow-up given the importance of surface lipoproteins

514   for immunity, pathogenesis, and *Bbss*-host interactions (reviewed in [1,55]). Particular alleles of

515   DbpA (BB_A24), and specific members of the Erp (BB_M38, BB_L39) and Mlp (BB_Q35)

516   (supplemental data file 2, Figures 6C and 6D) families are associated with dissemination and

517   represent potential candidates for evaluation in follow-up studies. Both the specific list of ORFs

518   strongly associated with OspC type A and the more general pattern of variation across strains

519   provides clues into enhanced virulence. Among those ORF groups associated with the OspC

520   genotype, allelic variants of DpbA have been shown to promote dissemination and alter tissue

521   tropism in a mouse model of Lyme disease [72]. Multiple genes in linked blocks probably

522   contribute to pathogenesis. For example, In OspC type A strains, DbpA is strongly linked to

21

523    OspC type A. Allelic variation in OspC alters binding to extracellular matrix components,

524    promotes joint invasion, and modulates joint colonization[73]; OspC has also been shown to

525    promote resistance in serum killing assays [74], and its role in causing infection can be, under

526    certain circumstances, partially complemented by other surface lipoproteins [75,76].

527        Our data suggest that copy number among multi-copy gene families may be linked to

528    dissemination. Given that Erps are divided into three families that each bind to distinct host

529    components (extracellular matrix, complement component, or complement regulatory protein)

530    [58,77–80]; it is possible that the strain-variable clusters of Erps (Figure S7B, Figure S7D-E)

531    may influence clinical manifestations by modulating strain-specific properties of tissue adhesion

532    or resistance to complement-mediated killing of spirochetes. The functions of Mlp proteins  and

533    many other strain-variable lipoproteins are still not well understood. The statistically-significant

534    relationship between lipoprotein number and probability of dissemination and the borderline-

535    significant relationships for copy number of Erps and Mlps (Figure S7D-E) suggest that varying

536    the amount and diversity of linked clusters of surface lipoproteins—which, individually or in

537    combination, may promote survival in the presence of immune defenses, binding to mammalian

538    host tissues and other pathogenic mechanisms— may be a general mechanism for strain-

539    specific virulence of *Bbss*.

540        Using unadjusted, univariate associate models, virtually all dissemination-associated

541    genes were found on plasmids. However, after correction for spirochete genetic structure due to

542    lineage, only weak locus-specific associations were observed. The block patterns of Figures 5

543    and 6 demonstrate why this is the case. Genes are inherited in blocks; the inheritance pattern of

544    genes within these blocks is strongly correlated such that only infrequently are genes from

545    within a block found in isolates that are outside the block. This pattern is also seen in plasmids,

546    and plasmids are a natural mechanism for this pattern of inheritance. An important

547    consequence of this finding is that it may not be possible to resolve individual loci beyond

548    correlated blocks of genes simply by increasing the number of samples or other methods to

549    improve statistical power because the near-complete correlation between individual loci makes

550    it statistically difficult to distinguish the individual effects among correlated genes. Thus, beyond

551    identifying genomic elements or groups of correlated genes associated with a phenotype,

552    further fine mapping will require biological experiments with reverse genetic tools. The results

553    shown in Figure 6 and 7, and Supplemental Data File 3 are helpful in narrowing down the

554    candidate loci and genetic elements that may predispose to or protect from dissemination.

555         Third, our analysis highlights how evolutionary history, geography, and differences in

556    strain genetic diversity interact in complex ways to contribute to clinical heterogeneity in Lyme

557    disease. In the context of known associations between genotype and clinical disease, the

558    difference in genetic markers across geographic areas may help explain why some clinical

559    phenotypes are more common in certain geographic locations. For example, Lyme arthritis is

560    more common in the US compared to Europe, probably because the infection in the US is due

561    predominantly to *Bb*ss strains which are more arthritogenic [81]). OspC type A strains appear to

562    be more common among patients in the US Northeast [26,27]. The intermixing of WGS groups

563    B and C in RST types 2 and 3 has not been a major issue in practice because the phenotypes

564    (for example, the relative rate of dissemination) of those groups appear more similar than the

565    genomically and phenotypically divergent RST1 / WGS A group. Similarly, OspC genotyping

566    has its limitations. The large number of OspC types (at least 30) makes phenotypic associations

567    with specific OspC genotypes challenging. More importantly, the discordance between OspC

568    sequences and whole-genome phylogenies—a discrepancy observed since the earliest OspC

569    sequences were published [82] and likely related to the fact that the OspC locus is a known

570    recombination hotspot on cp26 [83]—may make OspC unreliable as a genetic marker of

23

571    phenotypic traits. In this regard, WGS serves as a gold standard against which other typing

572    methods can be compared, facilitated here by our sequenced and fully-typed set of isolates.

573        WGS also offers new insight into evolutionary history and population divergence of *Bbss.*

574    Estimation of divergence times suggests a remote (at least hundreds of thousands of years)

575    origin for human infectious *Bbss.* The similarity in TMRCA estimates for samples from Slovenia,

576    the US Northeast, and the US Midwest indicate that the common ancestry for sequences

577    currently circulating in these populations is also remote; however, the strong lineage structure

578    and history of multiple exchanges suggests that the local history of distinct lineages is also

579    complex likely with multiple inter-region migration events. The consistent directional differences

580    in ORF number by region also suggest that adaptive evolution to local environments has

581    occurred, exploiting mechanisms of gene loss/gain on plasmids.

582        This report has several limitations. First, plasmids pose a unique challenge for assembly

583    and annotation [10,12]. As others have shown [13], complete plasmid assembly with short read

584    sequences is not possible. We devised two bioinformatic methods to overcome these changes

585    and infer plasmid presence/absence from short read sequencing, but neither is perfect. Our

586    PFam32 analysis is limited by an uncertainty as to which gene sequences are contained on the

587    plasmid associated with the PFam32 sequence. A complementary analysis based on the B31

588    reference sequence relies on a high-quality pre-existing assembly but cannot account for

589    genes/plasmids absent from the B31 reference. We also cannot exclude the possibility of

590    plasmid loss during culture, but isolates were passaged fewer than five times to minimize this

591    possibility.

592        Second, there are limitations due to analysis of isolates collected over time by different

593    groups at different sites. In particular, we may underestimate dissemination because an

594    assessment of spirochetemia (blood PCR or blood culture) was only available for 70.9% of

595    isolates (supplemental data file 3)  and the absence of positive culture or blood PCR from a

596    single time point does not rule out the possibility that dissemination from the initial skin lesion

597    may have occurred or may occur at a later time in untreated patients.

598         Third, there are statistical limitations related to the *Bbss* genome and study size. Models

599    that naively correlate a given gene with the phenotype of interest will produce spurious

600    associations due to the confounding effect of lineage and may overstate the effect from single

601    loci, a problem which is well known in human genome-wide association studies [84]. Corrections

602    for lineage and population structure are often applied to human [85,86] and bacterial [61,62]

603    association studies. However, *B. burgdorferi* underscores the challenges to these approaches,

604    both because lineages appear to be *defined* by the exchange of blocks of genes and because

605    the coarse tree structure differs for the core and accessory genomes, implying that a single

606    similarity measure to capture the pairwise dissimilarity between strains may not be adequate.

607    Larger studies with more isolates, statistical methods that incorporate the joint distribution

608    between genetic markers, and plasmid assemblies finished by long read sequencing are

609    required as a next step. The present study includes isolates collected by different investigators

610    over the past 30 years. Due to the logistical complexity and cost of collecting *Bbss* isolates from

611    patients in clinical studies, substantially larger studies of *Bbss* from patients may not be feasible

612    in the near term; however, long-read sequencing approaches have improved in accuracy,

613    availability, and cost, making finishing the genomes of existing isolates a logical next step.

614         Taken together, our results indicate that each *Bbss* genotype represents a tightly-linked

615    constellation of strain-specific variation that occurs primarily in plasmids, much of it involving

616    surface-exposed lipoproteins. OspC type A strains—with their enlarged pan-genome, distinct

617    set of plasmids, including lp28-1 and lp56, and variants of many surface lipoproteins, particularly

618    a unique subtype of DbpA—represent the most dramatic example of this genetic signature

619  associated with distinct phenotypes of Lyme disease in humans. Nevertheless, the pattern is

620  generalizable across genotypes and, given the strong linkage between microbial genotype and

621  phenotype for *Bbsl,* and similarities in genetic structure among *Bbsl* genomes, is likely true

622  broadly across all Lyme disease agents (*Bbsl*).

623

651

652

653

654
655 **Figure 1: A.** Counts of samples according to RST and OspC type. Top, middle, and lower
656 panels show samples from different geographic regions. X-axis gives OspC type. Bars are
657 colored according to RST type. **B.** Plots as in (A) but with bars colored according to the WGS
658 group. **C**. Multidimensional scaling (MDS) of 299 *Bbss* genomes, with WGS RST type
659 annotated. **D**. MDS of 299 *Bbss* genomes, with WGS type annotated.

660



**Figure 2: A-B.** Core genome phylogenetic tree with tips labeled by the three major WGS groups (A). **B.** Core genome phylogeny with WGS group B split into subgroups (B.1 and B.2). **C.** Core genome phylogenetic tree with tips labeled by region of collection. **D.** The core genome phylogenetic tree (left) compared to the accessory genome phylogenetic tree (right). Lines, colored by WGS groups, connect tips from identical samples. **E.** WGS tree with tips colored by OspC type. **F.** OspC tree with tips colored by OspC type. **G.** OspC tree with tips colored by WGS group. **H.** WGS tree (left) and OspC tree with identical tips connected by strain lines, colored by OspC type.

