

## Strathprints Institutional Repository

Holt, K. and Baker, Stephen and Weill, François-Xavier and Holmes, Edward C. and Kitchen, Andrew and Yu, Jun and Sangal, Vartul and Brown, Derek J. and Coia, John E. and Kim, Dong Wook and Choi, Seon Young and Kim, Su Hee and da Silveira, Wanderley D. and Pickard, Derek J. and Farrar, Jeremy J. and Parkhill, Julian and Dougan, Gordon and Thomson, Nicholas R. (2012) *Shigella sonnei* genome sequencing and phylogenetic analysis indicate recent global dissemination from Europe. *Nature Genetics*, 44 (9). 1056–1059. ISSN 1061-4036

Strathprints is designed to allow users to access the research output of the University of Strathclyde. Copyright © and Moral Rights for the papers on this site are retained by the individual authors and/or other copyright owners. You may not engage in further distribution of the material for any profitmaking activities or any commercial gain. You may freely distribute both the url (<http://strathprints.strath.ac.uk/>) and the content of this paper for research or study, educational, or not-for-profit purposes without prior permission or charge.

Any correspondence concerning this service should be sent to Strathprints administrator: <mailto:strathprints@strath.ac.uk>



## Manuscript Information

**Journal name:** Nature genetics

**Manuscript #:** 49152

**Manuscript Title:** Shigella sonnei genome sequencing and phylogenetic analysis indicate recent global dissemination from Europe

**Principal Investigator:** J Yu (jun.yu@strath.ac.uk)

**Submitter:** Nature Publishing Group (repositorynotifs@nature.com)

## Grant/Project/Contract/Support Information

Name	Support ID#	Title
J Yu	MRC_G0800173 (86345)	Genotyping epidemic Shigella sonnei

## Manuscript Files

Type	Fig/Table #	Filename	Size	Uploaded
manuscript	1	article_1.doc	276480	2012-07-06 08:00:36
figure	1	figure_1.ai	3313183	2012-07-06 08:00:53
supplement	1	supp_info_1.pdf	7713830	2012-07-06 08:01:30
supplement	table 1	supp table 1.xls	167936	2012-07-09 13:31:40

This PDF receipt will only be used as the basis for generating UK PubMed Central (UKPMC) documents. UKPMC documents will be made available for review after conversion (approx. 2-3 weeks time). Any corrections that need to be made will be done at that time. No materials will be released to UKPMC without the approval of an author. Only the UKPMC documents will appear on UK PubMed Central -- this PDF Receipt will not appear on UK PubMed Central.

1 ***Shigella sonnei* genome sequencing and phylogenetic analysis indicate recent**  
2 **global dissemination from Europe**

3 Kathryn E. Holt<sup>1</sup>, Stephen Baker<sup>2</sup>, François-Xavier Weill<sup>3</sup>, Edward C. Holmes<sup>4,5</sup>,  
4 Andrew Kitchen<sup>4</sup>, Jun Yu<sup>6</sup>, Vartul Sangal<sup>6</sup>, Derek J. Brown<sup>7</sup>, John E. Coia<sup>7</sup>, Dong  
5 Wook Kim<sup>8,9</sup>, Seon Young Choi<sup>8</sup>, Su Hee Kim<sup>8</sup>, Wanderley D. da Silveira<sup>10</sup>, Derek J.  
6 Pickard<sup>11</sup>, Jeremy J. Farrar<sup>2</sup>, Julian Parkhill<sup>11</sup>, Gordon Dougan<sup>11</sup>, Nicholas R.  
7 Thomson<sup>11</sup>

- 8 1. University of Melbourne, Department of Microbiology and Immunology, Royal Parade,  
9 Melbourne, Victoria, 3010, Australia
- 10 2. The Hospital for Tropical Diseases, Wellcome Trust Major Overseas Programme, Oxford  
11 University Clinical Research Unit, Ho Chi Minh City, Vietnam
- 12 3. Institut Pasteur, Unité des Bactéries Pathogènes Entériques, Paris, France
- 13 4. Center for Infectious Disease Dynamics, Department of Biology, The Pennsylvania State  
14 University, University Park, PA 16802, USA
- 15 5. Fogarty International Center, National Institutes of Health, Bethesda, MD 20892, USA
- 16 6. Strathclyde Institute of Pharmacy and Biomedical Sciences, University of Strathclyde,  
17 Glasgow, G4 0RE, UK
- 18 7. Scottish *Salmonella*, *Shigella* and *Clostridium difficile* Reference Laboratory, Stobhill  
19 Hospital, 133 Balornock Road, Glasgow, UK
- 20 8. Molecular Biology Laboratory, International Vaccine Institute (IVI), Seoul, Republic of  
21 Korea
- 22 9. Department of Pharmacy, College of Pharmacy, Hanyang University  
23 Ansan, Kyeonggi- do, 426-791, Korea
- 24 10. Department of Genetics, Evolution and Bioagents, Biology Institute, Campinas State  
25 University – UNICAMP, Brazil
- 26 11. Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK  
27 CB10 1SA

28

29 Abstract/First paragraph

30 *Shigella* are human-adapted *Escherichia coli* that have gained the ability to invade the  
31 human gut mucosa and cause dysentery<sup>1,2</sup>, spreading efficiently via low-dose fecal-  
32 oral transmission<sup>3,4</sup>. Historically, *S. sonnei* has been predominantly responsible for  
33 dysentery in developed countries, but is now emerging as a problem in the developing  
34 world, apparently replacing the more diverse *S. flexneri* in areas undergoing economic  
35 development and improvements in water quality<sup>4-6</sup>. Classical approaches have shown  
36 *S. sonnei* is genetically conserved and clonal<sup>7</sup>. We report here whole-genome  
37 sequencing of 132 globally-distributed isolates. Our phylogenetic analysis shows that  
38 the current *S. sonnei* population descends from a common ancestor that existed less  
39 than 500 years ago and has diversified into several distinct lineages with unique  
40 characteristics. Our analysis suggests the majority of this diversification occurred in  
41 Europe, followed by more recent establishment of local pathogen populations in other  
42 continents predominantly due to the pandemic spread of a single, rapidly-evolving,  
43 multidrug resistant lineage.

