# Origin of modern syphilis and emergence of a pandemic *Treponema pallidum* cluster

Natasha Arora[1,2*], Verena J. Schuenemann[3], Günter Jäger[4], Alexander Peltzer[3,4†], Alexander Seitz[4], Alexander Herbig[3,4†], Michal Strouhal[5], Linda Grillová[5], Leonor Sánchez-Busó[6,7], Denise Kühnert[8], Kirsten I. Bos[3†], Leyla Rivero Davis[1‡], Lenka Mikalová[5], Sylvia Bruisten[9], Peter Komericki[10], Patrick French[11], Paul R. Grant[12], María A. Pando[13], Lucía Gallo Vaulet[14], Marcelo Rodríguez Fermepin[14], Antonio Martinez[15], Arturo Centurion Lara[16], Lorenzo Giacani[16], Steven J. Norris[17] David Šmajs[5], Philipp P. Bosshard[18], Fernando González-Candelas[6*], Kay Nieselt[4*], Johannes Krause[3†*] and Homayoun C. Bagheri[18§*]

**Affiliations:**

[1]Institute for Evolutionary Biology and Environmental Studies, University of Zurich, Zurich, Switzerland.

[2]Zurich Institute of Forensic Medicine, University of Zurich, Zurich, Switzerland.

[3]Institute for Archaeological Sciences, University of Tübingen, Tübingen, Germany.

[4]Center for Bioinformatics, University of Tübingen, Tübingen, Germany.

[5]Department of Biology, Faculty of Medicine, Masaryk University, Brno, Czech Republic.

[6]Unidad Mixta Infección y Salud Pública FISABIO/Universidad de Valencia. CIBER in Epidemiology and Public Health, Spain.

[7]Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, United Kingdom

[8]Institute of Integrative Biology, Department of Environmental Systems Science, ETH Zürich, Switzerland.

[9]Public Health Laboratory, GGD Amsterdam, Department of Infectious Diseases, Amsterdam, the Netherlands.

[10]Department of Dermatology, Medical University of Graz, Graz, Austria.

[11]The Mortimer Market Centre CNWL, Camden Provider Services, UK.

[12]Department of Clinical Microbiology and Virology, University College London Hospitals NHS Foundation Trust, London, UK.

[13]Instituto de Investigaciones Biomédicas en Retrovirus y SIDA (INBIRS), Universidad de Buenos Aires-CONICET, Buenos Aires, Argentina.

[14]Universidad de Buenos Aires, Facultad de Farmacia y Bioquímica, Departamento de Bioquímica Clínica, Microbiología Clínica, Buenos Aires, Argentina.

[15]Servicio de Dermatología. Hospital General Universitario de Valencia, Spain.

34    [16]University of Washington, Department of Medicine, Division of Allergy and Infectious Diseases, and
35    Department of Global Health, Seattle (WA), USA.

36    [17]Department of Pathology and Laboratory Medicine, UTHealth McGovern Medical School, Houston,
37    TX USA.

38    [18]Department of Dermatology, University Hospital of Zurich, Zurich, Switzerland.

39    [†]Current address: Department of Archaeogenetics, Max Planck Institute for the Science of Human
40    History, Jena, Germany.

41    [‡]Current address: Department of Infectious Disease Epidemiology, Imperial College London, UK.