670
671 **Figure 3: A.** Number of ORFs by geographic region in different WGS groups. * denotes p <
672 0.05; ** denotes p < 0.01; *** denotes p < 0.001; **** denotes p < 0.0001; ns - not significant. **B.**
673 Number of ORFs by WGS group in different geographic regions. **C.**Probability of dissemination
674 by genomic group. Each point represents a sample. Points are colored by WGS group. The
675 samples that disseminated have been plotted at y = 1; those that did not have been plotted at y
676 = 0. Random noise has been added to the x- and y- coordinate to display the points. The mean
677 +/- 95% binomial confidence interval is shown for each group with error bars.

678



679

680

681 **Figure 4: A.** Core genome maximum likelihood phylogeny with tips colored by OspC type. The
682 clade corresponding to RST1 is shaded in light blue and the clade corresponding to OspC type
683 A is shaded in green. **B.** The presence/absence matrix at the right shows the presence or
684 absence of individual plasmids using the presence or absence of Pfam32 plasmid-compatibility
685 genes as a proxy. The columns of the matrix have been clustered using hierarchical clustering.
686 The rows of the matrix are ordered according to the midpoint rooted maximum likelihood
687 phylogeny shown at left. **C.** Odds ratio of dissemination and confidence interval by plasmid,
688 inferred by Pfam32 sequences. **D.** Volcano plot displaying the -log10 P value (as calculated

689    using Fisher's exact test) and the odds ratio of dissemination for each plasmid, inferred by
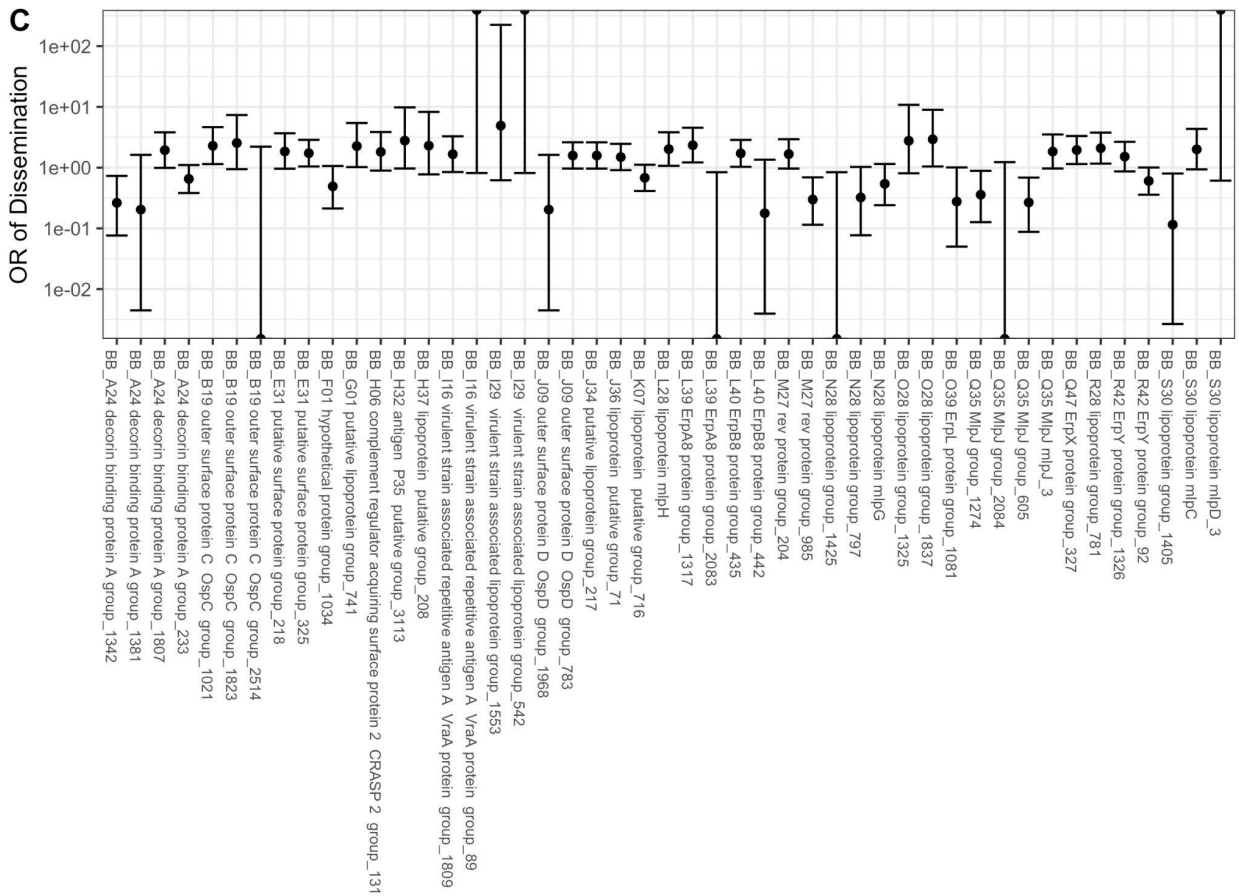690    Pfam32 sequences.

691



692

693

694 **Figure 5: A.** Core genome phylogeny with tips colored by OspC type. **B.** The phylogeny is
695 plotted alongside a matrix of presence (blue) or absence (white) for genes in the accessory
696 genome. The rows of the matrix are ordered by the phylogenetic tree in **A**. The columns of the
697 matrix are ordered using hierarchical clustering such that genes with similar patterns of
698 presence/absence across the sequenced isolates are grouped close together. **C.** Odds ratio
699 (OR) of dissemination and 95% confidence interval for ortholog groups encoding surface-
700 exposed lipoproteins and for which the unadjusted p-value for association with dissemination
701 (by Fisher's exact test) is < 0.15.

