

1 *Short title*

2 ***Staphylococcus chromogenes* multilocus sequence typing**

3 *Long title*

4 **Characterization of genetic diversity and population structure within *Staphylococcus***
5 ***chromogenes* by multilocus sequence typing**

6

7 Rebeca Huebner^{1,#a}, Robert Mugabi^{2,#b}, Gabriella Hetesy^{2,#c}, Lawrence Fox³, Sarne De Vlieghe⁴,
8 Anneleen De Visscher⁴, John W Barlow^{2*}, and George Sensabaugh¹

9

10 ¹ Division of Infectious Diseases and Vaccinology, School of Public Health, University of
11 California, Berkeley, USA

12 ² Department of Animal and Veterinary Sciences, University of Vermont, Burlington, Vermont,
13 USA

14 ³ College of Veterinary Medicine, Washington State University, Pullman, WA, USA

15 ⁴ M-team *and* Mastitis and Milk Quality Research Unit, Department of Reproduction, Obstetrics,
16 and Herd Health, Faculty of Veterinary Medicine, Ghent University, Merelbeke, Belgium

17 ^{#a} Division of Infectious Diseases and Vaccinology, School of Public Health, University of
18 California, Berkeley.

19 ^{#b} Department of Veterinary Diagnostic and Production Animal Medicine, College of Veterinary
20 Medicine, Iowa State University

21 ^{#c} Royal Veterinary College, Camden campus, 4 Royal College St, London NW1 0TU

22 ^{*}Corresponding Author

23 E-mail: john.barlow@uvm.edu (JWB)

24 **Abstract**

25 *Staphylococcus chromogenes* is a common skin commensal in cattle and has been identified as
26 a frequent cause of bovine mastitis and intramammary infections. To better understand the
27 extent of strain diversity within this species and to facilitate study of strain variation as a factor
28 in pathogenicity, we have developed a seven locus Multilocus Sequence Typing (MLST) scheme.
29 The scheme was tested on 120 isolates collected from three geographic locations, Vermont and
30 Washington State in the United States and Belgium. A total of 46 sequence types (STs) were
31 identified with most of the STs being location specific. The utility of the typing scheme is
32 indicated by a discrimination power of 95.6% for all isolates and greater than 90% for isolates
33 from each of the three locations. Phylogenetic analysis placed 39 of the 46 STs into single core
34 group consistent with a common genetic lineage; the STs in this group differ by less than 0.5%
35 at the nucleotide sequence level. Most of the diversification in this lineage group can be
36 attributed to mutation; recombination plays a limited role. This lineage group includes two
37 clusters of single nucleotide variants in starburst configurations indicative of recent clonal
38 expansion; nearly 50% of the isolates sampled in this study are in these two clusters. The
39 remaining seven STs were set apart from the core group by having alleles with highly variable
40 sequences at one or more loci. Recombination had a higher impact than mutation in the
41 diversification of these outlier STs. Alleles with hypervariable sequences were detected at five
42 of the seven loci used in the MLST scheme; the average sequence distances between the
43 hypervariable alleles and the common core alleles ranged from 12 to 34 nucleotides. The extent

44 of these sequence differences suggests the hypervariable alleles may be remnants of an
45 ancestral genotype.

46 **Introduction**

47 *Staphylococcus chromogenes* was first recognized by Devriese et al. [1] as one of two
48 subspecies of *Staphylococcus hyicus*, and was subsequently elevated to a novel species based
49 on chemical, physiological and DNA-DNA re-association binding experiments [2]. Phylogenetic
50 analyses by multi-locus and whole genome sequencing place *S. chromogenes* in a cluster with *S.*
51 *hyicus* and *Staphylococcus agnetis* [3, 4]. The habitat of *S. chromogenes* is described as the body
52 surface of cattle, pigs and poultry [2].

53 *S. chromogenes* is most commonly identified as a skin commensal and opportunistic
54 mammary pathogen in cattle, sheep, goats and milking buffalo. It is a frequent cause of bovine
55 mastitis [5], and reported as a skin pathogen of pigs and goats and as a cause of caprine mastitis
56 [6-8]. *S. chromogenes* is recognized as one of the most frequent species of non-*aureus*
57 staphylococci causing subclinical (asymptomatic) intramammary infections in dairy cattle in
58 Europe and the United States [reviewed in 5 and 9]. *S. chromogenes* has been identified as a
59 cause of persistent intramammary infections in dairy cattle [10-12], and infections appear to be
60 associated with increased milk somatic cell counts (i.e. intramammary inflammation or
61 subclinical mastitis) [12-14]. The organism has also been identified from extra-mammary skin
62 swabs of cattle, including udder skin, teat apex, and streak canal [15-18], and compared to
63 other non-*aureus* *Staphylococcus* species *S. chromogenes* is less commonly isolated from
64 environmental sites in surveys of dairy farm environmental sources (e.g. barn air, surfaces and

65 bedding) [19]. Some authors have suggested that *S. chromogenes* intramammary infection or
66 colonization of teat skin may have a protective effect against *S. aureus* mastitis [20, 21]. Using
67 PFGE, multiple strains (pulsotypes) of *S. chromogenes* have been isolated from intramammary
68 infections and extramammary skin sites of dairy cattle within individual herds [15, 18]. *S.*
69 *chromogenes* has been isolated infrequently from nasal swabs of humans in close contact with
70 cattle [9]. Development of portable sequence-based strain typing systems has been
71 recommended to improve our understanding of the epidemiology of *S. chromogenes* [5].

72 In this paper we report the development of a multilocus sequence typing (MLST)
73 scheme that provides a practical, portable, sequence-based approach for the identification of
74 strain types and the characterization of relationships between clonal lineages in *S.*
75 *chromogenes*. MLST schemes have been developed for a number of staphylococcal species
76 including *S. aureus*, *S. epidermidis*, *S. haemolyticus*, *S. hominis*, *S. lugdunensis*, *S.*
77 *pseudintermedius*, and *S. carnosus* [22-28]. This *S. chromogenes* MLST scheme is based on the
78 detection of genetic variation in seven housekeeping genes in 120 isolates collected from dairy
79 cattle (*Bos taurus*) in three geographic locations, Vermont and Washington State in the United
80 States and Belgium. This sample population allows assessment of both genetic and geographic
81 diversity present in this species. The scheme has been designed such that the seven loci are
82 well separated around the ca. 2.34 Mb genome of *S. chromogenes* to maximize the opportunity
83 to evaluate the extent to which recombination may play a role in shaping diversity at the
84 population level.

85 **Material and methods**

86 **Bacterial strains and DNA Isolation**

87 A total of 120 isolates were investigated in this study; these isolates originated from the
88 collections of three laboratories. The isolates were originally collected from dairy cattle in
89 Vermont (n=46) and Washington (n=24) in the USA and from Belgium (n=48), and pigs in
90 Vermont (n=2). The isolates from Belgium were collected from dairy cattle teat apex swabs
91 (n=20) and individual mammary quarter milk samples from apparent healthy quarters (n=28)
92 [14, 17]. The Washington isolates were collected from quarter milk samples of dairy cows with
93 intramammary infections [29]. The Vermont isolates were collected from 5 dairy farms from
94 either quarter milk samples of cows with intramammary infections (n=31), cow teat orifice
95 swabs (n=1), cow hock skin swabs (n=3), and bulk tank milk (n=11); two isolates originated from
96 pig nasal swabs collected from one of the 5 Vermont farms. This study was carried out in strict
97 accordance with the recommendations in the Guide for the Care and Use of Laboratory Animals
98 of the National Institutes of Health. The protocol was approved by the Committee on the Ethics
99 of Animal Experiments of the University of Vermont (Protocol Number: 13-033).