44

45

46 To establish an accurate population framework we sequenced the whole genomes of  
47 132 *S. sonnei* isolated between 1943 and 2008, spanning four continents  
48 (Supplementary Table 1). We detected 10,111 chromosomal single nucleotide  
49 polymorphisms (SNPs) randomly distributed around the *S. sonnei* chromosome,  
50 approximately one per 430 bp (0.23% nucleotide divergence) (Supplementary Fig. 1).  
51 To investigate the population structure of *S. sonnei*, we analysed these chromosomal  
52 SNPs using multiple phylogenetic methods. Maximum likelihood (ML) phylogenetic  
53 analysis (Supplementary Fig. 2) revealed a strong correlation between root-to-tip

54 branch lengths and the known dates of isolation for the sequenced *S. sonnei*,  
55 indicative of rapid, clock-like evolution (Supplementary Fig. 3). There appears to be  
56 some rate variation between lineages, possibly associated with differences in effective  
57 population size or in the mean number of generations per year (replication rate),  
58 which may in turn be associated with different lifestyles or niches. We used a  
59 Bayesian approach (BEAST<sup>8</sup>) to infer the evolutionary dynamics of the global *S.*  
60 *sonnei* population as a whole. Importantly, this yielded the same tree topology as the  
61 ML analysis, while also providing estimates of nucleotide substitution rates and  
62 divergence times for key *S. sonnei* lineages (Fig. 1). Interestingly, the phylogenies  
63 identified four distinct *S. sonnei* lineages, three encompassing isolates spanning the  
64 1940s through the 2000s and another comprising a single isolate from France (Fig. 1).  
65 These lineages each had 100% ML bootstrap support, 100% Bayesian posterior  
66 support (BEAST) and were also recovered using a Bayesian clustering analysis (see  
67 Online Methods). Whilst these lineages are uniquely characterized by hundreds of  
68 SNPs they display only minor differences in gene content and were correlated with  
69 traditional typing methods used to subdivide *S. sonnei* (biotypes a-g<sup>9</sup> and CRISPR  
70 types<sup>10</sup>) (Supplementary Note, Supplementary Fig. 2, Supplementary Table 3). We  
71 estimated a mean substitution rate of  $2.0 \times 10^{-4}$  site<sup>-1</sup> year<sup>-1</sup> among the 10,111  
72 chromosomal SNP loci [95% Highest Posterior Density (HPD)  $1.6 \times 10^{-4} - 2.3 \times 10^{-}$   
73 <sup>4</sup>], corresponding to the accumulation of approximately 2.2 SNPs chromosome<sup>-1</sup> year<sup>-</sup>  
74 <sup>1</sup> ([95% HPD 1.8 – 2.6], excluding repeated and phage regions). This scales to a  
75 genome-wide substitution rate of  $6.0 \times 10^{-7}$  substitutions site<sup>-1</sup> year<sup>-1</sup> [95% HPD =  $5.2$   
76  $\times 10^{-7} - 6.7 \times 10^{-7}$ ], which likely represents the upper bound of the true genome-wide  
77 substitution rate and is similar to that calculated for the enteric pathogen *Vibrio*  
78 *cholerae* ( $8 \times 10^{-7}$  site<sup>-1</sup> year<sup>-1</sup>)<sup>11</sup> but lies between the rates estimated for *Yersinia*

79 *pestis* ( $2 \times 10^{-8}$ )<sup>12</sup> and *Staphylococcus aureus* ( $3 \times 10^{-6}$ )<sup>13</sup>. From BEAST analysis, we  
80 estimated the most recent common ancestor (MRCA) of all contemporary *S. sonnei*  
81 existed less than 500 years ago [median calendar year for divergence date, 1669; 95%  
82 HPD, 1554 - 1763] (Fig. 1). Similarly, we estimate the MRCA for each of Lineages I  
83 and II existed in the early 19<sup>th</sup> century and that all Lineage III isolates descend from a  
84 hypothetical ancestor that existed around the turn of the 20<sup>th</sup> century (Fig. 1).  
85 Critically, these data indicate that though the extant *S. sonnei* population descends  
86 from a single ancestor existing in the 17<sup>th</sup> century, by the late 19<sup>th</sup> century *S. sonnei*  
87 had become segregated into at least four distinct lineages that still persist today.  
88  
89 There was strong evidence for regional clustering of *S. sonnei* within the phylogenetic  
90 tree (Fig. 1), indicating significant geographic structure in the global bacterial  
91 population ( $p < 1 \times 10^{-5}$  for association between phylogeny and geographic region<sup>14</sup>).  
92 Interestingly, the European population shows the richest diversity, with isolates  
93 distributed across all four lineages (31% lineage I, 35% lineage II, 31% lineage III,  
94 sole lineage IV isolate) and occupying basal branches in each lineage (Fig. 1). In  
95 contrast, *S. sonnei* isolates from Asia, Africa and America were mainly from lineage  
96 III (67-77%) with fewer lineage II representatives (22-26%) and just two from  
97 Lineage I. Furthermore, ancestral state reconstruction analysis indicated a >50%  
98 likelihood of a European common ancestor for each of the lineages I, II and III (Fig.  
99 1). The data also indicate Lineage III has been more successful at global dispersal  
100 than other lineages, with only low numbers of Lineage I or II detected outside Europe  
101 (Fig. 1). In particular, a recently derived clade within Lineage III (Global III, MRCA  
102 = 1972 [95% HPD = 1964-1979 C.E.]) has been particularly successful at global  
103 dissemination, comprising 49% of all isolates sampled since 1995 and detected in all

104 regions represented in our collection (Fig. 1). Unlike the European isolates, isolates  
105 from non-European countries form tight shallow-rooted phylogenetic clusters,  
106 consistent with and suggestive of contemporary dispersal (Fig. 1). In many cases,  
107 these clusters contain multiple isolates from the same country, indicating localized  
108 clonal expansions (Fig. 1). For example, isolates from Korea formed two subclades  
109 within lineages II and III that likely represent separate introductions of *S. sonnei* into  
110 Korea during the 1960s and 1970s, each followed by local clonal expansions (Fig. 1).  
111 Similarly, isolates originating in Vietnam form two subclades, indicating the local  
112 establishment of Lineage III clones in Vietnam in the 1990s (Fig. 1). At a regional  
113 level, there appears to have been an establishment of a Lineage III subclade in South  
114 America during the 1950s to which isolates from Brazil and Peru could be traced,  
115 followed by dissemination of the Global III clade into Africa and America in the early  
116 1980s (Fig. 1).