42    [§]Current address: Repsol Technology Center, Madrid, Spain

43    [*]Corresponding authors

44    **Introductory paragraph:** The abrupt onslaught of the syphilis pandemic starting in the late
45    15[th] century established this devastating infectious disease as one of the most feared in
46    human history [1]. Surprisingly, despite the availability of effective antibiotic treatment since
47    the mid-20[th] century, this bacterial infection caused by *Treponema pallidum* subsp. *pallidum*
48    (TPA), has been re-emerging globally in the last few decades with an estimated 10.6 million
49    cases in 2008 [2]. While resistance to penicillin has not yet been identified, an increasing
50    number of strains fail to respond to the second-line antibiotic azithromycin [3]. Little is known
51    about the genetic patterns in current infections or the evolutionary origins of the disease
52    due to the low quantities of treponemal DNA in clinical samples, and difficulties to cultivate
53    the pathogen [4]. Here we used DNA capture and whole genome sequencing to successfully
54    interrogate genome-wide variation from syphilis patient specimens, combining it with
55    laboratory samples of TPA and two other subspecies. Phylogenetic comparisons based on
56    the sequenced genomes indicate that the TPA strains examined share a common ancestor
57    after the 15[th] century, within the early modern era. Moreover, most contemporary strains
58    are azithromycin resistant and members of a globally dominant cluster named here as SS14-
59    Ω. This cluster diversified from a common ancestor in the mid-20[th] century subsequent to
60    the discovery of antibiotics. Its recent phylogenetic expansion and global presence point to
61    the emergence of a pandemic strain cluster.

62    **Main Text:** The first reported syphilis outbreaks in Europe occurred during the War of Naples
63    in 1495 [5], prompting unresolved theories on a post-Columbian introduction [6,7].
64    Subsequently, the epidemic spread to other continents, remaining a severe health burden
65    until treatment with penicillin five centuries later enabled incidence reduction. The striking
66    present-day resurgence is poorly understood, particularly the underlying patterns of genetic
67    diversity. Much of our molecular understanding of treponemes comes from propagating
68    strains in laboratory animals to obtain sufficient DNA. The few published whole genomes
69    were obtained after amplification through rabbit passage [4,8–10], and represent limited
70    diversity for phylogenetic analyses. These sequences suggest that the TPA genome of 1.14
71    Mb is genetically monomorphic. Potential genetic diversity remains unexplored because
72    clinical samples are mostly typed by PCR amplification of only 1-5 loci [11,12]. These
73    epidemiological strain typing studies are motivated by the limitations of serologic or
74    microscopic tests to distinguish among TPA strains or among the subspecies *Treponema*
75    *pallidum* subsp. *pertenue* (TPE) and *Treponema pallidum* subsp. *endemicum* (TEN), which
76    cause the diseases yaws and bejel, respectively. While all three diseases are transmitted
77    through skin contact and show an overlap in their clinical manifestations, syphilis is
78    geographically more widespread and generally transmitted sexually. The precise
79    relationships among the bacteria are still debated, particularly regarding the evolutionary
80    origin of syphilis.

81    The paucity of molecular studies and the focus on typing of a few genes means that we have
82    limited information regarding the evolution and spread of epidemic TPA. In this study, we
83    interrogated genome-wide variation across geographically widespread isolates. In total, we
84    obtained 70 samples from 13 countries, including 52 syphilis swabs collected directly from
85    patients between 2012 and 2013, and 18 syphilis, yaws, and bejel samples collected from
86    1912 onwards and propagated in laboratory rabbits (Supplementary Table 1). Through
87    comparative genome analyses and phylogenetic reconstruction, we shed light on the
88    evolutionary history of TPA and identify epidemiologically relevant haplotypes.

89    Due to the large background of host DNA, samples were enriched for treponemal DNA prior
90    to Illumina sequencing [13,14]. The resultant reads were mapped to the Nichols TPA reference

1

genome (RefSeq NC_021490; Supplementary Table 3) [4,15]. Genomic coverage ranged from 0.13-fold to over 1000-fold. As expected, the highest mean coverage was found in strains propagated in rabbits, while high variation in mean coverage was observed in samples collected directly from patients (0.13-fold to 223-fold) (Supplementary Table 2). This heterogeneity could potentially affect our inferences. Therefore, we restricted the genome-wide analyses to the 28 samples where at least 80% of the genome was covered by a minimum of three reads (highlighted in Supplementary Table 2). Across the 28 samples, the average proportion of genome coverage with at least 3-fold or 10-fold depth was 97% and 82%, respectively (Supplementary Table 4).

*De novo* assemblies for the four highest covered syphilis swab samples (NE17, NE20, CZ27, AU15) and one Indonesian yaws isolate (IND1) show no significant structural changes in the five genomes (Fig. 1a; Supplementary Table 5), except for the deletion in IND1 of gene TP1030, which potentially encodes a virulence-factor [17]. The deletion was shared across all the yaws infection isolates (Supplementary Methods), consistent with other studies [18].