100 All isolates were verified as *S. chromogenes* by sequence analysis of the *tuf* and *rpoB*
101 gene amplicon fragments with > 97 % sequence identity [30, 31]. The draft genome sequence of
102 *S. chromogenes* strain MU970 was downloaded from the NCBI microbial genome database
103 (GenBank accession JMJF01000000) to provide a genome reference sequence [12]. All isolates
104 were shared between the University of California Berkeley (UCB) and University of Vermont
105 (UVM) laboratories, and DNA extraction, amplification and analysis procedures were replicated
106 in both labs. Strain MU970 (gift from J. Middleton, University of Missouri) was cultured for
107 DNA extraction and sequence amplification in the UVM lab.

108 In the University of California Berkeley (UCB) lab, isolates were grown aerobically
109 overnight in tryptic soy broth (TSB, BBL) at 37°C without shaking and then plated on tryptic soy
110 agar plates with 5% sheep blood (TSA, BBL) and incubated aerobically at 37°C. Single colonies
111 were passed from each TSA plate and grown overnight at 37°C in 1 mL of TSB. To achieve mid-
112 phase growth, 100µL of overnight growth was combined with 5mL TSB and incubated at 37°C
113 with shaking for 4 hours. Cell pellets were collected from 3 mL of mid-phase growth by
114 centrifugation at 10,000 rpm for 5 minutes. Alternatively in parallel, at the University of
115 Vermont lab (UVM), a pure primary culture was grown aerobically for 48hrs on TSA, and single
116 colonies from this growth were inoculated to 5 ml TSB, grown aerobically overnight at 37°C and
117 cell pellets were collected by centrifugation directly from 1.8 ml of overnight TSB culture. The
118 cell pellets were then frozen at -20°C until DNA extraction could be performed (UCB) or held at
119 4°C and processed within 48 hours (UVM).

120 DNA extraction was performed using a DNeasy Blood & Tissue kit (Qiagen, Valencia, CA,
121 USA) according to manufacturer's instructions with the modification that the initial lysis buffer
122 was supplemented with lysostaphin (22 U/ml; Sigma-Aldrich). DNA yield and quality was
123 assessed by electrophoresis using a 0.75% agarose gel containing the DNA stain GelStar (Lonza,
124 Rockland, ME, USA) in 1X Tris Borate EDTA (TBE) buffer. Aliquots of DNA were stored at -20°C.

125 **Selection of target loci for MLST**

126 Fifteen loci were assessed initially as potential candidates for the *S. chromogenes* MLST
127 scheme. Nine loci originated from MLST schemes used for *S. aureus*, *S. epidermidis*, and *S.*
128 *saprophyticus* [<http://pubmlst.org> and unpublished] and six additional loci were selected from
129 the GenBank annotation listing for the *S. chromogenes* MU970 draft genome. PCR primers for

130 each of the candidate loci were designed using Primer-BLAST to yield sequence segments 700-
131 850 bp in length. The candidate loci were evaluated for sequence variation using a panel of 27
132 isolates from Vermont and those exhibiting the greatest sequence diversity were selected for
133 further evaluation. As a final filter, we sought to assess whether the loci were relatively evenly
134 distributed around the genome to maximize the potential for diversity resulting from inter-locus
135 recombination. Having determined strong synteny between the contigs in the draft MU970
136 genome and the complete genome sequence of the closely related species, *S. hyicus* (GenBank
137 accession CP008747.1), we determined that the seven loci selected were at least 250Kbp apart,
138 thus satisfying the criterion. Details for the seven loci selected for the MLST scheme are provided
139 in Supporting Information S1 Table.

140 **Target gene amplification and nucleotide sequencing**

141 Target genes were amplified using the polymerase chain reaction employing a master
142 mix containing 16.75 μ L DNase-free H₂O, 2.5 μ L PCR buffer, 1.5 μ L 50mM MgCl₂, 0.5 μ L 10mM
143 dNTP, and 0.25 μ L 0.5 U/ μ L Taq polymerase (Invitrogen) per reaction. Master mix was aliquoted
144 into tubes for each locus being amplified and 0.25 μ L each of 25 μ M forward and reverse primer
145 was added for each reaction. Three microliters of DNA was added to 22 μ L of the master mix
146 and primers. PCR cycling included heating to 95°C for 7 minutes, followed by 35 cycles of 94°C
147 for 45 seconds, 55°C for 45 seconds and 2 minutes of 72°C. On the last cycle, the samples were
148 heated to 72°C for 5 min. The samples were then held at 4°C until further analysis could be
149 completed.

150 Amplification products (3 μ L) were evaluated by electrophoresis at 150V for
151 approximately one hour on 1.5% agarose gel containing the DNA stain GelStar. PCR products

152 exhibiting a single band at the predicted amplicon size were processed for sequence analysis by
153 the UC Berkeley sequencing facility. PCR amplification and target gene sequencing were
154 replicated at UVM and the replicate sets of amplicons were processed for sequence analysis by
155 the University of Vermont DNA sequencing facility.

156 **Analysis of MLST sequence data**

157 Raw sequence data containing both forward and reverse reads were recorded in FASTA
158 format for analysis. Sequences were aligned using the MUSCLE function in MEGA 6.06 [32].
159 After alignment, single strand overhangs and any ambiguous reads were trimmed from the
160 ends of each sequence. Any missing or ambiguous nucleotides were resolved by reviewing the
161 trace data using the FinchTV 1.4.0. viewer. Once a consensus sequence was determined for
162 each sample, all sequences for a single locus were combined into a single FASTA file. The
163 sequences were trimmed again to obtain a standard length. Sites exhibiting single nucleotide
164 polymorphism (SNP) were identified in MEGA. For each of the seven loci, the gene sequence
165 present in the MU970 reference sequence was arbitrarily defined as allele 1 and new alleles
166 were identified by pairwise comparison of SNP sites; each new allele was assigned a number.
167 Sequence types (STs) were defined by unique allelic profiles at the seven loci. MEGA and DnaSP
168 5.10 were used for assessment of population genetics parameters such as nucleotide diversity
169 and discrimination index [32, 33]. DnaSP was also used to concatenate the sequences of the
170 seven loci for each ST.

171 The 7-locus concatenated nucleotide sequence data were used for the construction of
172 phylogenetic trees generated by the neighbor joining (NJ) algorithm in MEGA with 1000
173 bootstrap replications. Clonal clusters of sequence types were identified at the level of single

174 locus variants (SLVs) and double locus variants (DLVs) using the web program eBurst3 as
175 implemented in goeBURST [34, 35]. Evidence for recombination was assessed by surveying
176 allelic sequences at each locus within clonal subgroups delineated by eBURST and phylogenetic
177 analyses: alleles at a locus within a clonal subgroup differing at a single nucleotide site were
178 scored as mutations whereas alleles differing at multiple nucleotide sites and alleles shared
179 between different clonal subgroups were scored as recombination events [36]. Recombination
180 between loci was assessed using the four-gamete test [37] as implemented in DnaSP; this test
181 detects the minimum number of recombination events (RM) in the history of the sample. For
182 this test, the concatenated sequences were constructed with the first locus (arcC) sequences
183 appended to the end of the 7-locus concatenation to detect possible recombination between
184 the last and first locus in the circular genome. The pairwise homoplasmy index (PHI) for
185 recombination was measured using the program implemented in SplitsTree [38, 39].