117

118 Critically, the phylogeographic analysis indicates that all contemporary *S. sonnei*  
119 infections are caused by a small number of clones that have recently become globally  
120 dispersed (Fig. 1). The distribution of antimicrobial resistance genes and mutations  
121 within the *S. sonnei* phylogeny suggest that selection for multiple drug resistance  
122 (MDR) played a pivotal role in driving this global dissemination (Fig. 1,  
123 Supplementary Fig. 2, Supplementary Table 1). In particular, the establishment of  
124 local *S. sonnei* Lineage III populations outside Europe is intimately associated with  
125 the carriage of transposon Tn7 and class II integrons (In2) encoding resistance to  
126 multiple antimicrobials (Fig. 1). All three major Lineage III subgroups carry a distinct  
127 In2 variant, which is either plasmid-encoded (South America III) or integrated into  
128 the chromosome adjacent to *glmS* (Central Asia IIIa, Global III), suggesting

129 independent acquisitions of the integron in each group during the 1960s-1970s  
130 followed by clonal expansion and subsequent international spread (Fig. 1). Studies  
131 from Europe, Asia, Africa, South America and Australia have reported a high  
132 prevalence of In2-bearing, MDR, biotype g *S. sonnei*, often associated with local  
133 epidemics<sup>15</sup>. Our data demonstrate biotype g is a marker for Lineage III due to a  
134 conserved nonsense mutation in rhamnose regulatory gene *rhaR* (Supplementary Fig.  
135 2) and indicate that the global distribution of MDR biotype g/In2 *S. sonnei* is the  
136 result of global dissemination of multiple In2-bearing subclades of Lineage III *S.*  
137 *sonnei*. Half of the In2-bearing Lineage III isolates also harboured the small MDR  
138 plasmid spA<sup>2</sup> containing *tetAR*, *strAB* and *sul2* genes, which confer additional  
139 resistance to tetracycline, streptomycin and sulfonamides (Fig. 1). All quinolone  
140 resistant isolates harboured one of three point mutations in the chromosomal DNA  
141 gyrase gene, *gyrA*, known to confer quinolone resistance (Fig. 1, Supplementary  
142 Table 1; we detected no plasmid-mediated quinolone resistance genes). The  
143 distribution of *gyrA* mutations within the phylogeny shows these resistance mutations  
144 have arisen independently on at least nine occasions among our *S. sonnei* collection,  
145 including two separate mutations within the clonal group Korea II, indicative of  
146 surprisingly strong selection for quinolone resistance even among MDR isolates (Fig.  
147 1). To investigate other signals of selection, we examined the clustering of SNPs  
148 within genes and chromosomal regions (Supplementary Note). We found evidence of  
149 phage and transposase insertions and a single case of homologous recombination  
150 affecting the *sitABCD* operon in isolate 31382, but identified only two genes  
151 displaying amino acid variation significantly higher than expected under a random  
152 distribution of SNPs. Neither of these genes (*rpoS* and *mreB*) encodes an extracellular  
153 protein, suggesting a lack of immune selection, in common with another human



154 restricted pathogen *Salmonella* Typhi (typhoid fever)<sup>16</sup>. However, we detected a large  
155 number of nonsynonymous SNPs (nsSNPs) and a high rate of nonsynonymous to  
156 synonymous substitutions per site ( $d_N/d_S$ ) in the drug efflux pump component genes  
157 *acrD* (8 nsSNPs,  $d_N/d_S = 2.5$ ) and *acrB* (12 nsSNPs,  $d_N/d_S = 1.8$ ). Currently,  
158 antimicrobial treatment is recommended for the management of dysentery<sup>17</sup>, but may  
159 not significantly impact the resolution of *S. sonnei* or *S. flexneri* infections<sup>18,19</sup>.  
160 However, there is evidence such treatment can prevent shedding of *S. sonnei* after the  
161 resolution of symptoms<sup>20</sup>. Thus, while antimicrobial resistance may have only minor  
162 implications for dysentery treatment, this phenotype may be important in sustaining *S.*  
163 *sonnei* transmission within human populations and our data indicates there is a strong  
164 selective pressure for its maintenance. It has been hypothesized that free-living  
165 amoebae may represent an environmental reservoir for *Shigella*, which are able to  
166 survive intracellularly within *Acanthamoeba*<sup>21,22</sup>. This could potentially provide  
167 another niche in which selective pressure for antibiotic resistance may be exerted,  
168 although intracellular *Shigella* are likely to be protected from most antibiotics by their  
169 amoebae hosts<sup>23,24</sup>.