Prior to phylogenetic reconstruction we checked for signatures of recombination. While *T. pallidum* is considered to be a clonal species [19], previous studies suggest recombinant genes in a Mexican syphilis and a Bosnian bejel strain [10,16]. We screened for putative recombinants across the 978 annotated genes in our 28 sequenced genomes and the 11 publicly available genomes from laboratory strains (Supplementary Table 3). Genes were selected as candidates if they had unexpectedly high SNP densities, incongruent topologies with the genome-wide tree and more than 4 homoplasies in a pair of branches (Supplementary Methods). We identified 4 genes coding for outer membrane proteins (Supplementary Table 6), one of which (TP0136) is used in typing studies [8].

After excluding the 4 putative recombinant genes, the genome alignment for all 39 genomes contained 2,235 variable positions. We used the Bayesian framework implemented in BEAST [20] to reconstruct a phylogenetic tree (Fig. 1b). The tree topology revealed a marked separation between TPA and TPE/TEN (100% Bayesian posterior support), with TPA forming a monophyletic lineage. The distinction of the two lineages was robust even with the inclusion of putative recombinant genes (Supplementary Fig. 2). Analyses of divergence between the two lineages yielded an average mean distance of 1225 nucleotide differences. By contrast, within each of the lineages we found considerably less diversity (124.6 average pairwise mutations within the TPA lineage and 200.2 within TPE/TEN). A heat map (Supplementary Fig. 3) to show shared variation for pairs of samples with respect to the Nichols reference genome, confirms the divergence between the lineages. The underlying SNP matrix yielded 443 SNPs specific to TPA genomes and 1703 to TPE/TEN genomes. Previous studies have found cross-subspecies groupings when relying on a limited set of markers [21]. Our results, incorporating genome-wide data from clinical samples, not only establish a clear separation between the two lineages, in agreement with studies examining genomic data from rabbit propagated samples [10,18], but also illustrate the need for a careful choice of taxonomic markers when genome-wide data is not available.

Using the sample isolation dates as tip calibration and applying the Birth Death Serial Skyline model [23], we obtained a mean evolutionary rate of $3.6 \times 10^{-4}$ (rate variance $3.8 \times 10^{-8}$; 95% HPD $1.86 \times 10^{-4}$ - $5.73 \times 10^{-4}$). This estimate is equivalent to scaled mean rate of $6.6 \times 10^{-7}$ substitutions per site per year for the whole genome, in line with estimates for other clonal human pathogens such as *Shigella sonnei* ($6.0 \times 10^{-7}$) and *Vibrio cholerae* (O1 lineage; $8.0 \times 10^{-7}$) [24,25]. Our divergence analyses for TPA samples provide a time to the most recent common

137  ancestor (TMRCA) less than 500 years ago (mean calendar year 1744, 95% HPD 1611-1859;
138  Fig. 1B).

139  Within the TPA lineage the samples group in two clades named after the SS14 and Nichols
140  reference genomes (with 100% and 82% posterior support values respectively). The Nichols
141  clade consists almost exclusively of samples collected from patients in North America from
142  1912 to 1986 and passaged in rabbits prior to sequencing, with the exception of one patient
143  sample from 2013 (NE20). In contrast, the SS14 clade has a geographically widespread
144  distribution, encompassing European, North American and South American samples
145  collected from infections between 1951 and 2013. We investigated the TPA clades further by
146  generating a median-joining (MJ) network to illustrate the mutational differences among the
147  TPA samples (Fig. 2a). As underscored by distances in the network, greater nucleotide
148  diversity is found within the Nichols clade ($\pi$=0.05) compared to the SS14 clade ($\pi$= 0.01).
149  Three closely related sequences derive from the original Nichols sample isolated from the
150  cerebrospinal fluid of a patient in 1912 and propagated in the lab: NIC_REF, the reference
151  genome re-sequenced by Pětrošová et al. [15], and NIC-1 and NIC-2, which we sequenced
152  following independent propagation of the strains in Houston and Seattle, respectively,
153  during different time periods (Supplementary Table 1 and Supplementary Table 3). These
154  three group together with another three sequences in a cluster labelled Nichols- α (Fig. 2a),
155  with a TMRCA at the turn of the 19[th] century (Fig. 1a). The less diversified SS14 clade
156  contains a dominant central haplotype (labelled as SS14-Ω) from which the other sequences
157  radiate (Fig. 2A). Critically, the cluster associated with the SS14-Ω haplotype contains all but
158  one of the recent patient samples from 2012-2013 (n=17) that were captured and
159  sequenced directly, in addition to samples from 1977 (n=1) and 2004 (n=2). The genetic
160  variation within the SS14-Ω cluster is found primarily as singleton mutations (95.5%), with no
161  evidence for geographical structuring. Bayesian analyses estimate a median coalescence for
162  the SS14-Ω cluster in 1963 (95% HPD 1948-1974; Fig. 1a), at a time when incidence was
163  reduced due to the introduction of antibiotics. The star-like topology of this cluster observed
164  in both the tree and the network is suggestive of a recent and rapid clonal expansion.