186

187 **Results and Discussion**

188 **Genetic diversity in *S. chromogenes***

189 The MLST scheme for *S. chromogenes* is based on characterization of nucleotide
190 sequence variation in fragments of seven housekeeping genes in 120 isolates. Overall, 216
191 nucleotide substitutions at 213 sites were identified in the 4563 bp of genome sequence
192 covered by the scheme (Table 1). The 216 nucleotide substitutions resulted in 57 amino acid
193 replacements, a replacement rate of 26.4%. The number of alleles detected at the seven loci
194 ranged from 9 to 21; the majority of alleles differ in amino acid sequence as well as nucleotide

195 sequence. The extent of single nucleotide polymorphism (SNP) per locus among the 120
196 isolates in the sample population is indicated by the nucleotide diversity (π_p). The allelic
197 diversity (H_d) reflects the probability that any pair of isolates drawn from the sample
198 population will carry different alleles at a locus; it is a measure of the discrimination power of
199 the locus in the typing system (40). The *arcC* locus exhibits the greatest nucleotide diversity
200 followed by *dnaJ* and *glpF* among the 7 MLST loci; however, *glpF* is superior to *arc* and *dnaJ*
201 with regard to discrimination power. Sequences of the alleles at the each of the 7 loci are
202 available at the PubMLST database (<https://pubmlst.org/schromogenes>).

203 **Table 1. Characterization of allelic sequence variation observed in 120 unique isolates of *S. chromogenes***

Locus	Gene	Sequence length (bp)	No. Alleles	S (η)	Amino Acid Substitutions	Isolates (n=120)		STs (n=46)	
						π_p	Hd	π_s	Hd
<i>arcC</i>	Carbamate kinase	588	21	70 (72)	19	0.01545	0.690	0.01825	0.805
<i>hutU</i>	Urocanate hydase	693	9	17	6	0.00195	0.330	0.00251	0.388
<i>fumC</i>	Fumerate hydratase	636	14	16	4	0.00207	0.569	0.00250	0.698
<i>dnaJ</i>	chaperone protein dnaJ	747	18	48	11	0.00678	0.613	0.00828	0.760
<i>glpF</i>	glycerol uptake facilitator	612	17	30 (31)	8	0.00626	0.836	0.00737	0.871
<i>menF</i>	Isochorismate synthase	597	11	11	5	0.00095	0.341	0.00140	0.456
<i>pta</i>	Phosphate acetyl transferase	690	10	21	4	0.00275	0.591	0.00382	0.654
		4563	46	213 (216)	57	0.00507	0.956	0.00619	-

204 S (η): number of polymorphic sites (number of mutations when different from number of polymorphic sites)205 π_p : nucleotide diversity per site in the population of 120 isolates206 π_s : nucleotide diversity per site in 46 STs

207 Hd: Allelic Diversity

208 A total of 46 distinct Sequence Types (STs) were identified in the sample population; the
209 7-locus allelic profiles of the 46 STs are listed along with their geographic origins in Table 2 and
210 at the PubMLST database (<https://pubmlst.org/schromogenes>). By convention, the allele
211 sequences in the reference strain MU970 were defined as allele 1 with the corresponding 7-
212 locus allelic profile of ST1 for MU970. The average nucleotide diversity for the 46 STs is
213 0.00507; this is somewhat lower than the values 0.0068, 0.0064, and 0.010 for *S. aureus*, *S.*
214 *epidermidis*, and *S. hominis* respectively [25] but higher than the 0.0021 value for *S. carnosus*
215 [28] and much higher than the 0.00035 value for *S. haemolyticus* [24].

216

217 **Table 2. MLST Profiles of 46 STs and Isolate Origins.**

<u>ST</u>	<u>arcC</u>	<u>hutU</u>	<u>fumC</u>	<u>dnaJ</u>	<u>glpF</u>	<u>isoC</u>	<u>pta</u>	<u>N*</u>	<u>Vermont</u>	<u>Wash.</u>	<u>Belgium</u>
ST1	1	1	1	1	1	1	1	18	11	6	1
ST2	1	1	1	1	1	1	3	1	1		
ST3	1	1	1	1	1	1	6	1	1		
ST4	1	1	1	1	1	1	8	1		1	
ST5	1	1	1	1	4	1	1	7	7		
ST6	1	1	1	3	2	1	2	3	1		2
ST7	1	1	1	5	4	1	1	1	1		
ST8	1	1	1	8	1	1	1	1		1	
ST9	1	1	1	11	1	1	1	1		1	
ST10	1	1	5	3	2	1	2	4	4		
ST11	1	2	1	1	1	1	1	4	4		
ST12	1	6	1	1	1	1	3	1		1	
ST13	2	3	3	2	1	1	2	4	4		
ST14	3	1	1	1	3	3	2	1		1	
ST15	3	1	1	1	3	5	2	8	6	2	
ST16	3	1	1	1	3	6	2	1		1	
ST17	3	1	1	1	10	1	2	4		4	
ST18	4	4	4	4	5	4	4	2	1		1
ST19	5	1	1	1	6	5	2	2	2		
ST20	6	1	6	2	7	1	2	2	2		
ST21	7	1	6	1	3	1	5	1	1		
ST22	8	5	14	6	3	8	2	1	1		
ST23	9	1	8	7	8	7	7	1	1		

ST24	10	1	6	1	9	1	5	1	1
ST25	11	7	2	9	11	2	4	1	1
ST26	12	8	7	10	12	1	2	3	3
ST27	13	1	1	1	1	1	1	1	1
ST28	1	1	1	1	2	1	2	11	11
ST29	1	1	1	1	2	9	2	1	1
ST30	1	1	1	3	15	1	2	4	4
ST31	1	1	1	3	17	1	2	1	1
ST32	1	1	6	14	3	1	2	2	2
ST33	1	1	6	14	13	1	2	1	1
ST34	1	1	9	1	1	1	1	1	1
ST35	1	1	10	15	3	13	2	1	1
ST36	3	1	1	12	10	1	2	1	1
ST37	8	8	7	13	3	1	2	2	2
ST38	14	1	6	1	3	1	2	4	4
ST39	15	1	13	16	3	1	2	3	3
ST40	16	1	6	1	2	1	2	1	1
ST41	16	8	17	1	3	12	2	1	1
ST42	17	1	1	18	2	1	2	2	2
ST43	19	9	11	17	14	2	10	3	3
ST44	20	1	6	1	5	1	5	2	2
ST45	21	1	1	3	2	1	2	1	1
ST46	23	1	6	3	16	1	9	2	2

218 *The number of isolates (N) detected for each ST and their geographic origins are indicated in

219 the right hand columns

220

221 ST1 was the most common sequence type observed in the sample population; it was
222 detected in isolates from all three source locations though primarily (16 out of 17) from the two
223 US locales. Only three other STs were found in multiple locales: ST6 and ST18 in Vermont and
224 Belgium and ST15 in both US locales. The remaining 42 STs were detected in only one of the
225 three locales (Table 2).

226 ST1 plus three additional STs (ST28, ST15, & ST5) account for over 1/3 (n=44) of the
227 isolates in the sample population. At the other end of the frequency spectrum, 24 STs were
228 found only as single isolates. The remaining 52 isolates are distributed among 18 STs containing

229 2-4 isolates each. The discrimination power of the 7-locus MLST scheme for strain
230 characterization within the overall population is 95.6% (Table 1). The discrimination power for
231 the individual geographic populations ranged from 90.2% for the Vermont cohort to 93% for
232 the Belgian cohort. This indicates that each of the three sample populations is genetically
233 diverse despite apparent nearly complete genetic isolation from each other.

234 **Population structure and geographic origins**

235 Characterization of population structure using the eBURST algorithm groups STs
236 according to the number of allele differences at the 7 loci; this approach disregards the extent
237 of sequence difference between alleles. Initial analysis at the single locus variant (SLV) level
238 revealed two clonal clusters, one centered on ST1 with 11 satellite STs and the other centered
239 on ST6 with 7 satellite STs; in addition, there were several ST pairs and triplets. The ST6 cluster
240 included ST28, the second most common ST in the population with nearly four times as many
241 isolates than ST6, prompting the question of whether ST28 might be the founder of the cluster.
242 Investigation at double locus variant (DLV) level showed 33 STs connected in a single network
243 with ST28 at the central node with radiations leading to four secondary nodes centering on ST1,
244 ST6, ST15, and ST38 (Fig 1). The 33 STs in this core network account for 96 of the 120 isolates in
245 the sample population. The 13 STs not included in this network are separated from the network
246 and from each other by sequence differences at 3 or more loci.