170

171 Previous studies have proposed that the acquisition of virulence plasmid pINV B,  
172 encoding the *Plesiomonas shigelloides* related O antigen, was the defining event in  
173 the emergence of *S. sonnei*<sup>25</sup>. Unfortunately, the *S. sonnei* virulence plasmid is highly  
174 unstable on laboratory media and is commonly lost on sub-culturing<sup>26</sup> and, as a  
175 consequence, less than half of our isolates yielded sufficient virulence plasmid  
176 sequence data for analysis (46 isolates with >10x read depth). Phylogenetic analysis  
177 of the available virulence plasmid sequences (which contained 84 SNPs) identified

178 three distinct lineages (Supplementary Fig. 4). There was a parallel relationship  
179 between chromosomal and plasmid lineages, consistent with co-evolution of the  
180 plasmid and host chromosome, stable maintenance of the plasmid in the natural  
181 environment and no transfer of plasmid variants among host bacteria. It has also been  
182 proposed that exposure to *P. shigelloides* via contaminated water protects humans  
183 from *S. sonnei* infection<sup>5</sup> as the O antigens are indistinguishable and cross-react<sup>27,28</sup>.  
184 This may explain increases in *S. sonnei* incidence following economic development  
185 and water quality improvements, as the result of a decline in passive cross-protection  
186 by environmental immunization with *P. shigelloides*. If this cross-protection acts as a  
187 barrier to the establishment of *S. sonnei* in human populations, one would predict that  
188 *S. sonnei* infections would gradually increase following improvements in water  
189 quality, and that the geographical expansion of *S. sonnei* will be characterized by the  
190 introduction and expansion of novel clones moving into human populations with  
191 falling natural immunity previously obtained from exposure to *P. shigelloides*. Our  
192 model of recent dissemination out of Europe is remarkably consistent with these  
193 hypotheses. Transmission of *S. sonnei* into other continents has likely occurred  
194 sporadically over centuries through human migration, trade and travel; however the  
195 establishment of local *S. sonnei* populations – which we would observe as  
196 geographically clustered clonal groups outside Europe – is not evident until the last  
197 few decades.

198

199 Our findings have major implications for global public health and diarrheal infections.  
200 Improvement of drinking water, one of the Millenium Development Goals, is an  
201 undeniably important aim and is expected to reduce morbidity and mortality due to a  
202 diverse array of waterborne diseases. However, we predict that fulfilling this aim will

203 produce a concurrent increase in *S. sonnei* dysentery incidence in transitional  
204 countries. The combination of increased incidence and excessive antimicrobial  
205 resistance among globally disseminated *S. sonnei* indicates an anti-*S. sonnei* vaccine  
206 will be increasingly important for the control and long-term prevention of dysentery  
207 and associated morbidity and mortality. A suitable vaccine is an achievable goal,  
208 since all *S. sonnei* share a single O antigen that has proven to be a successful vaccine  
209 target<sup>29</sup>. Interestingly, the success of *S. sonnei* in the face of diminishing *S. flexneri*  
210 incidence suggests important epidemiological distinctions in transmission of the two  
211 pathogens. *S. sonnei* outbreaks have been associated with schools, care facilities,  
212 contaminated food and insects moving between fecal waste and food preparation  
213 areas<sup>30-32</sup>. These modes of transmission are considerably more direct than waterborne  
214 transmission and may explain the persistence of *S. sonnei* even when water  
215 infrastructure is improved, implying that vaccination and improved hygiene standards  
216 will be pivotal in eliminating *S. sonnei* infections in industrializing countries.

217

#### 218 **URLs**

219 Illumina sequence data provided at <http://www.ebi.ac.uk/ena/data/view/ERP000182>

220 TreeStat: <http://tree.bio.ed.ac.uk/software/treestat/>

221 Velvet Optimiser: <http://www.ebi.ac.uk/~zerbino/velvet/>

222 **Acknowledgements**

223 This work was supported by the Wellcome Trust (#0689); and a Victorian Life  
224 Sciences Computation Initiative (VLSCI) grant (#VR0082) on its Peak Computing  
225 Facility at the University of Melbourne, an initiative of the Victorian Government,  
226 Australia. KEH was supported by a Fellowship from the NHMRC of Australia  
227 (#628930); SB is supported by an OAK Foundation Fellowship through Oxford  
228 University (#OAKF9) and the Li Ka Shing foundation (#LG13); FXW was partially  
229 funded by the Institut de Veille Sanitaire; JY was supported by a MRC grant  
230 (#G0800173); DWK was partially supported by grant RTI05-01-01 from the Ministry  
231 of Knowledge and Economy (MKE). We thank Myron Levine (University of  
232 Maryland School of Medicine, Center for Vaccine Development, USA) and Christoph  
233 Tang (University of Oxford, UK) for their kind gift of *Shigella sonnei* strain 53G.

234

235 **Author Contributions**

236 KEH, NRT, ECH and AK analysed the data and performed phylogenetic analysis.  
237 NRT, GD, JY, SB, JJF, KEH and JP were involved in the study design. FXW, DJB,  
238 JEC, JY, VS, DWK, SYC, SHK, WDS and DJP were involved in isolate collection,  
239 DNA analysis and resistance phenotyping. KEH, SB, NRT, GD, AK, ECH and FXW  
240 contributed to the manuscript writing.

241

242 **Accession Numbers**

243 The finished genome of *S. sonnei* 53G is available under EMBL accessions  
244 HE616528 (chromosome) and HE616529, HE616530, HE616531 and HE616532  
245 (plasmids). Sequence reads for the 132 Illumina-sequenced *S. sonnei* are deposited in  
246 the European Nucleotide Archive under accession ERP000182.

247

248 The authors declare no competing financial interests.

249 **Figure Legends**

250

251 **Figure 1. Bayesian maximum clade credibility phylogeny for *S. sonnei*.** Branches  
252 defining major lineages in bold (each had 100% posterior support); pie charts indicate  
253 ML estimates for geographic origin of major nodes, according to inset legend (lower  
254 left). Time (x-axis) is relative to the Common Era; divergence dates (median estimate  
255 and 95% HPD) are given in blue for major nodes. Distribution of antimicrobial  
256 resistance determinants is indicated in the heatmap according to the legends provided,  
257 which reflect percentage of bases in each gene sequence that are covered by reads  
258 from each isolate (top right). Geographically localised clonal expansions are  
259 highlighted on the right, labeled with their median estimated divergence date.