702 **A**

703
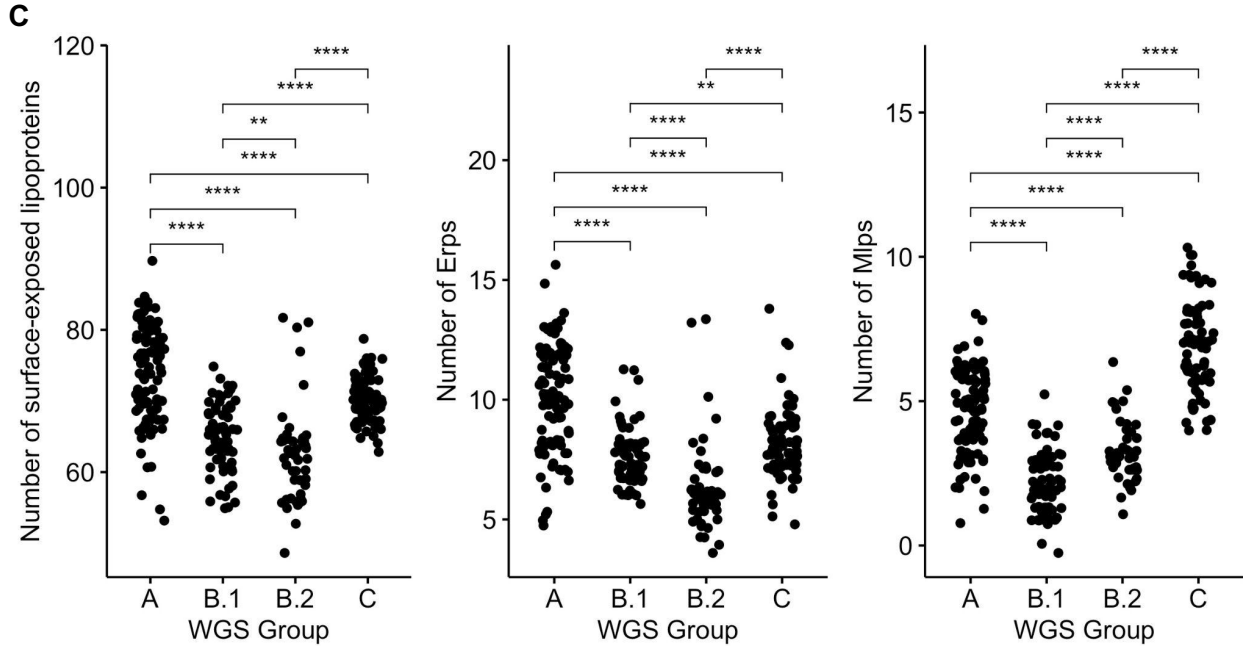
**B**

BB_K01 lipoprotein, putative group_543
BB_N28 lipoprotein group_2012
BB_M38 ErpK protein group_1372
BB_N28 lipoprotein group_1373
BB_B19 outer surface protein C (OspC) group_1044
BB_R42 ErpY protein group_1367
BB_O28 lipoprotein group_1837
BB_Q35 MlpJ group_1057
BB_N28 lipoprotein group_797
BB_Q35 MlpJ group_1293
BB_B19 outer surface protein C (OspC) group_1687
BB_L40 ErpB8 protein group_445
BB_N38 ErpP protein group_590
BB_N38 ErpP protein group_1048
BB_Q35 MlpJ group_1274
BB_C10 rev protein group_178
BB_L28 lipoprotein mlpH
BB_L39 ErpA8 protein group_1317
BB_S30 lipoprotein group_1027
BB_R42 ErpY protein group_92
BB_Q35 MlpJ group_1333
BB_S30 lipoprotein group_1332
BB_B19 outer surface protein C (OspC) group_1823
BB_L40 ErpB8 protein group_752
BB_M38 ErpK protein group_986
BB_R42 ErpY protein group_1024
BB_N28 lipoprotein mlpI_3
BB_F01 hypothetical protein group_768
BB_Q35 MlpJ group_757
BB_Q47 ErpX protein group_327
BB_O40 ErpM protein group_143
BB_O39 ErpL protein group_59
BB_B19 outer surface protein C (OspC) group_742
BB_I39 putative surface antigen group_991
BB_L28 lipoprotein mlpH_1
BB_N38 ErpP protein group_1316
BB_L40 ErpB8 protein group_435
BB_K07 lipoprotein, putative group_716
BB_M27 rev protein group_204
BB_N28 lipoprotein mlpI_2
BB_Q05 antigen, P35, putative group_169
BB_J34 putative lipoprotein group_584
BB_Q89 hypothetical protein group_981
BB_A24 decorin-binding protein A dbpA
BB_S41 ErpG protein group_90
BB_A69 putative surface protein group_102
BB_E31 putative surface protein group_218
BB_Q35 MlpJ group_771
BB_B19 outer surface protein C (OspC) ospC
BB_Q35 MlpJ group_769
BB_O39 ErpL protein group_434
BB_L40 ErpB8 protein group_250
BB_N38 ErpP protein group_767
BB_L28 lipoprotein mlpH_2
BB_Q35 MlpJ mlpJ_1
BB_M27 rev protein group_770
BB_P28 surface protein, mlp lipoprotein family mlpA
BB_Q35 MlpJ mlpJ_3
BB_O39 ErpL protein group_203
BB_M28 lipoprotein mlpD
BB_Q35 MlpJ group_1609
BB_O40 ErpM protein group_734
BB_N28 lipoprotein mlpG
BB_A24 decorin-binding protein A group_233
BB_S30 lipoprotein mlpC
BB_M27 rev protein group_961
BB_E31 putative surface protein group_325
BB_J36 lipoprotein, putative group_71
BB_J34 putative lipoprotein group_217
BB_J09 outer surface protein D (OspD) group_783
BB_N38 ErpP protein group_1329
BB_Q35 MlpJ mlpJ_2
BB_Q05 antigen, P35, putative group_1327
BB_I16 virulent strain-associated repetitive antigen A (VraA protein) group_1809
BB_B19 outer surface protein C (OspC) group_1021
BB_A24 decorin-binding protein A group_1807
BB_S41 ErpG protein group_133
BB_M38 ErpK protein group_569
BB_I38 putative surface antigen group_159
BB_R42 ErpY protein group_1326
BB_R28 lipoprotein group_781
BB_K48 immunogenic protein P37, putative group_714
BB_K50 immunogenic protein P37 group_194
BB_E31 putative surface protein group_1266
BB_J41 putative antigen P35 group_1659
BB_N38 ErpP protein group_409
BB_K48 immunogenic protein P37, putative group_116
BB_K52 putative lipoprotein group_1270
BB_N28 lipoprotein group_1253
BB_N39 ErpQ protein group_565
BB_S30 lipoprotein group_1405
BB_N28 lipoprotein mlpI
BB_O39 ErpL protein group_1081
BB_M27 rev protein group_985
BB_A24 decorin-binding protein A group_1342
BB_Q35 MlpJ group_605
BB_O28 lipoprotein group_1063
BB_B19 outer surface protein C (OspC) group_1348
BB_F01 hypothetical protein group_1034
BB_O28 lipoprotein group_1035
BB_B19 outer surface protein C (OspC) group_1015
BB_S30 lipoprotein group_1324
BB_O39 ErpL protein group_2015
BB_Q35 MlpJ group_1018
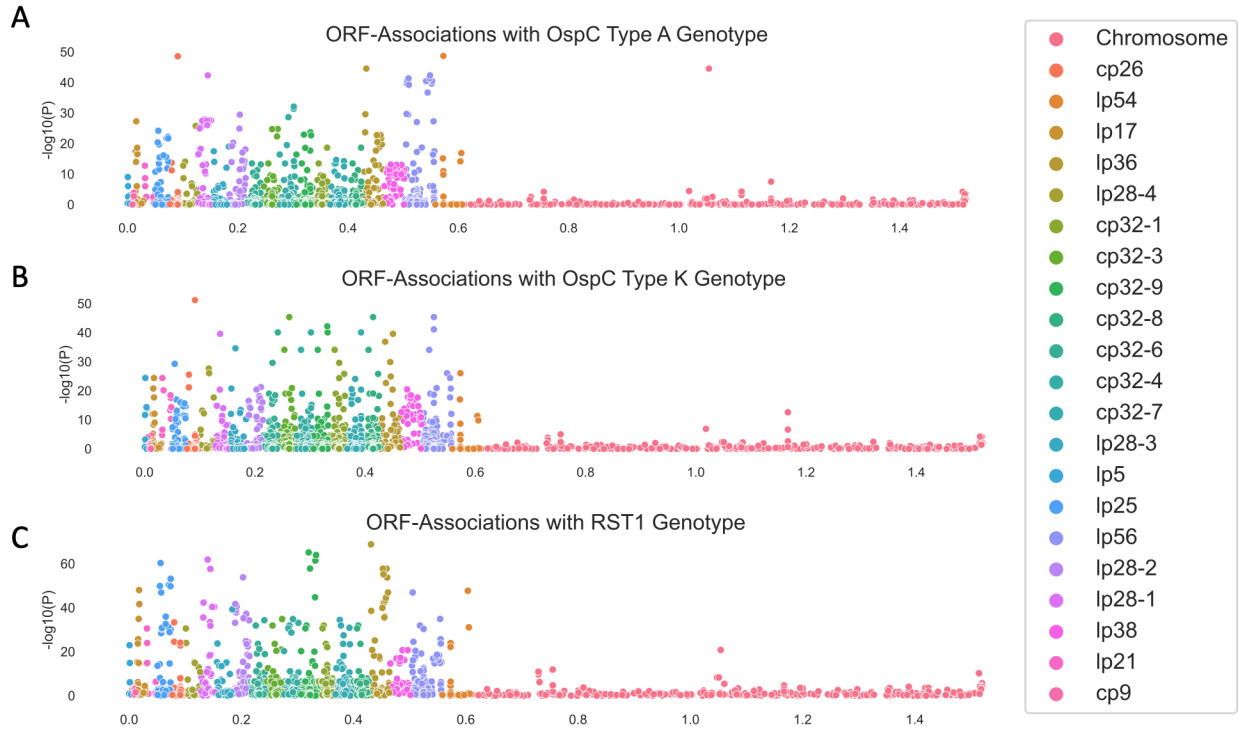BB_O40 ErpM protein group_1323

705
706
707

708 **C**



709
710 **Figure 6 : A.** *Bbss* core surface lipoproteome: Core genome phylogeny with tips colored by
711 OspC type (colored according to the scheme in Figure 5) with a matrix of presence (blue) or
712 absence (white) for surface lipoproteins. Surface-exposed lipoproteins present in at least 80% of
713 strains were considered to be part of the core lipoproteome. **B.** *Bbss* strain-variable (accessory)
714 surface lipoproteome: Core genome phylogeny with tips colored by OspC type (colored
715 according to the scheme in Figure 5) with a matrix of presence (blue) or absence (white) for
716 surface lipoproteins. Surface-exposed lipoproteins present in between 5% and 80% of strains
717 were considered to be part of the strain-variable (accessory) lipoproteome. **C.** The number of
718 surface-exposed lipoproteins (left panel), Erps (middle panel), and Mlps (right panel) by WGS
719 group. ** denotes p < 0.01; *** denotes p < 0.001; **** denotes p < 0.0001.

720



A

ORF-Associations with Dissemination, no Lineage Correction

B

ORF-Associations with Dissemination, with Lineage Correction

C

ORF-Associations with Dissemination, by Lineage

721

D

722

723 **Figure 7:** Manhattan Plots showing the association of individual ORF ortholog groups with the
724 phenotype of dissemination. **A.** P-values from univariate logistic regression by genomic position
725 for each ORF. **B.** P-values from regression estimates that include principal components
726 distance matrix between strains. **C.** Manhattan plot showing loci associated with each lineage
727 for the lineages associated with phenotype. **D.** Odds ratios (OR) (exp(beta)) with 95%
728 confidence interval are shown for dissemination for the lineage-adjusted model. ORFs with p <
729 0.1 and allele frequency > 0.1 and < 0.9 are displayed.
730
731

732
733
734 **Figure 8:** Manhattan Plots showing the association of individual ORF ortholog groups with
735 OspC type A (panel **A**), Osp C type K (panel **B**), and RST1 (panel **C**).
736

737     Supplemental Figures:

738



739
740     **Supplemental Figure 1:** Multidimensional scaling (MDS) reveals the population structure of US
741     and Slovenian *Bbss* isolates.