165  To determine whether the dominance of SS14 clade sequences applies across other
166  countries for which genetic data is available, we examined sequences from the widely typed
167  TP0548 gene in worldwide epidemiological studies [11]. Phylogenies for the TP0548 typing
168  regions separate the SS14 from the Nichols clade for the TPA samples, while not
169  distinguishing the TPA and TPE/TEN lineages (Supplementary Methods; Supplementary Fig.
170  3). Across 1353 worldwide TP0548 sequences from clinical samples, including the 78 from
171  patients in this study, we found that 94% of them grouped in the SS14 clade (Supplementary
172  Tables 8-9; Supplementary Fig. 5), consistent with a probable recent spread of the epidemic
173  cluster. The wide geographical distribution of the SS14 clade establishes it as representative
174  of the present worldwide epidemic. While studies to date have focused on the Nichols strain
175  [26,27], our results indicate that further work on the SS14 clade is warranted.

176  Critically, typing of samples over multiple years in the Czech Republic, San Francisco, British
177  Columbia and Seattle indicate that macrolide antibiotic resistance has increased over time
178  [3,12,28–30]. We queried the presence of the two mutations (A2058G and A2059G) in the 23S
179  rRNA genes associated with azithromycin resistance [3,31,32]. As observed in the MJ network,
180  the resistance marker is a dominant characteristic of the SS14-Ω cluster (Fig. 2a), although it
181  is also found in a recent patient sample (NE20) of the Nichols clade. Extending our analyses
182  of the 23S rRNA gene to all sequenced samples from our study, including the 42 with lower
183  coverage, revealed the mutations in 90% of the SS14 (n=51) and 25% of the Nichols (n=12)

3

184  samples, indicating that neither resistance nor sensitivity is clade-specific (Supplementary
185  Table 8). Hence resistance was probably not an ancestral characteristic of the SS14 clade. A
186  likely scenario is that the extensive usage of azithromycin to treat syphilis and a wide range
187  of bacterial infections, including co-infections with other sexually-transmitted diseases
188  (STDs) such as chlamydia, has played an important role in the selection and subsequent
189  spread of resistance [33,34].

190  Results here represent the first reported set of whole genome sequences successfully
191  obtained directly from syphilis patients, enabling us to disentangle evolutionary relationships
192  at high resolution, and paving the way for further clinical sequencing from current
193  epidemics. Given our identification of putative recombinant genes in *Treponema*, and
194  previous reports on genes involved in homologous recombination [4,35], further detailed
195  analyses on the potential mechanisms of recombination will be necessary. Our phylogenetic
196  reconstruction indicates that all TPA samples examined to date share a common ancestor
197  that was infecting populations in the 1700s, within the early centuries of the modern era,
198  and that was successful in leaving descendants until today. This date is posterior to the
199  colonization of the Americas, and therefore potentially compatible with the post-Columbian
200  model for the emergence of syphilis in Europe. Nonetheless, our work does not exclude the
201  possibility that older TPA lineages had previously existed in Europe but went extinct.
202  Obtaining more patient sample genomes with high coverage could potentially refine our
203  detection of putative recombinants and our phylogenetic inferences. In addition, sequencing
204  from ancient skeletal material would help to further ascertain the history of syphilis.
205  Interestingly, we observed a time difference between the first reported syphilis outbreak in
206  1495 and the last common ancestor of modern strains dated to the 1700s. While this
207  difference could stem from imprecision in the divergence estimates, an alternative scenario
208  is the eventual establishment of a specific lineage due to selection. For instance, it has been
209  hypothesized that the symptoms of syphilis became less severe after the first reported
210  outbreaks in Europe because of the evolution of strains with lower virulence and higher
211  transmission rates [36]. In this scenario, the 18th century provided the context for the origin
212  and propagation of a lineage that successfully outcompeted other lineages.