260

261

262

263 **Online Methods**

264

265 ***Bacterial isolates and sequencing***

266 Bacterial isolates analysed in this study are detailed in Supplementary Table 1. DNA  
267 was prepared using the Wizard Genomic DNA Kit (Promega, Madison, WI) or phenol  
268 extraction. Index-tagged paired end Illumina sequencing libraries were prepared using  
269 one of 12 unique indexing tags as previously described<sup>13</sup>. These were combined into  
270 pools each containing 11-12 uniquely tagged libraries and sequenced on the Illumina  
271 Genome Analyzer GAII according to manufacturer's protocols to generate tagged 54  
272 bp paired-end reads.

273

274 ***Read alignment and SNP detection***

275 Reads from each isolate were mapped to the *S. sonnei* reference genome (strain Ss046  
276 chromosome, NC\_007384; strain Ss046 plasmids, NC\_007385, NC\_009347,  
277 NC\_009346, NC\_009345; plasmid pEG356, NC\_013727) using BWA<sup>33</sup> with default  
278 parameters. Average read depths are given in Supplementary Table 1. SNPs were  
279 identified using SamTools<sup>34</sup>. SNPs in the previously sequenced *S. sonnei* strain 53G  
280 were identified using the same mapping procedure to analyse reads simulated from  
281 the finished genome (chromosome: HE616528; plasmids: HE616529, HE616530,  
282 HE616531 and HE616532) using SamTools' wgsim algorithm. SNPs called in phage  
283 regions or repetitive sequences (10.2% of bases and 15.5% of genes in the Ss046  
284 reference chromosome) were excluded<sup>16</sup>, resulting in a final set of 10,111  
285 chromosomal SNP loci. The allele at each locus in each isolate was determined by  
286 reference to the consensus base in that genome (using SamTools pileup and removing

287 low confidence alleles with consensus base quality  $\leq 20$ , read depth  $\leq 5$  or a  
288 heterozygous base call).

289

290 The SNP calling procedure was repeated using *S. sonnei* 53G (Lineage II) as the  
291 reference for mapping. This resulted in an identical tree topology with near-identical  
292 branch lengths (Pearson correlation coefficient = 0.995,  $p < 1 \times 10^{-15}$ ), demonstrating the  
293 robustness of the method and its independence from the choice of reference genome.  
294 The Ss046-mapped data was used for all analyses reported, since the Ss046 genome  
295 has been widely used in previous comparative studies while the 53G genome is  
296 reported here for the first time.

297

298 The same procedures were followed to identify SNPs in the invasion plasmid. The  
299 analysis was restricted to strains with a mean plasmid read depth of  $\geq 10x$  and the 137  
300 kbp of non-repetitive plasmid sequence (63% of the *S. sonnei* pSs046 reference  
301 plasmid sequence).

302

303 Alleles in outgroup genomes were determined using the same approach to analyse  
304 reads simulated from other *Shigella* and *E. coli* reference genomes (Supplementary  
305 Table 2) using wgsim (distributed with SamTools).

306

### 307 ***Phylogenetic and temporal analyses***

308 Chromosomal SNP alleles were concatenated for each strain to generate a multiple  
309 alignment of all SNPs (where high confidence base calls could not be determined, the  
310 allele was recorded as a gap character). Clusters of SNPs introduced via horizontal  
311 transfer (see *SNP distribution* section below) were removed from the alignment. The



312 resulting alignment was further filtered to remove loci at which alleles were unknown  
313 for >40% of isolates (indicating the site is not conserved) and an ML phylogeny was  
314 estimated using RAxML<sup>35</sup>. The BEAST package<sup>8</sup> was utilized for the Bayesian  
315 inference of phylogeny and divergence dates. Additionally, we used the *BAPS*  
316 program (Bayesian Analysis of Population Structure)<sup>36</sup> to examine clustering of  
317 isolates based on SNP data.

318

319 For ML analysis, RAxML was run ten times using the generalized time-reversible  
320 model with a  $\Gamma$  distribution to model site-specific rate variation (i.e., the GTR+ $\Gamma$   
321 substitution model; GTRGAMMA in RAxML). 1000 bootstrap pseudo-replicate  
322 analyses were performed to assess support for the ML phylogeny. The final result  
323 (Supplementary Fig. 2) is the tree with the highest likelihood across all ten runs, with  
324 ML estimates of branch length and confidence in major bipartitions calculated using  
325 the bootstrap values across all runs. This phylogeny was rooted using *E. coli* and  
326 *Shigella* outgroups (Supplementary Table 2).

327

328 Root-to-tip branches were extracted from the ML tree using the program TreeStat (see  
329 URLs). The relationship between root-to-tip distances, year of isolation and lineage  
330 were analysed using linear regression. Plots and regression lines are shown in  
331 Supplementary Figure 3, along with Pearson correlation coefficients.

332

333 For BEAST analysis, we also used the GTR+ $\Gamma$  substitution model and defined tip  
334 dates as the year of isolation (restricting the analysis to those sequences with recorded  
335 dates). We performed multiple analyses using both constant size and Bayesian skyline  
336 demographic models, in combination with either a strict molecular clock or a relaxed

337 clock (uncorrelated lognormal distribution). BEAST (v1.6) uses a Markov chain  
338 Monte Carlo (MCMC) method for sampling the posterior probability distributions.  
339 Analyses of all model combinations (demographic and clock) were performed using  
340 ten chains of 100 million generations each to ensure convergence, with samples taken  
341 every 1,000 MCMC generations. Parameters were estimated after combining all  
342 replicate analyses, totaling 900 million MCMC generations post-burnin, with all  
343 reported parameter estimates (i.e., medians and 95% Highest Probability Densities –  
344 HPDs) calculated using the program Tracer v1.5. The relaxed clock models provided  
345 much better fit to the data (Bayes Factor > 100; using the harmonic mean estimator of  
346 the marginal likelihood) and the standard deviation of inferred substitution rates  
347 across branches was 0.45 [95% HPD = 0.38 - 0.52], providing additional strong  
348 support for a relaxed molecular clock. Bayesian skyline plots indicated a constant  
349 population size through time and estimates under a constant population model yielded  
350 very similar results to that under a Bayesian skyline model. Therefore, all parameter  
351 estimates quoted are from analyses using relaxed clock and Bayesian skyline  
352 demographic models. To test the validity of the temporal signal in the data, we  
353 performed 20 additional BEAST runs (of 200 million MCMC generations each) with  
354 identical substitution (GTR+ $\Gamma$ ), clock (relaxed), and demographic (Bayesian skyline)  
355 models, but with randomized tip dates (Supplementary Fig. 5). This randomization  
356 procedure produces a null set of tipdate and sequence correlations that may be  
357 analysed to produce null substitution rate distributions, which can then be compared  
358 with empirical rate estimates.