247

248 **Fig. 1. Population structure of *S. chromogenes* as indicated by eBURST at the double locus**

249 **variant (DLV) level.** Each of the 33 STs in the eBURST network is represented by a box, the size

250 of which corresponds to the number of isolates in the ST. Heavy (black) lines represent single
251 locus variants, light (grey) lines represent DLVs.

252

253 The ST1 and ST6 nodes connect to ST28 directly, ST1 as DLV and ST6 as a SLV. In terms
254 of nucleotide distances, ST28 and ST1 differ at 4 SNP sites (3 in *glpF* and 1 in *pta*); ST28 and ST6
255 differ at 7 SNP sites in *dnaJ* allele 3. The *glpF* and *dnaJ* allele differences are more likely the
256 result of recombination events given that the sequence motifs of both alleles are present in STs
257 within and outside the common network. The ST15 and ST38 nodes, in contrast, connect to
258 ST28 via intermediary DLV STs: ST17 and ST40 respectively. Despite this, the two nodes are
259 relatively close to ST28 in nucleotide distance, differing at 6 SNP sites for ST15 and 7 SNP sites
260 for ST38, again likely involving recombination events.

261 The cluster around ST1 consists entirely of single locus variants (SLVs), each bearing a
262 different single nucleotide substitution. This starburst pattern is indicative of a recent clonal
263 expansion with ST1 as the founder. Of the 11 STs in the ST1 cluster, all but one, ST34, originate
264 from Vermont or Washington farms. The cluster around ST 15 is also primarily associated with
265 isolates from Vermont and Washington farms. Unlike the ST1 cluster, the ST15 cluster consists
266 mostly of DLVs. Despite this, the average nucleotide distance between ST15 and its satellites is
267 2.2.

268 The STs in the ST6 cluster and the ST38 cluster are predominantly of Belgian origin. The
269 ST6 cluster consists of SLVs in which all but one of the linkages involves loci differing at a single
270 SNP site. Again, this starburst pattern is indicative of a recent clonal expansion with ST6 as the
271 founder. The cluster around ST38 contains more DLVs than SLVs. One ST in the ST38 cluster,

272 ST44, separates itself from the other STs in the cluster by differing from them at an average of
273 58 SNP sites.

274 The geographic partitioning of isolates within the core ST network suggests that *S.*
275 *chromogenes* populations are relatively isolated, more so between Belgium and the United
276 States than between Vermont and Washington within the U.S. That Belgium is home to more
277 STs in the core network than either of the other two source locations and that the central node
278 of the core network, ST28, is of Belgian origin suggests a European origin for the dominant
279 populations of *S. chromogenes* found on both sides of the Atlantic. This hypothesis can be
280 tested by characterizing MLST databases representing more geographically diverse sample
281 populations should they be available in the future.

282

283 **Phylogenetic analysis distinguishes core and outlier STs**

284 Phylogenetic analysis based on overall nucleotide sequence variation between the 46
285 STs provides an alternative perspective on the population structure of *S. chromogenes* (Fig 2).
286 As shown in Fig. 2a, this analysis clusters 39 of the 46 STs into one large group with 100%
287 bootstrap support. The remaining 7 STs are placed on separate branches with deeper roots.
288 This topography is maintained when *S. hyicus*, the nearest neighbor species to *S. chromogenes*,
289 is used as an outgroup, indicating the validity of the topography (data not shown).

290

291 **Fig. 2. Phylogenetic tree of the 46 STs of *S. chromogenes*.** The trees were constructed using
292 the concatenated sequences of the seven MLST loci; bootstrap values at indicated at the branch
293 points and the scale bar is in units of nucleotide differences. Fig. 2a characterizes phylogenetic

294 relationships of all 46 STs. Fig. 2b elaborates the relationships among the 39 STs in the large
295 undifferentiated group in Fig.2a. Nodal clusters identified in the eBURST analysis are specified
296 as ST groups, e.g., ST Gp1, ST Gp6, etc. STs differing at 3 or more loci from the eBURST groups
297 are identified TLV+.

298

299 The large cluster (Fig.2b) includes 32 of the 33 STs in the eBURST core network plus 7
300 STs differing at 3 or more loci from those in the eBURST core group (STs 13, 20, 22, 35, 37, 41, &
301 46). The one member of the eBURST core network that placed outside the phylogenetically
302 defined large group was ST44, previously noted as differing substantially at the sequence level
303 from the other STs in the eBURST network. Within the large cluster, only the STs in the ST Gp1
304 and ST Gp6 appear as unified clusters with strong bootstrap support; the remaining STs,
305 including the 7 STs noted above, are interspersed on variably supported branches. The mean
306 pairwise nucleotide distance between the 39 STs in the group is 9.6 (range 1-22), a relatively
307 small increase over the mean distance of 8.1 (range 1-16) between the 32 STs in the eBurst core
308 network. This increase in nucleotide distance is accounted for by the additional sequence
309 variation present in the STs varying at three or more loci compared to the STs in the eBURST
310 network which are single or double locus variants. The conjoining of the 32 STs in the eBURST
311 network with the seven additional STs in the phylogenetic analysis is thus consistent with all 39
312 STs sharing a common genetic lineage that has undergone diversification. This grouping
313 includes 105 of the 120 isolates in the total population set.

314 The placement of the remaining seven STs as outliers to the common core cluster does
315 not reflect meaningful phylogenetic relationships. Rather it is the consequence of these seven

316 STs carrying hypervariable allelic variants at one or more of the MLST loci. Pairwise comparisons
317 of allele sequences at each of the 7 MLST loci show alleles at five loci partition into two classes,
318 one consisting of alleles typically differing at 4 or fewer nucleotides and a second smaller group
319 differing by 10 or more nucleotides from the first group. The alleles in the first group comprise
320 the allelic composition of the 39 STs in the common core cluster and are designated here as
321 common core alleles; alleles in the second group are designated hypervariable (HV). Table 3
322 compares the relationship of the two classes of alleles in terms of the average pairwise
323 nucleotide distances within and between the classes. It is clear the distances between the two
324 classes are substantially greater than the within-class distances. Two loci, *hutU* and *menF*, loci
325 lack HV alleles.
326

327 **Table 3. Comparison of common core alleles and hypervariable (HV) alleles.** Alleles at each locus were partitioned into common
 328 core and hypervariable groups; each group was characterized independently.

Locus	Common Core Alleles			Hypervariable alleles			Average Pairwise Distance. (nt)		
	No. Alleles	SNP sites	a.a. subs.	No. Alleles	SNP Sites	a.a. subs	Core	HV	Core vs. HV
<i>arcC</i>	14	14	11	7	58	8	3.2	20.3	34.5
<i>hutU</i>	6	5	3	3	3	2	1.7	2.0	12.2
<i>fumC</i>	14	16	4	--	--	--	3.6	--	--
<i>dnaJ</i>	14	15	7	4	32	4	4.9	17.0	20.0
<i>glpF</i>	14	16	8	3	21	2	3.3	2.7	17.8
<i>menF</i>	11	11	5	--	--	--	2.7	--	--
<i>pta</i>	7	6	4	3	6	0	1.9	4.0	13.1

330

331 To illustrate the effect of a single HV allele in a MLST profile, STs 26 & 39 have HV alleles
332 only at the *arcC* locus and are clear outliers in the 7-locus phylogeny (Fig. 2a) but a phylogeny
333 built on the six loci excluding *arcC* results in a repositioning of these two STs within the
334 common core cluster (data not shown). The outliers ST44 and ST23 differ from the common
335 core with HV alleles at two and three loci respectively. The remaining three outlier STs (STs 18,
336 25, & 43) have HV alleles at the five loci and fall into a well-supported group with average
337 nucleotide distances of 105.5 to 108.6 separating these three from the 39 STs in the common
338 core cluster. Notably, these three STs also differ significantly from each other with an average
339 pairwise nucleotide difference of 30.7 between them. These three STs represent 6 isolates of
340 which 4 originate from Belgium and one each from Vermont and Washington State.