213  Critical to our epidemiological understanding of contemporary syphilis is our observation of
214  an epidemic cluster (SS14-Ω) that emerged after the discovery of antibiotics. The relatively
215  recent phylogenetic expansion of the SS14-Ω cluster and its global presence point to the
216  emergence of a pandemic azithromycin-resistant cluster. The genome-wide data in this
217  study will be useful to determine a suitable set of typing loci, since typing remains a more
218  accessible method for most laboratories. Further characterization of the genomic diversity of
219  TPA across the globe can prove instrumental in understanding the genetic and
220  epidemiological basis for the spread of SS14-Ω strains.

**Methods**

**Sample collection, DNA extraction and library preparation**

Samples from 64 syphilis infections, 5 yaws infections and 1 bejel infection were collected from numerous countries across the globe (Supplementary Table 1). Syphilis infection samples were classified as either clinical, if obtained from patients directly, or as laboratory strains, if passaged in rabbits after isolation from patients. Clinical samples were obtained after swabbing lesions from patients at sexual health clinics, dermatological clinics or hospitals. Flocked swabs (from Copan Diagnostics, Brescia, Italy) or Nylon swabs were used according to local laboratory instructions. Laboratory strains were obtained as DNA extracts from Masaryk University (Brno, Czech Republic) and the University of Washington (Seattle, USA).DNA extractions were carried out in the participating laboratories using in-house protocols. At the University of Zurich the QIAmp DNA mini kit and QIAmp DNA blood min kit (Qiagen) were used following the manufacturer's protocols.

Library preparation was conducted following a modified Illumina protocol for ancient DNA [14,37], at the University of Tübingen (Supplementary Materials and Methods). Libraries were barcoded with double indices.

**Genome-wide enrichment and sequencing**

Target enrichment for *Treponema pallidum* subsp. *pallidum* was carried out through two rounds of capture hybridization on a 1 million Agilent SureSelect array following the protocol detailed by Hodges et al. [13]. The probes on the array were based on two reference genomes (Nichols, here abbreviated as NIC_REF, GenBank ID CP004010.2/RefSeq ID NC_021490.2, and SS14, GenBank ID CP000805.1/RefSeq ID NC_010741.1). High-throughput sequencing of the enriched libraries was performed on an Illumina Hiseq 2500 platform.

**Sequencing analyses and genome reconstruction**

We applied EAGER [38], our own developed pipeline for read preprocessing (adapter clipping, merging of corresponding paired-end reads in the overlapping regions and quality trimming), mapping, variant identification and genome reconstruction, to all sequenced samples (for full details see Supplementary Materials and Methods). All reads (merged and unmerged) were treated as single-end reads and mapping was performed using the BWA-MEM algorithm [39] with default parameters, using the Nichols genome as a reference. Subsequently, we selected the samples which had at least 80% coverage of the Nichols genome and a minimum of 3 reads (n= 28 samples, Supplementary Table 3). For each of these samples, we used the Genome Analysis Toolkit (GATK) [40] to generate a mapping assembly, applying the UnifiedGenotyper module of GATK to call reference bases and variants from the mapping. The reference base was called if the genotype quality of the call was at least 30 and the position was covered by at least 3 reads. A variant position (SNP) was called if the following criteria 3 were met: i) the position was covered by at least 3 reads; ii) the genotype quality of the call was at least 30 and iii) the minimum SNP allele frequency was 90%. If neither of the requirements for a reference base call nor the requirements for a variant call were met, the character 'N' was inserted at the respective position. For the generation of draft genome sequences we used an in-house tool (VCF2Genome), which reads a VCF file such as produced by the GATK UnifiedGenotyper and incorporates for each row, and thus for each call, one nucleotide into the new draft sequence.