359

360 *Phylogeographic analysis*

361 The geographic region of isolation of each *S. sonnei* was analysed as a discrete  
362 character trait using two complementary methods. Phylogeographic analyses were  
363 performed using the 126 isolates which had complete information on both year and  
364 geographic region of isolation (see Supplementary Table 1). First, the association  
365 between the phylogenetic relationships of *S. sonnei* isolates (inferred by BEAST) and  
366 their geographic region of isolation was tested using the Bayesian Tip-Significance  
367 software (BaTS<sup>14</sup>). A random selection of 50,000 trees sampled during the Bayesian  
368 phylogenetic analysis described above were used as input, and 1,000 randomizations  
369 were used to generate a null distribution for significance testing. Second, ancestral  
370 state reconstruction of the geographic origin of hypothetical common ancestors (i.e.,  
371 internal nodes in the phylogeny) was performed using the ‘ace’ function implemented  
372 in the ‘ape’ package for R<sup>37</sup>. The percent probability estimates quoted, and illustrated  
373 by pie charts in Figure 1, are scaled likelihoods for the discrete character trait (i.e.,  
374 region of isolation) at each node.

375

### 376 ***Gene content analysis***

377 Each read set was assembled using the *de novo* short read assembler Velvet<sup>38</sup> and  
378 Velvet Optimiser (see URLs). Contigs less than 100 bp in size were excluded from  
379 further analysis. The *S. sonnei* 53G genome (chromosome: HE616528; plasmids:  
380 HE616529, HE616530, HE616531 and HE616532) and *de novo* assembled contig  
381 sets were mapped iteratively to the pan-genome reference set (initialized as the  
382 concatenation of *S. sonnei* Ss046 chromosome, NC\_007384; Ss046 plasmids,  
383 NC\_007385, NC\_009347, NC\_009346, NC\_009345; plasmid pEG356, NC\_013727)  
384 using MUMmer (nucmer algorithm)<sup>39</sup>. At each iteration *i*, sequences not aligning to  
385 the current pan-genome  $P_{i-1}$  set were incorporated into an extended pan-genome,  $P_i$ .

386 The final pan-genome, *P*, was annotated using a combination of annotation transfer  
387 (for *S. sonnei* reference sequences) and *de novo* annotation using the RAST  
388 annotation server<sup>40</sup> for novel sequences assembled from reads. The latter included  
389 1.67 Mbp of sequence in 862 contigs, in which 2,422 genes were annotated  
390 (incorporating 80.5% of bases), resulting in a total of 6,852 genes.

391

392 *S. sonnei* read sets were then aligned to the pan-genome using BWA<sup>27</sup> with default  
393 mapping parameters. A pileup was generated for each aligned read set using  
394 SamTools<sup>28</sup> and used to summarize, for each annotated gene in the pan-genome *P*, the  
395 coverage (% of bases covered) and presence of inactivating mutations (nonsense  
396 SNPs or non-triplet indels resulting in frameshifts) in each genome. The results were  
397 used to identify genes whose presence or inactivation was associated with specific  
398 lineages (Supplementary Note, Supplementary Fig. 6).

399

#### 400 ***Resistance gene analysis***

401 The presence of resistance genes was initially determined from mapping data  
402 described above. The genetic context of resistance genes was examined by blastn  
403 search of each contig set with known resistance, transposase or integrase genes as  
404 query sequences. The resulting contigs were compared to the NCBI non-redundant  
405 nucleotide database to annotate the resistance genes and mobile elements. Mapping  
406 was then repeated using annotated mobile elements to generate the gene coverage  
407 maps shown in Figure 1 and Supplementary Figure 2, which indicate the proportion  
408 of bases in each gene sequence that are covered by reads from each isolate (reference  
409 sequences are provided in Supplementary Fig. 2).

410 **References**

- 411 1. Pupo, G.M., Lan, R. & Reeves, P.R. Multiple independent origins of Shigella  
412 clones of Escherichia coli and convergent evolution of many of their characteristics.  
413 *Proc Natl Acad Sci U S A* **97**, 10567-72 (2000).
- 414 2. Yang, F. *et al.* Genome dynamics and diversity of Shigella species, the  
415 etiologic agents of bacillary dysentery. *Nucleic Acids Res* **33**, 6445-58 (2005).
- 416 3. DuPont, H.L., Levine, M.M., Hornick, R.B. & Formal, S.B. Inoculum size in  
417 shigellosis and implications for expected mode of transmission. *J Infect Dis* **159**,  
418 1126-8 (1989).
- 419 4. Kotloff, K.L. *et al.* Global burden of Shigella infections: implications for  
420 vaccine development and implementation of control strategies. *Bull World Health*  
421 *Organ* **77**, 651-66 (1999).
- 422 5. Sack, D.A., Hoque, A.T., Huq, A. & Etheridge, M. Is protection against  
423 shigellosis induced by natural infection with *Plesiomonas shigelloides*? *Lancet* **343**,  
424 1413-5 (1994).
- 425 6. Vinh, H. *et al.* A changing picture of shigellosis in southern Vietnam: shifting  
426 species dominance, antimicrobial susceptibility and clinical presentation. *BMC Infect*  
427 *Dis* **9**, 204 (2009).
- 428 7. Karaolis, D.K., Lan, R. & Reeves, P.R. Sequence variation in Shigella sonnei  
429 (Sonnei), a pathogenic clone of Escherichia coli, over four continents and 41 years. *J*  
430 *Clin Microbiol* **32**, 796-802 (1994).
- 431 8. Drummond, A.J. & Rambaut, A. BEAST: Bayesian evolutionary analysis by  
432 sampling trees. *BMC Evol Biol* **7**, 214 (2007).