341 **Evidence of Recombination in *S. chromogenes***

342 Both mutation and recombination are drivers of genetic diversity in bacterial species
343 [36, 41]. Species that undergo very low rates of recombination have population structures
344 characterized by clonal lineages that diversify slowly by the accumulation of point mutations. At
345 the other end of the spectrum, species that undergo frequent recombination can exhibit a level
346 of genetic diversity that complicates phylogenetic analysis and reconstruction of population
347 structure.

348 The pairwise homoplasmy index (PHI) was used to gain an initial assessment of
349 recombination among the concatenated sequences of the 32 STs in the eBurst network, the 39
350 STs in the common core, and the 46 STs in the full data set. No statistically significant evidence
351 of recombination was detected for the core 32 STs ($p=0.80$), but recombination was indicated

352 for the 39 STs in the common core ($p=0.018$) and very strong evidence for recombination was
353 found for the full set of 46 STs ($p<0.0001$). To characterize the distribution of recombination
354 events within and between loci, the “four gametes” test of Hudson and Kaplan [37] was used;
355 this test yields the minimum number of recombination events between SNP positions in the
356 concatenated ST sequences. Detection of recombination between MLST loci is of particular
357 interest for it indicates expansion of genomic diversity beyond that provided by allele sequence
358 variation. For the 32 ST sequences in the eBURST core complex, four inter-locus and no intra-
359 locus recombination events were detected; the inter-locus recombinants were *arcC/fumC*,
360 *fumC/dnaJ*, *dnaJ/glpF*, and *glpF/menF*. Analysis of the 39 ST sequences in the common core
361 group added one more inter-locus recombination event, *menF/arcC*, plus an intra-locus
362 recombination event in *dnaJ*. Analysis of all 46 ST sequences added 10 more within-locus
363 events: 6 in *arcC*, 2 in *fumC*, and 2 in *glpF* for a total of 16 minimum recombination events
364 overall. Analysis of the outlier 7 ST sequences accounted for 11 of these events, the five
365 between loci and six within loci. These findings are consistent in showing that recombination
366 contributes to genetic diversification in *S. chromogenes*, particularly in the STs with HV alleles.

367 To assess the relative contributions of mutation and recombination events at the allele
368 level, allelic sequence changes were surveyed at each locus within the nodal subgroups
369 delineated by eBURST. Alleles differing at a single nucleotide site were scored as mutations
370 whereas alleles differing at multiple nucleotide sites and alleles shared between different clonal
371 subgroups were scored as recombination events [36]. Notably, the defining allelic signature of
372 three of the four nodal subgroups can be attributed to recombination events contributing one
373 or more new alleles to the allelic profile of ST28, the central node. The nodes of the ST1 and

374 ST6 nodal subgroups differ from ST28 by recombined alleles at the *glpF* and *dnaJ* loci
375 respectively. In contrast, the allelic differences between the STs within each nodal cluster are
376 single nucleotide substitutions in keeping with the starburst topologies of these two nodal
377 clusters. The ST15 nodal subgroup differs from ST28 with recombinant alleles at both the *arcC*
378 and *glpF* loci; single nucleotide variants account for the remainder of the variation within this
379 nodal cluster. The nodal subgroup around ST38 presents a different picture. Of the 10 allele
380 changes occurring within the six STs in this subgroup, four can be attributed to recombination
381 and the remaining six to mutation; thus both single site substitutions and recombination events
382 contribute to the differences between STs within the subgroup. Overall, this assessment
383 indicates the ratio of recombination to mutation to be about 8:32 in the 32 STs comprising the
384 eBURST clonal network. Notably, there is only one example of allele sharing between STs in
385 different nodal subgroups in the eBURST network: the variant allele *glpF-3* is shared between
386 multiple STs in nodal subgroups ST15 and ST38. This allele is also shared with multiple STs
387 outside the eBURST clonal network, validating its status as recombinant.

388 In contrast to the predominance of mutation over recombination in the eBURST clonal
389 network, recombination events predominate in the seven STs containing HV alleles. Indeed,
390 that these seven STs are comprised of mixtures of common core and HV alleles is indicative of
391 recombination. Comparison of HV allele sequences at each of the five loci with HV alleles
392 provides an estimated recombination to mutation ratio of 13:5. The phylogenetically
393 supported branch containing ST18, ST25, and ST43 (Fig. 2a) allows direct comparison at the ST
394 sequence level and yields a recombination to mutation ratio of 8:2. The predominance of

395 recombination to mutation among the HV alleles is indicative of the deeper ancestry of these
396 alleles compared to those of the common core.

397

398 **Hypervariable Alleles – Remnants of a Relict Genotype?**

399 The extreme sequence variation in the HV alleles relative to the common core alleles
400 prompts the question of the origin of these alleles. To test the possibility the HV alleles are
401 introgressions from other species, BLAST searches were done querying representative HV
402 alleles against all genomes in the genus *Staphylococcus*; no hits above 80% sequence identity
403 were observed for any species other than *S. chromogenes*. Additionally, both the common core
404 and HV allele sets are equidistant from the corresponding genes in *S. hyicus* reference
405 sequences, consistent with expectation for common ancestry. An alternative hypothesis is that
406 the HV alleles are remnants of a relict genotype. The large average pairwise SNP distances
407 separating HV and common core alleles is indicative of an early time of divergence between the
408 two classes of alleles (Table 3). The hypothesis that the HV alleles are remnants of a lineage
409 older than the common core alleles is supported by the larger number of variant sites per allele
410 for the HV alleles than for the common core alleles (6 vs. 1.04) and the higher average
411 frequency of synonymous site variants in HV alleles than in the common core alleles (86.7% vs.
412 62.6%). Additional support for this hypothesis is the increased incidence of recombination
413 relative to mutation in the HV alleles compared to the common core alleles; sequence variation
414 due to recombination tends to accumulate over time.

415 The MLST profiles of the 7 outlier STs contain both common core and HV alleles, ranging
416 from one to 5 HV alleles in an MLST profile. These mixed profiles are most likely to have arisen
417 via recombination; mutational variation is not a plausible alternative. It is not possible to
418 ascertain from the MLST data alone whether the mixtures are a result of introgression of non-
419 HV alleles into an HV genome or the other way around. However, the apparent recent origin of
420 the common core alleles and the possible relict origin of the HV alleles suggest the mixtures are
421 relatively recent. A more detailed picture of the population history of *S. chromogenes* awaits
422 further study using whole genome sequence data.

423

424 **Conclusions**

425 The MLST scheme described in this paper provides a tool for the differentiation and
426 identification of strains within *S. chromogenes*. With a power of discrimination between strain
427 types exceeding 90% in geographically localized populations and greater than 95% overall, this
428 MLST scheme has potential for use in epidemiological investigations of pathologies associated
429 with this species and the ecological relationships between microbe and host. The geographic
430 distribution of strain types indicated a high degree of genetic isolation between locales, posing
431 a question of the historical and genetic factors accounting for this separation. Phylogenetic
432 analysis of strain types identified by the scheme showed most to be contained within a single
433 large and genetically diversified lineage which included strains arising from mutation driven
434 clonal expansions and more varied strains generated by recombination events. The MLST
435 analysis also revealed that some strain types were differentiated by having alleles with highly
436 variable sequences at one or more of the loci in the 7-locus MLST scheme; these highly variable

437 alleles were posited to be remnants of a relic genotype of *S. chromogenes*. These features of
438 the population structure of this species provide a prospectus for future studies.