In order to apply our analysis pipeline also to those samples for which complete genomic sequences are available in GenBank (Supplementary Table 2), we produced artificial

266 reads in these cases using an in-house tool (Genome2Reads), and then applied the same
267 mapping, SNP calling and genome reconstruction procedure as for the sequenced samples in
268 order to obtain consistent and comparable results.

269 To investigate conservation of structure and gene order in the genomes, in addition to
270 the mapping assembly, we also performed a *de novo* assembly for the 5 samples with
271 highest coverage (Supplementary Table 5). Our *de novo* assembly pipeline started with the
272 merged reads and in a first step utilized the short read assembler software SOAPdenovo2
273 using ten different k-mer sizes (k = 37 + i·10, i=0,…,9). Different k-mer sizes were used
274 because merging of read pairs into one single read results in very different lengths (between
275 30 and 190 bases). Next, all input reads were mapped back against the resulting contigs
276 using BWA-MEM [39]. Contigs that were not supported by any reads (no read mapped against
277 these contigs) were removed. In order to assemble the contigs resulting from the different k-
278 mers, the remaining contigs were subject to the overlap-based String Graph Assembler (SGA)
279 [41]. Finally, contigs smaller than 1,000 bp were removed before these contigs were mapped
280 against the Nichols reference genome for comparison of genome architectures.

281 Analyses to detect recombinants and reconstruct evolutionary relationships using
282 genome-wide variation were conducted for the 28 sequenced samples meeting our genome-
283 wide coverage criteria (highlighted in the Supplementary Table 3) as well as the 11 published
284 genomes (Supplementary Table 2). Across the 39 whole genomes and draft genomes, 31
285 were TPA, 8 TPE and 1 TEN.

**Recombination detection**

287 Tests for the non-vertical transmission of genes were carried out on the TPA, TPE and
288 TEN genomes (n= 39) by identifying those genes that i) had an unexpectedly high number of
289 SNPs and ii) displayed patterns of transmission (i.e., phylogenies) incongruent with most
290 other genes. First, an expected substitution rate was computed by dividing the total number
291 of observed SNPs in the 978 annotated genes (n=2,098) by the total length of these genes
292 (1,046,421 bp). This rate was then used to calculate the expected number of polymorphisms
293 per gene according to its length. A total of 87 genes displayed at least twice the expected
294 number of polymorphisms. Second, for each of these 87 genes the gene sequence alignment
295 and the gene tree topology were tested against the maximum likelihood tree topology of the
296 draft genome in TREE-PUZZLE v5.2 [42,43]. Genes for which both the Expected Likelihood
297 Weight [44] and the Shimodaira-Hasegawa [45] test rejected the genome tree (p < 0.05) were
298 examined more closely. Third, genes within which we identified a minimum of 4 homoplasies
299 (identical mutations in separate lineages) in at least 2 branches of the tree were marked as
300 putative recombinants (Supplementary Table 6).

**Genome-wide variation and phylogenetic analyses**

302 We investigated genome-wide patterns of polymorphism and divergence using MEGA 6 [46]
303 and DnaSP v.5.10 to compute various measures of diversity including the average pairwise
304 nucleotide differences, Nei's Pi (π), and the number of singletons in each group. We also
305 estimated the number of SNPs private to particular groups. A comparison of the TPA and
306 TPE/TEN genomes revealed between 1 (NIC1) and 339 (AR2) SNPs observed in the TPA
307 samples and between 1091 (GHA1) and 1443 (Bosnia A) SNPs in the TPE/TEN strains
308 (Supplementary Table 4). Furthermore, we produced a heat map to display the number of
309 SNPs that any two genomes share (Supplementary Fig. 3).