- 433 9. Nastasi, A., Pignato, S., Mammina, C. & Giammanco, G. rRNA gene  
434 restriction patterns and biotypes of *Shigella sonnei*. *Epidemiol Infect* **110**, 23-30  
435 (1993).
- 436 10. Touchon, M. *et al.* CRISPR distribution within the *Escherichia coli* species is  
437 not suggestive of immunity-associated diversifying selection. *J Bacteriol* **193**, 2460-7  
438 (2011).
- 439 11. Mutreja, A. *et al.* Evidence for several waves of global transmission in the  
440 seventh cholera pandemic. *Nature* **477**, 462-5 (2011).
- 441 12. Morelli, G. *et al.* *Yersinia pestis* genome sequencing identifies patterns of  
442 global phylogenetic diversity. *Nat Genet* **42**, 1140-3 (2010).
- 443 13. Harris, S.R. *et al.* Evolution of MRSA during hospital transmission and  
444 intercontinental spread. *Science* **327**, 469-74 (2010).
- 445 14. Parker, J., Rambaut, A. & Pybus, O.G. Correlating viral phenotypes with  
446 phylogeny: accounting for phylogenetic uncertainty. *Infect Genet Evol* **8**, 239-46  
447 (2008).
- 448 15. Ranjbar, R. *et al.* Genetic relatedness among isolates of *Shigella sonnei*  
449 carrying class 2 integrons in Tehran, Iran, 2002-2003. *BMC Infect Dis* **7**, 62 (2007).
- 450 16. Holt, K.E. *et al.* High-throughput sequencing provides insights into genome  
451 variation and evolution in *Salmonella Typhi*. *Nat Genet* **40**, 987-93 (2008).
- 452 17. World Health Organization Guidelines for the control of shigellosis, including  
453 epidemics due to *Shigella dysenteriae* 1. WHO Document Production Services,  
454 Geneva, Switzerland (2005).
- 455 18. Christopher, P.R., David, K.V., John, S.M. & Sankarapandian, V. Antibiotic  
456 therapy for *Shigella* dysentery. *Cochrane Database of Systematic Reviews*, CD006784  
457 (2010).

- 458 19. Vinh, H. *et al.* A multi-center randomized trial to assess the efficacy of  
459 gatifloxacin versus ciprofloxacin for the treatment of shigellosis in Vietnamese  
460 children. *PLoS Negl Trop Dis* **5**, e1264 (2011).
- 461 20. Vinh, H. *et al.* Treatment of bacillary dysentery in Vietnamese children: two  
462 doses of ofloxacin versus 5-days nalidixic acid. *Trans Royal Soc Trop Med Hyg* **94**,  
463 323-6 (2000).
- 464 21. Jeong, H.J. *et al.* Acanthamoeba: could it be an environmental host of  
465 Shigella? *Exp Parasitol* **115**, 181-6 (2007).
- 466 22. Saeed, A., Abd, H., Edvinsson, B. & Sandstrom, G. Acanthamoeba castellanii  
467 an environmental host for Shigella dysenteriae and Shigella sonnei. *Arch Microbiol*  
468 **191**, 83-8 (2009).
- 469 23. Winiecka-Krusnell, J. & Linder, E. Free-living amoebae protecting Legionella  
470 in water: the tip of an iceberg? *Scand J Infect Dis* **31**, 383-5 (1999).
- 471 24. Greub, G. & Raoult, D. Microorganisms resistant to free-living amoebae. *Clin*  
472 *Microbiol Rev* **17**, 413-33 (2004).
- 473 25. Shepherd, J.G., Wang, L. & Reeves, P.R. Comparison of O-antigen gene  
474 clusters of Escherichia coli (Shigella) sonnei and Plesiomonas shigelloides O17:  
475 sonnei gained its current plasmid-borne O-antigen genes from P. shigelloides in a  
476 recent event. *Infect Immunity* **68**, 6056-61 (2000).
- 477 26. Sansonetti, P.J., Kopecko, D.J. & Formal, S.B. Shigella sonnei plasmids:  
478 evidence that a large plasmid is necessary for virulence. *Infect Immunity* **34**, 75-83  
479 (1981).
- 480 27. Van de Verg, L.L., Herrington, D.A., Boslego, J., Lindberg, A.A. & Levine,  
481 M.M. Age-specific prevalence of serum antibodies to the invasion plasmid and

482 lipopolysaccharide antigens of *Shigella* species in Chilean and North American  
483 populations. *J Infect Dis* **166**, 158-61 (1992).

484 28. Shimada, T. & Sakazaki, R. On the serology of *Plesiomonas shigelloides*. *Jap*  
485 *J Med Science Biol* **31**, 135-42 (1978).

486 29. Kaminski, R.W. & Oaks, E.V. Inactivated and subunit vaccines to prevent  
487 shigellosis. *Exp Review Vaccines* **8**, 1693-704 (2009).

488 30. Genobile, D. *et al.* An outbreak of shigellosis in a child care centre.  
489 *Communicable Dis Intell* **28**, 225-9 (2004).

490 31. Lewis, H.C. *et al.* Outbreaks of *Shigella sonnei* infections in Denmark and  
491 Australia linked to consumption of imported raw baby corn. *Epidemiol Infect* **137**,  
492 326-34 (2009).

493 32. Cohen, D. *et al.* Reduction of transmission of shigellosis by control of  
494 houseflies (*Musca domestica*). *Lancet* **337**, 993-7 (1991).

495 33. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-  
496 Wheeler transform. *Bioinformatics* **25**, 1754-60 (2009).

497 34. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools.  
498 *Bioinformatics* **25**, 2078-9 (2009).