744    **H**



**Supplemental Figure 2: A.** Maximum clade credibility (MCC) tree. Nodes with posterior
probability > 0.9 are colored. **B.** Maximum likelihood (left panel) and MCC tree, with identical
tips connected with lines colored according to WGS group. **C.** MCC tree with nodes with
posterior probability > 0.9 labeled. Tips from the US have been grouped and their most recent
common ancestor are colored blue; all others are colored red. **D.** MCC tree with nodes with
posterior probability > 0.9 labeled. Tips from outside the US Midwest have been grouped  and
their most recent common ancestor are colored blue; all others are colored red. **E.** Time-tree
with 95% credible interval of node heights plotted as gray bars. **F.** Density of time to most recent
common ancestry (TMRCA) for major subpopulations and the full sample set (root). An inset
boxplot gives the median and IQR. **G.** Density of time to most recent common ancestry
(TMRCA) for major subpopulations and the full sample set (root) under three different fixed-
clock models with the clock rate set at $1\times10^{-10}$ substitutions/site/yr (left panel), $1\times10^{-9}$
substitutions/site/yr (middle panel), or $1\times10^{-8}$ substitutions/site/yr (right panel). **H.** Core genome
phylogeny of 299 whole-genome sequences. The phylogeny is shown as a cladeogram (branch
length does not correspond to genetic distance). The tips are labeled with sample names. RST
type, OspC type, location, and MLST type are annotated. Whole genome sequences
recapitulates existing typing schemes while adding additional resolution. Geographic origin is
associated with different branches of the tree. For example, Slovenian isolates cluster in two
distinct branches.

766    **A**



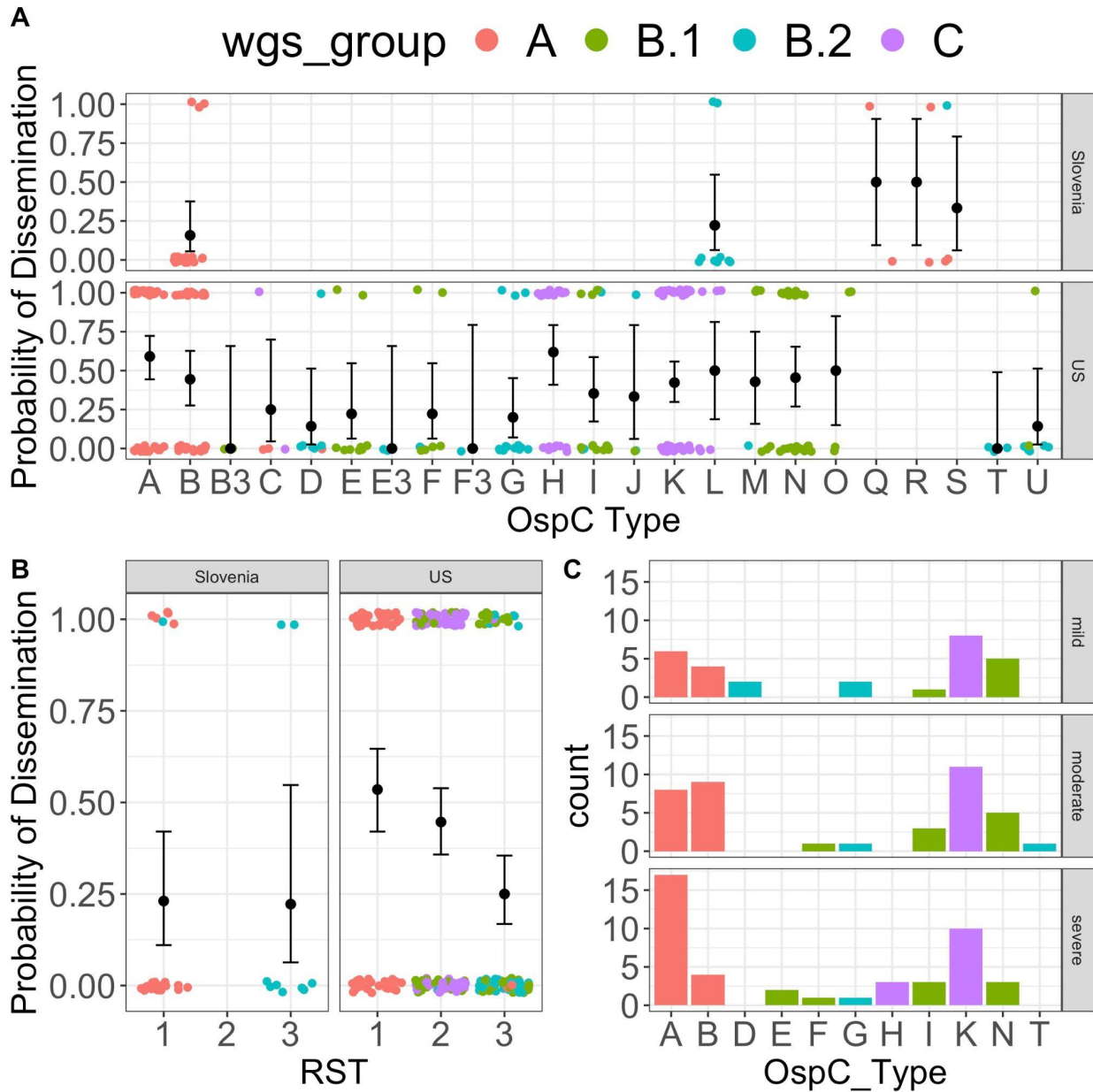767
768



769
770    **Supplemental Figure 3: A.** Core genome phylogenetic tree colored by WGS groups A-C with
771    group B divided into B.1 and B2; accessory genome presence/absence matrix is reproduced
772    from Figure 5 to highlight accessory genome elements that correlate with B.1 and B.2
773    sublineages. The clade corresponding to RST1 is shaded in light blue and the clade
774    corresponding to OspC type A is shaded in green. **B.** MDS plot with group B divided into B.1
775    and B.2. **C.** Probability of dissemination by genomic group using the four groups including B.1
776    and B.2.
777

**Supplemental Figure 4:** Probability of dissemination by **(A)** OspC type and **(B)** RST. **C.** Severity of Lyme disease by OspC type with WGS group shown by color.

784   **A**



785



786
787
788   **Supplemental Figure 5: A.** Inferred presence / absence of a plasmid based on alignment of
789   assembly contigs to the B31 reference. A plasmid is inferred as 'present' in the isolate if > 50%
790   of the length is covered by aligned contigs in the de novo assembly for the genome of the
791   corresponding isolate. The clade corresponding to RST1 is shaded in light blue and the clade
792   corresponding to OspC type A is shaded in green. **B.** Odds ratio of dissemination and
793   confidence interval by plasmid, inferred by Pfam32 sequences. **C.** Volcano plot displaying the -