439

440 **Acknowledgments**

441 We acknowledge the work of the staff The Vermont Integrative Genomics Resource DNA
442 Facility, who completed the automated DNA sequencing.

443

444 **Data Availability**

445 Sequences for the alleles and isolates from this study are available at
446 <https://pubmlst.org/schromogenes>

447

448 **References**

- 449 1. Devriese LA, Hajek V, Oeding P, Meyer SA, Schleifer KH. *Staphylococcus hyicus* (Sompolinsky, 1953)
450 comb. nov. and *Staphylococcus hyicus* subsp. *chromogenes* subsp. nov. Int. J. Syst. Bacteriol. 1978;
451 28: 482-490. doi: 10.1099/00207713-28-4-482
- 452 2. Hajek V, Devriese LA, Mordarski M, Goodfellow M, Pulverer G, Varaldo PE. Elevation of
453 *Staphylococcus hyicus* subsp. *chromogenes* (Devriese et al., 1978) to Species Status: *Staphylococcus*
454 *chromogenes* (Devriese et al., 1978) comb. nov. System. Appl. Microbiol. 1986; 8: 169-173. doi:
455 10.1016/S0723-2020(86)80071-6

- 456 3. Lamers RP, Muthukrishnan G, Castoe TA, Tafur S, Cole AM, Parkinson CL. Phylogenetic relationships
457 among *Staphylococcus* species and refinement of cluster groups based on multilocus data. BMC
458 Evol Biol. 2012; 12:171. doi: 10.1186/1471-2148-12-171.
- 459 4. Naushad S, Barkema HW, Luby C, Condas LA, Nobrega DB, Carson DA, De Buck J. Comprehensive
460 Phylogenetic Analysis of Bovine Non-*aureus* Staphylococci Species Based on Whole-Genome
461 Sequencing. Front Microbiol. 2016; 7:1990. doi: 10.3389/fmicb.2016.01990.
- 462 5. Vanderhaeghen W, Piepers S, Leroy F, Van Coillie E, Haesebrouck F, De Vliegher S. Identification,
463 typing, ecology and epidemiology of coagulase-negative staphylococci associated with ruminants.
464 The Veterinary Journal. 2015; 203: 44-51. doi: 10.1016/j.tvjl.2014.11.001.
- 465 6. Foster AP. Staphylococcal skin disease in livestock. Vet Dermatol. 2012; 23(4):342-51, e63. doi:
466 10.1111/j.1365-3164.2012.01093.x.
- 467 7. Andresen LO, Ahrens P, Daugaard L, Bille-Hansen V. Exudative epidermitis in pigs caused by
468 toxigenic *Staphylococcus chromogenes*. Vet Microbiol. 2005; 105(3-4):291-300. doi:
469 10.1016/j.vetmic.2004.12.006.
- 470 8. Gosselin VB, Lovstad J, Dufour S, Adkins PRF, Middleton JR. Use of MALDI-TOF to characterize
471 staphylococcal intramammary infections in dairy goats. J Dairy Sci. 2018; 101(7):6262-6270. doi:
472 10.3168/jds.2017-14224. 9.
- 473 9. Schmidt T, Kock MM, Ehlers MM. Diversity and antimicrobial susceptibility profiling of
474 staphylococci isolated from bovine mastitis cases and close human contacts. J Dairy Sci. 2015;
475 98:6256-6269. doi: 10.3168/jds.2015-9715.
- 476 10. Mørk T, Jørgensen HJ, Sunde M, Kvitle B, Sviland S, Waage S, Tollersrud T. Persistence of
477 staphylococcal species and genotypes in the bovine udder. Vet Microbiol. 2012; 159(1-2):171-80.
478 doi: 10.1016/j.vetmic.2012.03.034.

- 479 11. Bexiga R, Rato MG, Lemsaddek A, Semedo-Lemsaddek T, Carneiro C, Pereira H, Mellor DJ, Ellis KA,
480 Villela CL. Dynamics of bovine intramammary infections due to coagulase-negative staphylococci on
481 four farms. *J Dairy Res.* 2014; 81:208-214. doi: 10.1017/S0022029914000041.
- 482 12. Fry PR, Calcutt MJ, Foecking MF, Hsieh HY, Suntrup DG, Perry J, Stewart GC, Middleton JR. Draft
483 Genome Sequence of *Staphylococcus chromogenes* Strain MU 970, Isolated from a Case of Chronic
484 Bovine Mastitis. *Genome Announc.* 2014; 2(4): e00835-14. doi: 10.1128/genomeA.00835-14.
- 485 13. Supré K, Haesebrouck F, Zadoks RN, Vanechoutte M, Piepers S, De Vlieghe S. Some coagulase-
486 negative *Staphylococcus* species affect udder health more than others. *J Dairy Sci.* 2011;
487 94(5):2329-40. doi: 10.3168/jds.2010-3741.
- 488 14. De Visscher A, Piepers S, Haesebrouck F, De Vlieghe S. Intramammary infection with coagulase-
489 negative staphylococci at parturition: species-specific prevalence, risk factors, and effect on udder
490 health. *J of Dairy Sci.* 2015; 99:6457-6469. doi: 10.3168/jds.2015-10458.
- 491 15. Taponen S, Björkrot J, Pyörälä S. Coagulase-negative staphylococci isolated from bovine
492 extramammary sites and intramammary infections in a single dairy herd. *J Dairy Res.* 2008; 75:422-
493 429. doi: 10.1017/S0022029908003312.
- 494 16. White DG, Harmon RJ, Matos JE, Langlois BE. Isolation and identification of coagulase-negative
495 *Staphylococcus* species from bovine body sites and streak canals of nulliparous heifers. *J Dairy Sci.*
496 1989; 2(7):1886-92. doi: 10.3168/jds.S0022-0302(89)79307-3.
- 497 17. De Visscher A, Piepers S, Haesebrouck F, De Vlieghe S. Teat apex colonization with coagulase-
498 negative *Staphylococcus* species before parturition: distribution and species-specific risk factors. *J*
499 *Dairy Sci.* 2016; 99:1427-1439. doi: 10.3168/jds.2015-10326.
- 500 18. Adkins PRF, Dufour S, Spain JN, Calcutt MJ, Reilly TJ, Stewart GC, Middleton JR. Molecular
501 characterization of non-aureus *Staphylococcus* spp. from heifer intramammary infections and body
502 sites. *J Dairy Sci.* 2018; 101:5388-5403. doi: 10.3168/jds.2017-13910.