499 35. Stamatakis, A. RAxML-VI-HPC: maximum likelihood-based phylogenetic  
500 analyses with thousands of taxa and mixed models. *Bioinformatics* **22**, 2688-90  
501 (2006).

502 36. Tang, J., Hanage, W.P., Fraser, C. & Corander, J. Identifying currents in the  
503 gene pool for bacterial populations using an integrative approach. *PLoS Comp Biol* **5**,  
504 e1000455 (2009).

505 37. Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and  
506 Evolution in R language. *Bioinformatics* **20**, 289-90 (2004).



- 507 38. Zerbino, D.R. & Birney, E. Velvet: algorithms for de novo short read  
508 assembly using de Bruijn graphs. *Genome Res* **18**, 821-9 (2008).
- 509 39. Kurtz, S. *et al.* Versatile and open software for comparing large genomes.  
510 *Genome Biol* **5**, R12 (2004).
- 511 40. Aziz, R.K. *et al.* The RAST Server: rapid annotations using subsystems  
512 technology. *BMC Genomics* **9**, 75 (2008).
- 513 41. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular  
514 Biology Open Software Suite. *Trends Genet* **16**, 276-7 (2000).
- 515 42. Croucher, N.J. *et al.* Rapid pneumococcal evolution in response to clinical  
516 interventions. *Science* **331**, 430-4 (2011).

517  
518

519

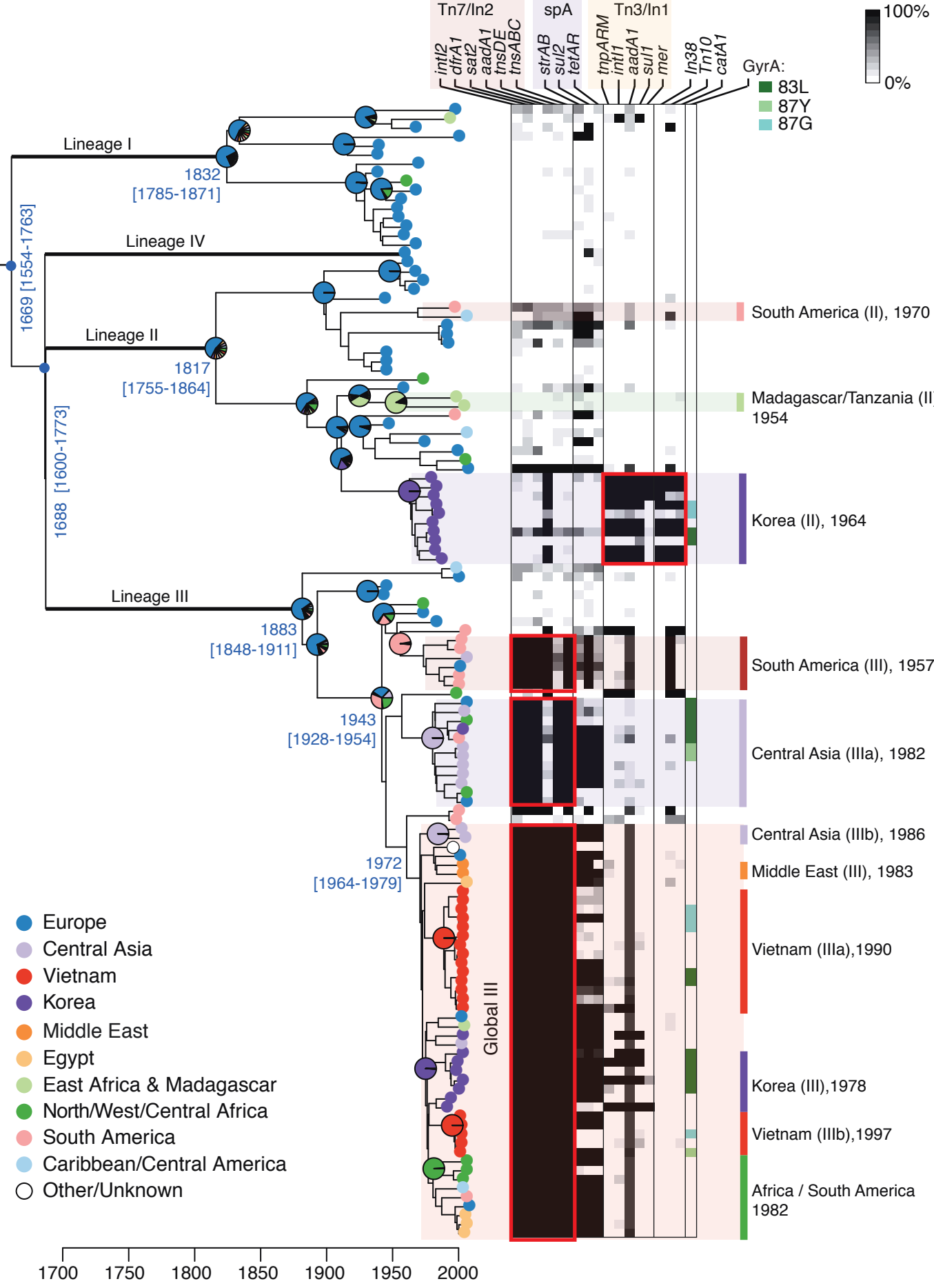
520 **Editorial Summary (AOP and Month, same):**

521 Nicholas Thomson and colleagues report whole-genome sequencing of 132 globally  
522 distributed isolates of *Shigella sonnei*, a cause of human dysentery. Their  
523 phylogeographic analyses suggest that the current *S. sonnei* population is under 500  
524 years old, and the authors are able to trace several distinct lineages that have spread  
525 out of Europe to other continents over the last few decades.

Type of file: figure

Label: 1

Filename: figure\_1.ai



UKPMC+ has received the file 'supp\_info\_1.pdf' as supplementary data. The file will not appear in this PDF Receipt, but it will be linked to the web version of your manuscript.

UKPMC+ has received the file 'supp table 1.xls' as supplementary data. The file will not appear in this PDF Receipt, but it will be linked to the web version of your manuscript.