794      log10 P value (as calculated using Fisher's exact test) and the odds ratio of dissemination for
795      each plasmid, inferred by alignment of assembled contigs to the B31 reference sequence.
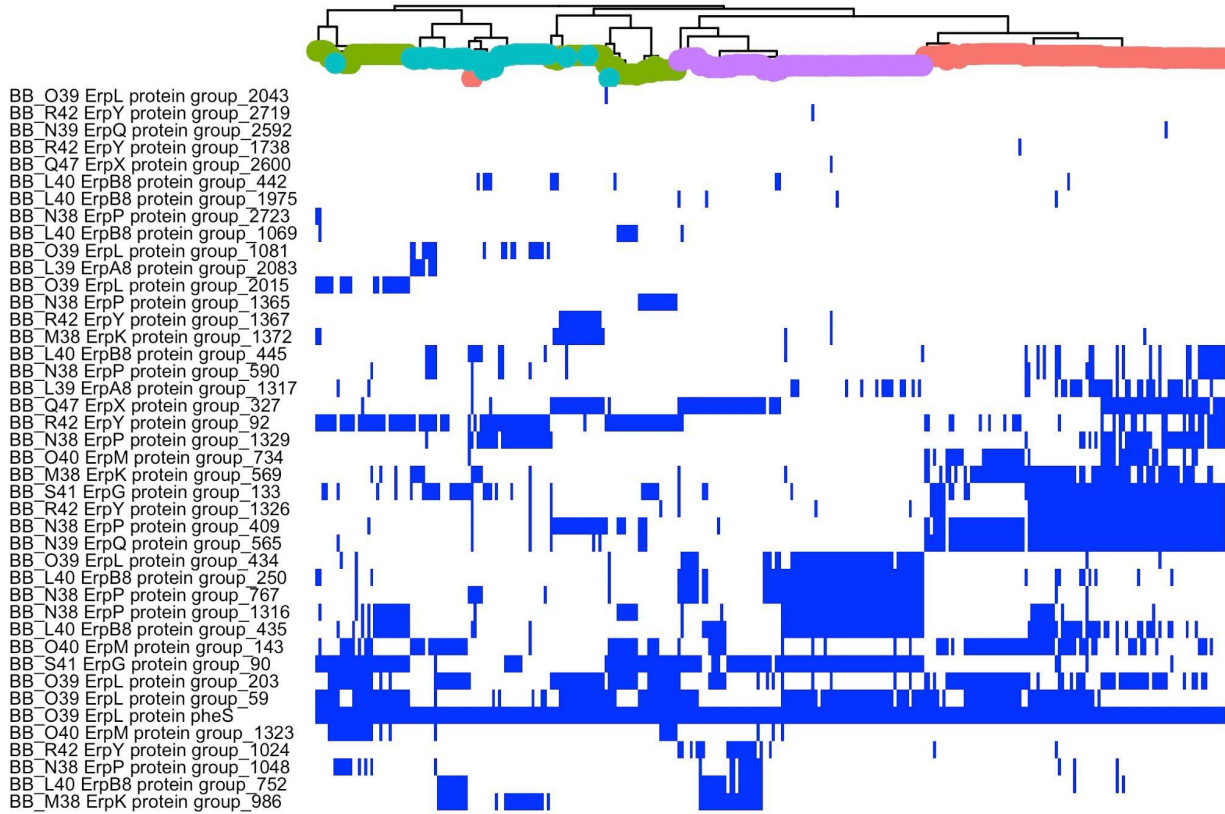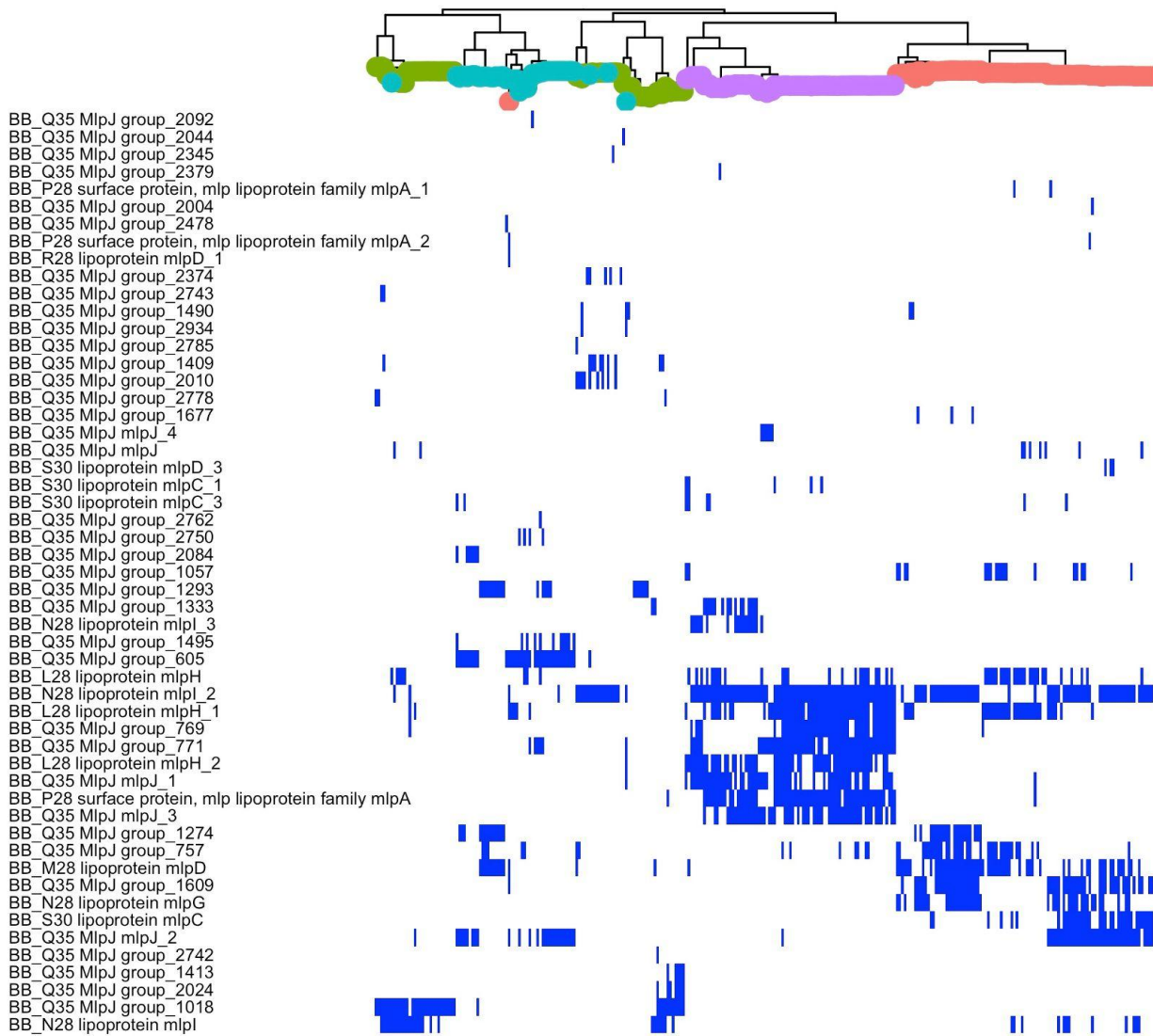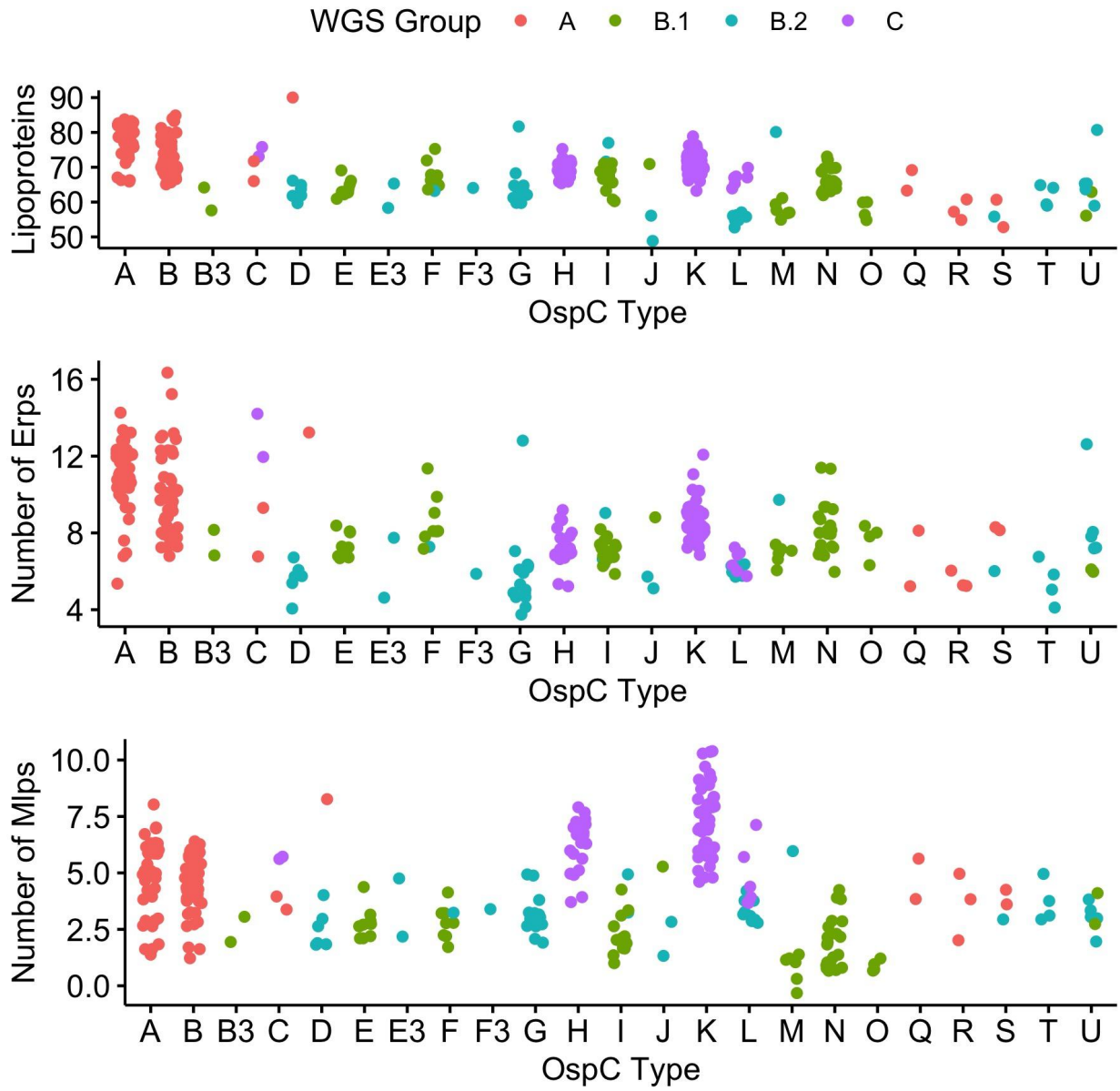
799 **B**