- 503 19. Piessens V, Van Coillie E, Verbist B, Supré K, Braem G, Van Nuffel A, De Vuyst L, Heyndrickx M, De
504 Vlieghe S. Distribution of coagulase-negative *Staphylococcus* species from milk and environment
505 of dairy cows differs between herds. *J Dairy Sci.* 2011; 94(6):2933-44. doi: 10.3168/jds.2010-3956.
- 506 20. Matthews KR, Harmon RJ, Smith BA. Protective effect of *Staphylococcus chromogenes* infection
507 against *Staphylococcus aureus* infection in the lactating bovine mammary gland. *J Dairy Sci.* 1990;
508 73(12):3457-62. doi: 10.3168/jds.S0022-0302(90)79044-3.
- 509 21. De Vlieghe S, Opsomer G, Vanrolleghem A, Devriese LA, Sampimon OC, Sol J, Barkema HW,
510 Haesebrouck F, de Kruif A. In vitro growth inhibition of major mastitis pathogens by
511 *Staphylococcus chromogenes* originating from teat apices of dairy heifers. *Vet Microbiol.* 2004;
512 101(3):215-21. doi: 10.1016/j.vetmic.2004.03.020.
- 513 22. Enright MC, Day NP, Davies CE, Peacock SJ, Spratt BG. Multilocus sequence typing for
514 characterization of methicillin-resistant and methicillin susceptible clones of *Staphylococcus*
515 *aureus*. *J Clin Microbiol.* 2000; 38: 1008–1015.
- 516 23. Thomas JC, Vargas MR, Miragaia M, Peacock SJ, Archer GL, et al. Improved multilocus sequence
517 typing scheme for *Staphylococcus epidermidis*. *J Clin Microbiol.* 2007; 45: 616–619. doi:
518 10.1128/JCM.01934-06.
- 519 24. Cavanagh JP, Klingenberg C, Hanssen A-M, Fredheim EA, Francois P, et al. Core genome
520 conservation of *Staphylococcus haemolyticus* limits sequence based population structure analysis. *J*
521 *Microbiol Methods.* 2012; 89: 159–166. doi: 10.1016/j.mimet.2012.03.014.
- 522 25. Zhang L, Thomas JC, Miragaia M, Bouchami O, Chaves F, et al. Multilocus Sequence Typing and
523 Further Genetic Characterization of the Enigmatic Pathogen, *Staphylococcus hominis*. *PLoS ONE.*
524 2013; 8(6): e66496. doi: 10.1371/journal.pone.0066496.

- 525 26. Chassain B, Lemee L, Didi J, Thiberge JM, Brisse S, Pons JL, Pestel-Caron M. Multilocus sequence
526 typing analysis of *Staphylococcus lugdunensis* implies a clonal population structure. J Clin Microbiol.
527 2012; 50: 3003–3009. doi: 10.1128/JCM.00988-12.
- 528 27. Solyman SM, Black CC, Duim B, Perreten V, van Duijkeren E, Wagenaar JA, Eberlein LC, Sadeghi LN,
529 Videla R, Bemis DA, Kania SA. Multilocus sequence typing for characterization of *Staphylococcus*
530 *pseudintermedius*. J Clin Microbiol. 2013; 51(1):306-10. doi: 10.1128/JCM.02421-12.
- 531 28. Bückle A, Kranz M, Schmidt H, Weiss A. Genetic diversity and population structure of food-borne
532 *Staphylococcus carnosus* strains. Syst Appl Microbiol. 2017; 40(1):34-41. doi:
533 10.1016/j.syapm.2016.11.005.
- 534 29. Park JY, Fox LK, Seo KS, McGuire MA, Park YH, Rurangirwa FR, Sischo WM, Bohach GA. Comparison
535 of phenotypic and genotypic methods for the species identification of coagulase-negative
536 staphylococcal isolates from bovine intramammary infections. Vet Microbiol. 2011; 147(1-2):142-8.
537 doi: 10.1016/j.vetmic.2010.06.020.
- 538 30. Drancourt M, Raoult D. rpoB gene sequence-based identification of *Staphylococcus* species. J. Clin.
539 Microbiol. 2002; 40(4):1333-1338. doi: 10.1128/jcm.40.4.1333-1338.2002.
- 540 31. Heikens E, Fleer A, Paauw A, Florijn A, Fluit AC. Comparison of genotypic and phenotypic methods
541 for species-level identification of clinical isolates of coagulase-negative staphylococci. J. Clin.
542 Microbiol. 2005; 43(5):2286-2290. doi: 10.1128/JCM.43.5.2286-2290.2005.
- 543 32. Tamura K, Stecher G, Peterson D, Filipski A, Kumar S. MEGA6: Molecular evolutionary genetics
544 analysis version 6.0. Molecular Biology and Evolution. 2013; 30(12), 2725-2729. doi:
545 10.1093/molbev/mst197.
- 546 33. Librado P, Rozas J. DnaSP v5: A software for comprehensive analysis of DNA polymorphism data.
547 Bioinformatics. 2009; 25: 1451-1452. doi: 10.1093/bioinformatics/btp187.

- 548 34. Feil EJ, Li BC, Aanensen DM, Hanage WP, Spratt BG. eBURST: inferring patterns of evolutionary
549 descent among clusters of related bacterial genotypes from multilocus sequence typing data. J
550 Bacteriol. 2004; 186(5):1518-30. doi: 10.1128/jb.186.5.1518-1530.2004.
- 551 35. Francisco AP, Bugalho M, Ramirez M, Carrico JA. Global Optimal eBURST analysis of Multilocus
552 typing data using a graphic matroid approach. BMC Bioinformatics. 2009; 10:152. doi:
553 10.1186/1471-2105-10-152.
- 554 36. Feil EJ, Enright MC, Spratt BG. Estimating the relative contributions of mutation and recombination
555 to clonal diversification: a comparison between *Neisseria meningitidis* and *Streptococcus*
556 *pneumoniae*. Res. Microbiol. 2000; 151:465-469. doi: 10.1016/s0923-2508(00)00168-6.
- 557 37. Hudson, R. R. and N. L. Kaplan. Statistical properties of the number of recombination events in the
558 history of a sample of DNA sequences. Genetics. 1985; 111: 147-164.
- 559 38. Huson DH and D Bryant. Application of Phylogenetic Networks in Evolutionary Studies. Mol. Biol.
560 Evol. 2006; 23(2):254-267.
- 561 39. Bruen TC, Philippe H, Bryant D. A simple and robust statistical test for detecting the presence of
562 recombination. Genetics. 2006; 172(4):2665-81. doi: 10.1534/genetics.105.048975.
- 563 40. Hunter PR, Gaston MA. Numerical index of the discriminatory ability of typing systems: an
564 application of Simpson's index of diversity. J Clin Microbiol. 1988; 26(11):2465-6.
- 565 41. Smith JM, Feil EJ, Smith NH. Population structure and evolutionary dynamics of pathogenic
566 bacteria. Bioessays. 2000; 22(12):1115-22. doi: 10.1002/1521-1878(200012)22:12<1115::AID-
567 BIES9>3.0.CO;2-R.
- 568

569 **Supporting Information**570 **S1 Table. Genetic loci and primer sequences used in MLST scheme**

571

Locus	Contig	Locus Tag	Gene	PCR Primers (5'→3')	PCR Product (bp)	Gene sequence segment used for MLST (length)
<i>arcC</i>	11	SCHR_09950	Carbamate kinase	F: CGGCGATTCGACAAACTC R: TGGCAACATCGACCCTTCTG	746	133 - 720 (588 bp)
<i>hutU</i>	4	SCHR_06535	Urocanate hydase	F: AAGGGGTTGTCATCGGTGTA R: GCATCGGAACCGTCTTTCAT	829	616 – 1308 (693 bp)
<i>fumC</i>	16	SCHR_11020	Fumerate hydratase	F: TGCATGTGCGACTATATCAC R: CATCAATATGTTCTCAATCG	756	496 – 1131 (636 bp)
<i>dnaJ</i>	3	SCHR_04712	chaperone protein dnaJ	F: AAAGGGAGCGATAGCATTGG R: CATCACCTAACGCAGCTTGT	869	46 – 792 (747 bp)
<i>glpF</i>	2	SCHR_03515	glycerol uptake facilitator	F: TACGGTTAGGCAAGGAGTCT R: AACGACCTGGTAGGCCAAT	759	25 – 636 (612 bp)
<i>menF</i>	1	SCHR_00545	Isochorismate synthase	F: TGTCACACCTGAAGAACAACA R: TAACGCTTGGTTACCTTGAATC	730	592 – 1188 (597 bp)
<i>pta</i>	1	SCHR_02435	Phosphate acetyl transferase	F: AACGCCCCCTTGAAAAGTC R: TGGATTTAGCGCCCGGTG	870	13 – 702 (690 bp)

Fig 1

bioRxiv preprint doi: <https://doi.org/10.1101/2020.11.30.403683>; this version posted November 30, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