800
801 **Supplemental Figure 6: A.** Phylogenetic tree created from the accessory genome with
802 accessory genome elements plotted according to their presence/absence in individual strains.
803 **B.** Phylogenetic tree created from the accessory genome with PFam32 plasmid compatibility
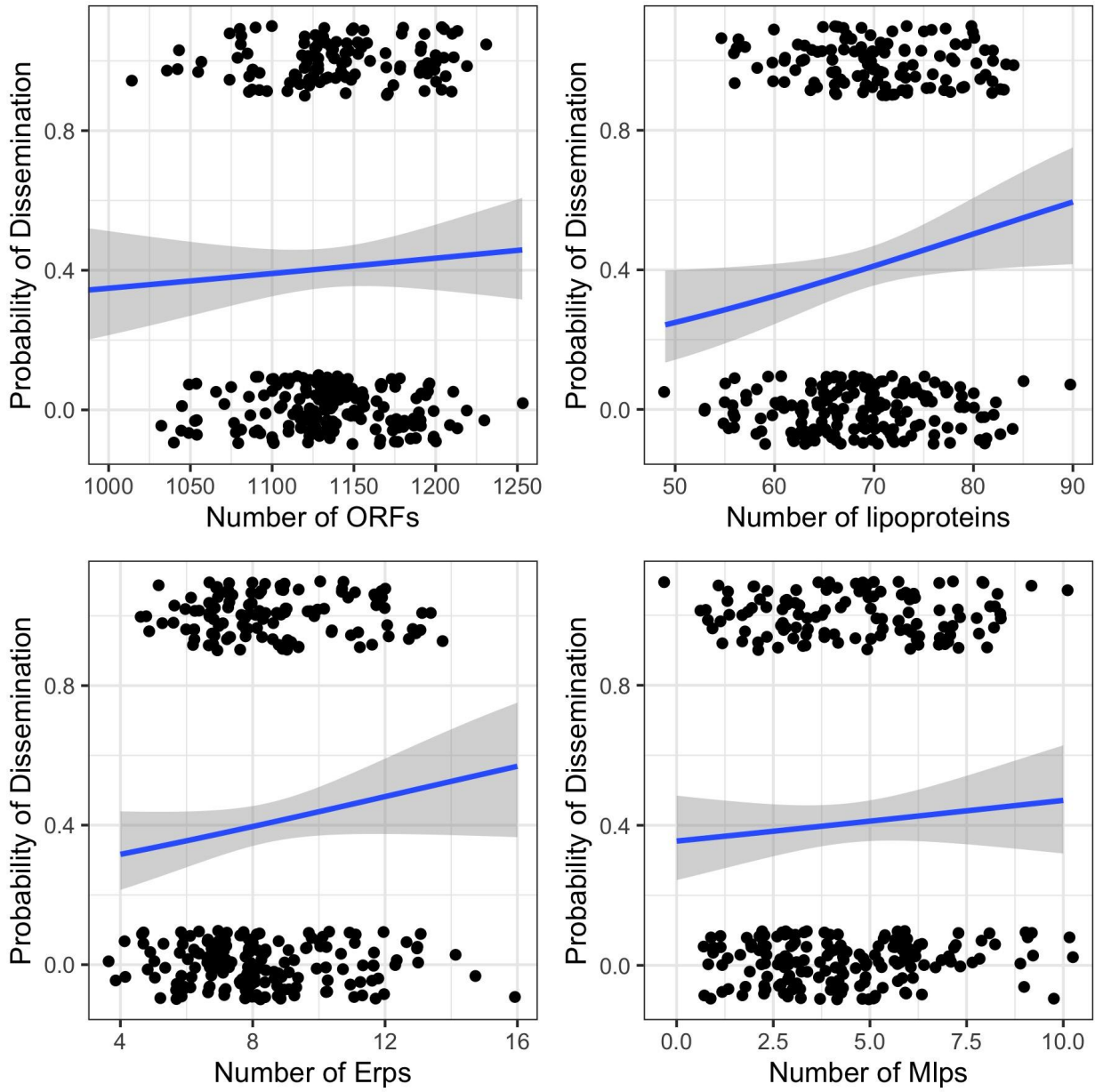804 sequences plotted according to the presence/absence in individual strains.
805

49

**A**

808　**B**



WGS Group
A ●
B.1 ●
B.2 ●
C ●

Surface Lipoprotein
Present ■
Absent

809
810

51

811    C

**D**

817　**E**
818



819
820　**Supplemental Figure 7: A** and **B.** Core genome phylogeny with presence/absence of Erp (C)
821　orthologs and Mlp (D) orthologs. **C.** The number of surface-exposed lipoproteins (top panel),
822　Erps (middle panel), and Mlps (bottom panel) by OspC type. **D.** Probability of dissemination by
823　number of ORF (top left, logistic regression coefficient for slope, $\beta_1$ = 0.002 +/- 0.002, p =
824　0.450), number of surface-exposed lipoproteins (top right, $\beta_1$ = 0.037 +/- 0.017, p = 0.03, logistic
825　regression), number of Erps (bottom left, $\beta_1$ = 0.087 +/- 0.053, p = 0.10, logistic regression), and
826　number of Mlps (bottom right, $\beta_1$ = 0.048 +/- 0.055　p = 0.38, logistic regression). **E.** For each
827　OspC type, mean probability of dissemination vs mean number of ORF (top left), mean number
828　of surface-exposed lipoproteins (top right), mean number of Erps (bottom left), and mean
829　number of Mlps (bottom right).
830

831



832

**Supplemental Figure 8:** Manhattan Plots showing the association of individual lipoproteins with OspC type A (top panel), Osp C type K (middle panel), and RST1 (bottom panel). Individual lipoproteins are annotated by their localization. P-IM: Periplasmic inner membrane. P-OM: Periplasmic outer membrane. S: surface.

837

838

**List of supplemental data files**

839

840    Supplemental Table 1: Summary table of isolates and phenotypes
841    Supplemental Table 2: List of isolates and phenotypes
842    Supplemental Table 3: Assembly statistics
843    Supplemental Table 4: Association statistics for plasmids, as inferred from PFam32 types.
844    Supplemental Table 5: Association statistics for plasmids, as inferred from B31 reference
845    Supplemental Table 5: Association statistics for lineage model
846    Supplemental Table 6: Association statistics for lineage model restricted to surface lipoproteins
847    Supplemental Table 7: Association statistics for OspC type A associations
848    Supplemental Table 8: Association statistics for OspC type K associations
849    Supplemental Table 9: Association statistics for RST1 associations
850    Supplemental Data File 1: List of ortholog groups with reference sequences
851    Supplemental Data File 2: High resolution version of presence/absence matrix in Figure 5B.
852

853

854   **Supplemental Note 1:**

855   The clock rate (in substitutions/site/year) for our initial model using a non-informative (CTMC

856   rate reference) prior failed to converge—resulting in posterior 95% posterior density range from

857   $5 \times 10^{-25}$ substitutions/site/year to $1.2 \times 10^{-8}$ substitutions/site/year—the implausibly small values

858   at the lower end of the range are indicative of an insufficient temporal signal associated with

859   genetic diversity in the core genome to establish an estimate without a priori assumptions.

860   However, the inferred clock rate posterior had a clear single mode and a reasonable posterior

861   mean ($1.8 \times 10^{-9}$ substitutions/site/year). To address this, we incorporated a priori information on

862   mutation (gamma prior with shape 2, scale $1 \times 10^{-9}$, for which 95% of the density is between 3.55

863   $\times 10^{-10}$ substitutions/site/year and $4.47 \times 10^{-9}$ substitutions/site/year, concordant with previous

864   suggestions that the rate is approximately $1 \times 10^{-9}$ substitutions/site/year[52]). This analysis

865   suggests that the common ancestry of circulating human-infectious populations was remote

866   (95% posterior density for Midwest strains: 380,000 - 11.8 million years; 95% posterior density

867   for Slovenian strains: 379,000 - 11.5 million years; all strains: 380,000 years, 11.8 million years)

868   (Figure 2E-F). We also ran models with a fixed rate across a variety of reasonable values (1e-

869   10 to 1e-8) (Figure 2G).

**References**

871  1.  Radolf JD, Strle K, Lemieux JE, Strle F. Lyme Disease in Humans. Curr Issues Mol Biol.
872      2021;42: 333–384.

873  2.  Lantos PM, Rumbaugh J, Bockenstedt LK, Falck-Ytter YT, Aguero-Rosenfeld ME,
874      Auwaerter PG, et al. Clinical practice guidelines by the Infectious Diseases Society of
875      America (IDSA), American Academy of Neurology (AAN), and American College of
876      Rheumatology (ACR): 2020 guidelines for the prevention, diagnosis and treatment of Lyme
877      disease. Clin Infect Dis. 2021;72: e1–e48.

878  3.  Strle K, Jones KL, Drouin EE, Li X, Steere AC. Borrelia burgdorferi RST1 (OspC type A)
879      genotype is associated with greater inflammation and more severe Lyme disease. Am J
880      Pathol. 2011;178: 2726–2739.

881  4.  Wormser GP, Brisson D, Liveris D, Hanincová K, Sandigursky S, Nowakowski J, et al.
882      Borrelia burgdorferi genotype predicts the capacity for hematogenous dissemination during
883      early Lyme disease. J Infect Dis. 2008;198: 1358–1364.

884  5.  Wormser GP, Liveris D, Nowakowski J, Nadelman RB, Cavaliere LF, McKenna D, et al.
885      Association of specific subtypes of Borrelia burgdorferi with hematogenous dissemination in
886      early Lyme disease. J Infect Dis. 1999;180: 720–725.

887  6.  Wang G, Ojaimi C, Wu H, Saksenberg V, Iyer R, Liveris D, et al. Disease severity in a
888      murine model of lyme borreliosis is associated with the genotype of the infecting Borrelia
889      burgdorferi sensu stricto strain. J Infect Dis. 2002;186: 782–791.

890  7.  Wang G, van Dam AP, Schwartz I, Dankert J. Molecular typing of Borrelia burgdorferi
891      sensu lato: taxonomic, epidemiological, and clinical implications. Clin Microbiol Rev.
892      1999;12: 633–653.

893  8.  Wang G, Ojaimi C, Iyer R, Saksenberg V, McClain SA, Wormser GP, et al. Impact of
894      genotypic variation of Borrelia burgdorferi sensu stricto on kinetics of dissemination and
895      severity of disease in C3H/HeJ mice. Infect Immun. 2001;69: 4303–4312.

896  9.  Fraser CM, Casjens S, Huang WM, Sutton GG, Clayton R, Lathigra R, et al. Genomic
897      sequence of a Lyme disease spirochaete, Borrelia burgdorferi. Nature. 1997;390: 580–586.

898  10. Casjens S, Palmer N, van Vugt R, Huang WM, Stevenson B, Rosa P, et al. A bacterial
899      genome in flux: the twelve linear and nine circular extrachromosomal DNAs in an infectious
900      isolate of the Lyme disease spirochete Borrelia burgdorferi. Mol Microbiol. 2000;35: 490–
901      516.

902  11. Casjens S, van Vugt R, Tilly K, Rosa PA, Stevenson B. Homology throughout the multiple
903      32-kilobase circular plasmids present in Lyme disease spirochetes. J Bacteriol. 1997;179:
904      217–227.

905  12. Casjens SR, Di L, Akther S, Mongodin EF, Luft BJ, Schutzer SE, et al. Primordial origin and
906      diversification of plasmids in Lyme disease agent bacteria. BMC Genomics. 2018;19: 218.

907  13. Margos G, Hepner S, Mang C, Marosevic D, Reynolds SE, Krebs S, et al. Lost in plasmids:

908    next generation sequencing and the complex genome of the tick-borne pathogen Borrelia
909    burgdorferi. BMC Genomics. 2017;18: 422.

910  14. Tyler S, Tyson S, Dibernardo A, Drebot M, Feil EJ, Graham M, et al. Whole genome
911    sequencing and phylogenetic analysis of strains of the agent of Lyme disease Borrelia
912    burgdorferi from Canadian emergence zones. Sci Rep. 2018;8: 10552.

913  15. Castillo-Ramírez S, Fingerle V, Jungnick S, Straubinger RK, Krebs S, Blum H, et al. Trans-
914    Atlantic exchanges have shaped the population structure of the Lyme disease agent
915    Borrelia burgdorferi sensu stricto. Sci Rep. 2016;6: 22794.

916  16. Walter KS, Carpi G, Caccone A, Diuk-Wasser MA. Genomic insights into the ancient
917    spread of Lyme disease across North America. Nat Ecol Evol. 2017;1: 1569–1576.

918  17. Di L, Pagan PE, Packer D, Martin CL, Akther S, Ramrattan G, et al. BorreliaBase: a
919    phylogeny-centered browser of Borrelia genomes. BMC Bioinformatics. 2014;15: 233.

920  18. Carpi G, Walter KS, Bent SJ, Hoen AG, Diuk-Wasser M, Caccone A. Whole genome
921    capture of vector-borne pathogens from mixed DNA samples: a case study of Borrelia
922    burgdorferi. BMC Genomics. 2015;16: 434.

923  19. Schutzer SE, Fraser-Liggett CM, Casjens SR, Qiu W-G, Dunn JJ, Mongodin EF, et al.
924    Whole-genome sequences of thirteen isolates of Borrelia burgdorferi. J Bacteriol. 2011;193:
925    1018–1020.

926  20. Wang G, Liveris D, Mukherjee P, Jungnick S, Margos G, Schwartz I. Molecular Typing of
927    Borrelia burgdorferi. Curr Protoc Microbiol. 2014;34: 12C.5.1–31.

928  21. Liveris D, Gazumyan A, Schwartz I. Molecular typing of Borrelia burgdorferi sensu lato by
929    PCR-restriction fragment length polymorphism analysis. J Clin Microbiol. 1995;33: 589–
930    595.

931  22. Liveris D, Wormser GP, Nowakowski J, Nadelman R, Bittker S, Cooper D, et al. Molecular
932    typing of Borrelia burgdorferi from Lyme disease patients by PCR-restriction fragment
933    length polymorphism analysis. J Clin Microbiol. 1996;34: 1306–1309.

934  23. Qiu W-G, Schutzer SE, Bruno JF, Attie O, Xu Y, Dunn JJ, et al. Genetic exchange and
935    plasmid transfers in Borrelia burgdorferi sensu stricto revealed by three-way genome
936    comparisons and multilocus sequence typing. Proc Natl Acad Sci U S A. 2004;101: 14150–
937    14155.

938  24. Bunikis J, Garpmo U, Tsao J, Berglund J, Fish D, Barbour AG. Sequence typing reveals
939    extensive strain diversity of the Lyme borreliosis agents Borrelia burgdorferi in North
940    America and Borrelia afzelii in Europe. Microbiology. 2004;150: 1741–1755.

941  25. Barbour AG, Travinsky B. Evolution and distribution of the ospC Gene, a transferable
942    serotype determinant of Borrelia burgdorferi. MBio. 2010;1. Available:
943    https://www.ncbi.nlm.nih.gov/pubmed/20877579

944  26. Cerar T, Strle F, Stupica D, Ruzic-Sabljic E, McHugh G, Steere AC, et al. Differences in
945    Genotype, Clinical Features, and Inflammatory Potential of Borrelia burgdorferi sensu
946    stricto Strains from Europe and the United States. Emerg Infect Dis. 2016;22: 818–827.

947    27.   Hanincova K, Mukherjee P, Ogden NH, Margos G, Wormser GP, Reed KD, et al. Multilocus
948         sequence typing of Borrelia burgdorferi suggests existence of lineages with differential
949         pathogenic properties in humans. PLoS One. 2013;8: e73066.

950    28.   Margos G, Gatewood AG, Aanensen DM, Hanincová K, Terekhova D, Vollmer SA, et al.
951         MLST of housekeeping genes captures geographic population structure and suggests a
952         European origin of Borrelia burgdorferi. Proc Natl Acad Sci U S A. 2008;105: 8730–8735.

953    29.   Strle K, Shin JJ, Glickstein LJ, Steere AC. Association of a Toll-like receptor 1
954         polymorphism with heightened Th1 inflammatory responses and antibiotic-refractory Lyme
955         arthritis. Arthritis Rheum. 2012;64: 1497–1507.

956    30.   Jones KL, Glickstein LJ, Damle N, Sikand VK, McHugh G, Steere AC. Borrelia burgdorferi
957         genetic markers and disseminated disease in patients with early Lyme disease. J Clin
958         Microbiol. 2006;44: 4407–4413.

959    31.   Dykhuizen DE, Brisson D, Sandigursky S, Wormser GP, Nowakowski J, Nadelman RB, et
960         al. The Propensity of Different Borrelia burgdorferi sensu stricto Genotypes to Cause
961         Disseminated Infections in Humans. Am J Trop Med Hyg. 2008;78: 806–810.

962    32.   Brisson D, Baxamusa N, Schwartz I, Wormser GP. Biodiversity of Borrelia burgdorferi
963         strains in tissues of Lyme disease patients. PLoS One. 2011;6: e22926.

964    33.   Ružić-Sabljić E, Maraspin V, Stupica D, Rojko T, Bogovič P, Strle F, et al. Comparison of
965         MKP and BSK-H media for the cultivation and isolation of Borrelia burgdorferi sensu lato.
966         PLoS One. 2017;12: e0171622.

967    34.   Wang G, Iyer R, Bittker S, Cooper D, Small J, Wormser GP, et al. Variations in Barbour-
968         Stoenner-Kelly culture medium modulate infectivity and pathogenicity of Borrelia burgdorferi
969         clinical isolates. Infect Immun. 2004;72: 6702–6706.

970    35.   Centers for Disease Control and Prevention (CDC). Recommendations for test
971         performance and interpretation from the Second National Conference on Serologic
972         Diagnosis of Lyme Disease. MMWR Morb Mortal Wkly Rep. 1995;44: 590–591.

973    36.   Mygland A, Ljøstad U, Fingerle V, Rupprecht T, Schmutzhard E, Steiner I, et al. EFNS
974         guidelines on the diagnosis and management of European Lyme neuroborreliosis. Eur J
975         Neurol. 2010;17: 8–16, e1–4.

976    37.   Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.
977         Bioinformatics. 2014;30: 2114–2120.

978    38.   Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a new
979         genome assembly algorithm and its applications to single-cell sequencing. J Comput Biol.
980         2012;19: 455–477.

981    39.   Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome
982         assemblies. Bioinformatics. 2013;29: 1072–1075.

983    40.   Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. Genome
984         Biol. 2019;20: 257.

985    41.   Jolley KA, Maiden MCJ. BIGSdb: Scalable analysis of bacterial genome variation at the

986     population level. BMC Bioinformatics. 2010;11: 595.

987   42. Murray KD, Webers C, Ong CS, Borevitz J, Warthmann N. kWIP: The k-mer weighted inner
988       product, a de novo estimator of genetic similarity. PLoS Comput Biol. 2017;13: e1005727.

989   43. Seemann T. Prokka: rapid prokaryotic genome annotation. Bioinformatics. 2014;30: 2068–
990       2069.

991   44. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, et al. Roary: rapid
992       large-scale prokaryote pan genome analysis. Bioinformatics. 2015;31: 3691–3693.

993   45. Price MN, Dehal PS, Arkin AP. FastTree 2--approximately maximum-likelihood trees for
994       large alignments. PLoS One. 2010;5: e9490.

995   46. Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, et al. Bioconductor:
996       open software development for computational biology and bioinformatics. Genome Biol.
997       2004;5: R80.

998   47. Ihaka R, Gentleman R. R: A Language for Data Analysis and Graphics. J Comput Graph
999       Stat. 1996;5: 299–314.

1000  48. Wickham H, Others. Tidyverse: Easily install and load the "tidyverse." R package version.
1001      2017;1: 2017.

1002  49. Yu G, Smith DK, Zhu H, Guan Y, Lam TT. ggtree : an r package for visualization and
1003      annotation of phylogenetic trees with their covariates and other associated data. McInerny
1004      G, editor. Methods Ecol Evol. 2017;8: 28–36.