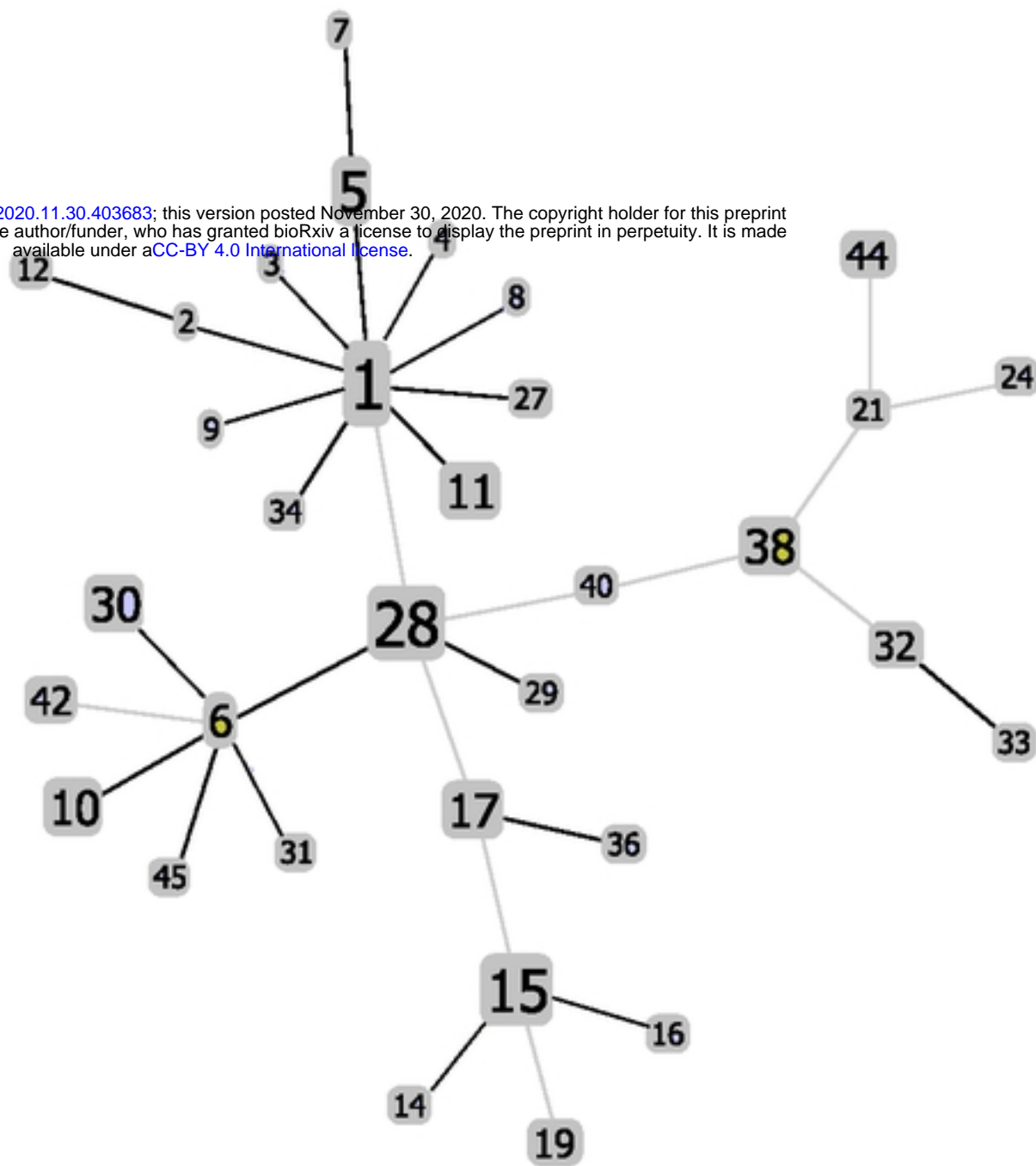


Fig 2a

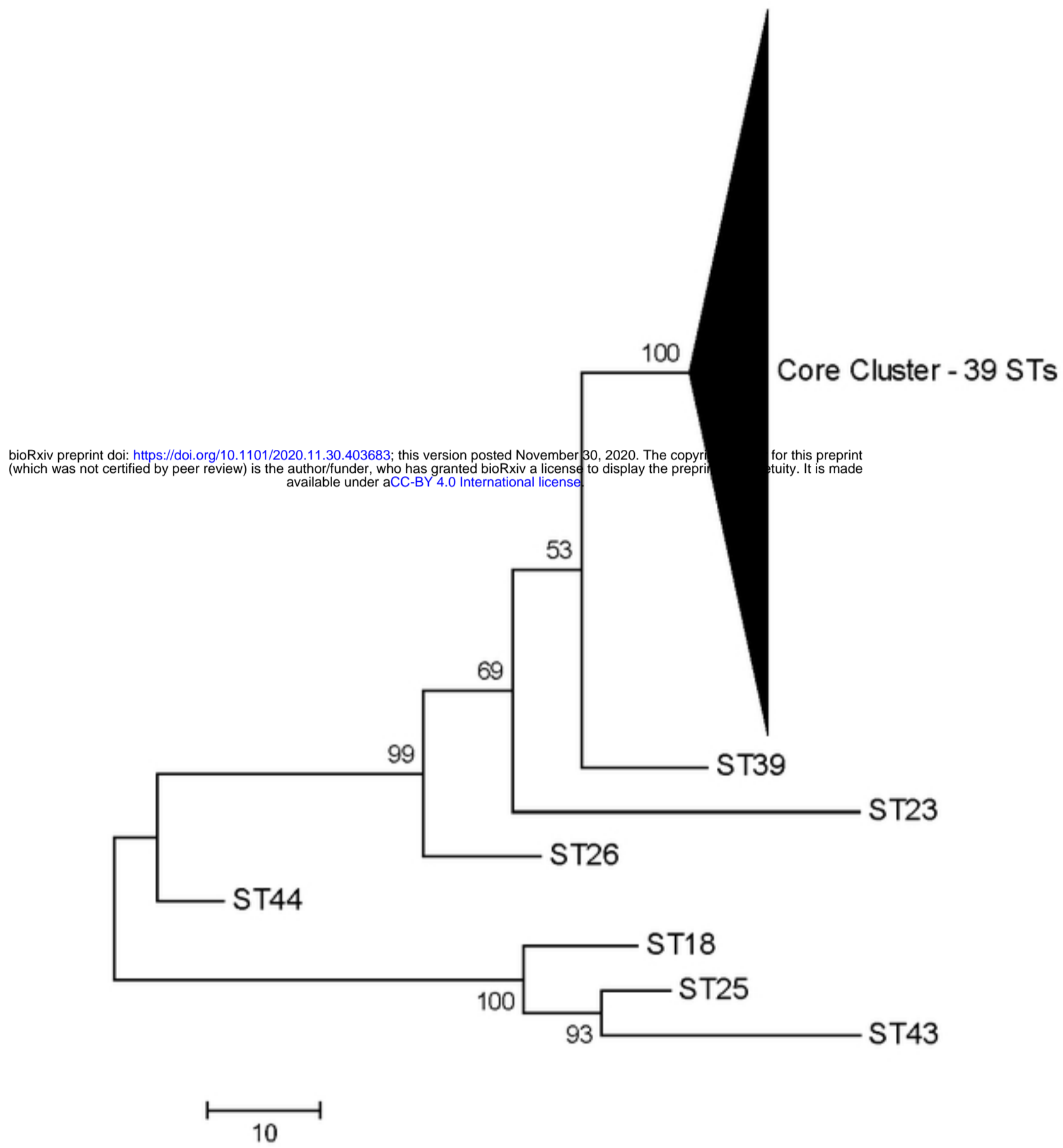


Fig 2b