1005  50. Eddy SR. A new generation of homology search tools based on probabilistic inference.
1006      Genome Inform. 2009;23: 205–211.

1007  51. Schwartz I, Margos G, Casjens SR, Qiu W-G, Eggers CH. Multipartite Genome of Lyme
1008      Disease Borrelia: Structure, Variation and Prophages. Curr Issues Mol Biol. 2021;42: 409–
1009      454.

1010  52. Hoen AG, Margos G, Bent SJ, Diuk-Wasser MA, Barbour A, Kurtenbach K, et al.
1011      Phylogeography of *Borrelia burgdorferi* in the eastern United States reflects multiple
1012      independent Lyme disease emergence events. Proc Natl Acad Sci U S A. 2009;106:
1013      15013–15018.

1014  53. Radolf JD, Caimano MJ, Stevenson B, Hu LT. Of ticks, mice and men: understanding the
1015      dual-host lifestyle of Lyme disease spirochaetes. Nat Rev Microbiol. 2012;10: 87–99.

1016  54. Brisson D, Drecktrah D, Eggers CH, Samuels DS. Genetics of Borrelia burgdorferi. Annu
1017      Rev Genet. 2012;46: 515–536.

1018  55. Steere AC, Strle F, Wormser GP, Hu LT, Branda JA, Hovius JWR, et al. Lyme borreliosis.
1019      Nat Rev Dis Primers. 2016;2: 16090.

1020  56. Dowdell AS, Murphy MD, Azodi C, Swanson SK, Florens L, Chen S, et al. Comprehensive
1021      Spatial Analysis of the Borrelia burgdorferi Lipoproteome Reveals a Compartmentalization
1022      Bias toward the Bacterial Surface. J Bacteriol. 2017;199. doi:10.1128/JB.00658-16

1023    57. Stevenson B, Tilly K, Rosa PA. A family of genes located on four separate 32-kilobase
1024         circular plasmids in Borrelia burgdorferi B31. J Bacteriol. 1996;178: 3508–3516.

1025    58. Brissette CA, Cooley AE, Burns LH, Riley SP, Verma A, Woodman ME, et al. Lyme
1026         borreliosis spirochete Erp proteins, their known host ligands, and potential roles in
1027         mammalian infection. Int J Med Microbiol. 2008;298: 257–267.

1028    59. Porcella SF, Popova TG, Akins DR, Li M, Radolf JD, Norgard MV. Borrelia burgdorferi
1029         supercoiled plasmids encode multicopy tandem open reading frames and a lipoprotein
1030         gene family. J Bacteriol. 1996;178: 3293–3307.

1031    60. Porcella SF, Fitzpatrick CA, Bono JL. Expression and Immunological Analysis of the
1032         Plasmid-Borne mlp Genes of Borrelia burgdorferiStrain B31. Infect Immun. 2000;68: 4992–
1033         5001.

1034    61. Earle SG, Wu C-H, Charlesworth J, Stoesser N, Gordon NC, Walker TM, et al. Identifying
1035         lineage effects when controlling for population structure improves power in bacterial
1036         association studies. Nat Microbiol. 2016;1: 16041.

1037    62. Lees JA, Vehkala M, Välimäki N, Harris SR, Chewapreecha C, Croucher NJ, et al.
1038         Sequence element enrichment analysis to determine the genetic basis of bacterial
1039         phenotypes. Nat Commun. 2016;7: 12797.

1040    63. Terekhova D, Iyer R, Wormser GP, Schwartz I. Comparative genome hybridization reveals
1041         substantial variation among clinical isolates of Borrelia burgdorferi sensu stricto with
1042         different pathogenic properties. J Bacteriol. 2006;188: 6124–6134.

1043    64. Qiu W-G, Bruno JF, McCaig WD, Xu Y, Livey I, Schriefer ME, et al. Wide distribution of a
1044         high-virulence Borrelia burgdorferi clone in Europe and North America. Emerg Infect Dis.
1045         2008;14: 1097–1104.

1046    65. Mongodin EF, Casjens SR, Bruno JF, Xu Y, Drabek EF, Riley DR, et al. Inter- and intra-
1047         specific pan-genomes of Borrelia burgdorferi sensu lato: genome stability and adaptive
1048         radiation. BMC Genomics. 2013;14: 693.

1049    66. Labandeira-Rey M, Skare JT. Decreased infectivity in Borrelia burgdorferi strain B31 is
1050         associated with loss of linear plasmid 25 or 28-1. Infect Immun. 2001;69: 446–455.

1051    67. Magunda PRH, Bankhead T. Investigating the potential role of non-vls genes on linear
1052         plasmid 28–1 in virulence and persistence by Borrelia burgdorferi. BMC Microbiol. 2016;16:
1053         180.

1054    68. Labandeira-Rey M, Seshu J, Skare JT. The absence of linear plasmid 25 or 28-1 of Borrelia
1055         burgdorferi dramatically alters the kinetics of experimental infection via distinct
1056         mechanisms. Infect Immun. 2003;71: 4608–4613.

1057    69. Purser JE, Norris SJ. Correlation between plasmid content and infectivity in Borrelia
1058         burgdorferi. Proc Natl Acad Sci U S A. 2000;97: 13865–13870.

1059    70. Lawrenz MB, Kawabata H, Purser JE, Norris SJ. Decreased electroporation efficiency in
1060         Borrelia burgdorferi containing linear plasmids lp25 and lp56: impact on transformation of
1061         infectious B. burgdorferi. Infect Immun. 2002;70: 4798–4804.

1062 71. Rego ROM, Bestor A, Rosa PA. Defining the plasmid-borne restriction-modification
1063   systems of the Lyme disease spirochete Borrelia burgdorferi. J Bacteriol. 2011;193: 1161–
1064   1171.

1065 72. Lin Y-P, Benoit V, Yang X, Martínez-Herranz R, Pal U, Leong JM. Strain-specific variation
1066   of the decorin-binding adhesin DbpA influences the tissue tropism of the lyme disease
1067   spirochete. PLoS Pathog. 2014;10: e1004238.

1068 73. Lin Y-P, Tan X, Caine JA, Castellanos M, Chaconas G, Coburn J, et al. Strain-specific joint
1069   invasion and colonization by Lyme disease spirochetes is promoted by outer surface
1070   protein C. PLoS Pathog. 2020;16: e1008516.

1071 74. Caine JA, Lin Y-P, Kessler JR, Sato H, Leong JM, Coburn J. Borrelia burgdorferi outer
1072   surface protein C (OspC) binds complement component C4b and confers bloodstream
1073   survival. Cell Microbiol. 2017;19. doi:10.1111/cmi.12786

1074 75. Xu Q, McShan K, Liang FT. Essential protective role attributed to the surface lipoproteins of
1075   Borrelia burgdorferi against innate defences. Mol Microbiol. 2008;69: 15–29.

1076 76. Tilly K, Bestor A, Rosa PA. Lipoprotein succession in Borrelia burgdorferi: similar but
1077   distinct roles for OspC and VlsE at different stages of mammalian infection. Mol Microbiol.
1078   2013;89: 216–227.

1079 77. El-Hage N, Babb K, Carroll JA, Lindstrom N, Fischer ER, Miller JC, et al. Surface exposure
1080   and protease insensitivity of Borrelia burgdorferi Erp (OspEF-related) lipoproteins.
1081   Microbiology. 2001;147: 821–830.

1082 78. Stevenson B, El-Hage N, Hines MA, Miller JC, Babb K. Differential binding of host
1083   complement inhibitor factor H by Borrelia burgdorferi Erp surface proteins: a possible
1084   mechanism underlying the expansive host range of Lyme disease spirochetes. Infect
1085   Immun. 2002;70: 491–497.

1086 79. Lin Y-P, Bhowmick R, Coburn J, Leong JM. Host cell heparan sulfate glycosaminoglycans
1087   are ligands for OspF-related proteins of the Lyme disease spirochete. Cell Microbiol.
1088   2015;17: 1464–1476.

1089 80. Pereira MJ, Wager B, Garrigues RJ, Gerlach E, Quinn JD, Dowdell AS, et al. Lipoproteome
1090   screening of the Lyme disease agent identifies inhibitors of antibody-mediated complement
1091   killing. Proc Natl Acad Sci U S A. 2022;119: e2117770119.

1092 81. Grillon A, Scherlinger M, Boyer P-H, De Martino S, Perdriger A, Blasquez A, et al.
1093   Characteristics and clinical outcomes after treatment of a national cohort of PCR-positive
1094   Lyme arthritis. Semin Arthritis Rheum. 2018. doi:10.1016/j.semarthrit.2018.09.007

1095 82. Livey I, Gibbs CP, Schuster R, Dorner F. Evidence for lateral transfer and recombination in
1096   OspC variation in Lyme disease Borrelia. Mol Microbiol. 1995;18: 257–269.

1097 83. Qiu W-G, Martin CL. Evolutionary genomics of Borrelia burgdorferi sensu lato: findings,
1098   hypotheses, and the rise of hybrids. Infect Genet Evol. 2014;27: 576–593.

1099 84. Balding DJ. A tutorial on statistical methods for population association studies. Nat Rev
1100   Genet. 2006;7: 781–791.

1101  85. Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. FaST linear mixed
1102      models for genome-wide association studies. Nat Methods. 2011;8: 833–835.

1103  86. Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association studies.
1104      Nat Genet. 2012;44: 821–824.

1105

Click here to access/download
**Supporting Information - Compressed/ZIP File Archive**
Supplemental_Data.zip