310 The molecular clock hypothesis was tested with the maximum likelihood analysis in
311 MEGA 6.0 [46]. Tests were conducted for all TPA, TPE and TEN genomes (39 samples) using i)

multiple whole genome alignments and ii) alignments with only the variable positions, in both cases excluding the 4 putative recombinant genes. The molecular clock hypothesis was rejected at the 5% significance level.

Bayesian phylogenetic trees were produced in BEAST 2.3 [47] for the 28 sequenced samples and the 11 published samples. We compared the trees generated with the alignment of all variable positions in the TPA, TPE and TEN genomes (2,506) and the tree generated with the set of variable positions after excluding the 4 putative recombinant genes (2,235 positions). Additionally, rooted trees were generated with Maximum Parsimony by including *Treponema paraluiscuniculi* (NC_015714) as the outgroup.

As a calibration for the BEAST trees we used tip dates, that is, the isolation years of all samples. When not known with precision, we provided a range (for NIC_REF, NIC1, NIC2, and GAU). The two demographic models (coalescent tree prior under Constant Size and the Birth-Death Serial Skyline model (BDSS)) resulted in consistent parameter estimates. The relaxed clock model was chosen over the strict clock model based on marginal likelihood estimates obtained with PathSampler [47,48]. We provide results for the BDSS model run with the following specifications: uncorrelated lognormal relaxed clock-clock model, GTR plus gamma substitution model, 50 million generations with parameter sampling every 5,000 generations. The log file was viewed in Tracer 1.6 [49] to determine the appropriate burn-in period for adequate effective sample sizes. The annotated maximum clade credibility tree was visualized and edited using Figtree v1.4.2 [50]. Because TPA samples are the focus of this study and therefore more extensively sampled, we report mean branch rate and divergence estimates for the TPA lineage. The mean branch rate estimate obtained is in line with the number of mutations that differed between the samples NIC_REF and NIC 2 (n=15), which were isolated 15-20 years apart following continuous rabbit propagation. We also checked that a run with the same specifications but with only TPA samples (n=31) produced consistent results.

The phylogenetic relationships among the closely related TPA samples (n=31) were examined and visualized through a median joining (MJ) network analysis in Network 4.6 and Network Publisher [51,52] using all variable positions after excluding the putative recombinant loci and sites with missing data (resulting in a total of 628 variable positions).

**Clade classification**

*Samples from this study:* From the 70 TPA, TPE and TEN samples sequenced in this study, 28 fulfilled our criteria for genome-wide analyses (minimum 80% genome covered with at least 3 reads). For the remaining 42 samples, we implemented two classification strategies. First, we generated a new clade prediction strategy based on NGS reads to classify the genomes according to lineage (TPA or TPE/TEN), and within the TPA lineage, as part of the SS14 or the Nichols clade (details provided in the Supplementary Information). Second, we used a classification scheme based on the TP0548 gene. For the TP0548 classification scheme we carried out PCR and Sanger sequencing of the TP0548 gene region following the protocols and primers of Matějková et al.[31]. Single nucleotide polymorphisms (SNPs) in the TP0548 typing regions enable the distinction of an SS14 clade versus a Nichols clade. Indels enable the classification of TPE and TEN. Our NGS prediction strategy (detailed in the Supplementary Materials and Methods) was congruent with the TP0548 classification scheme wherever prediction strength was above 0.4, with the exception of 1 TEN sample.

*Samples from typing studies:* We put together all publicly available TP0548 sequences obtained in typing studies of syphilis infections around the world [12,53–60]. We additionally incorporated TP0548 sequences obtained for 34 Argentinian clinical samples by LGV at the

359    University of Buenos Aires, Argentina (Supplementary Table 8). All TP0548 sequences were
360    classified as part of the SS14 clade or part of the Nichols clade based on an ML tree
361    (Supplementary Fig. 5). Subtypes were distinguished through visual inspection
362    (Supplementary Table 8).

363    **Antibiotic resistance**
364    The two mutations associated with resistance to the macrolide azithromycin, A2058G
365    and A2059G on the 23S ribosomal RNA operon (with positions referring to coordinates in the
366    23S ribosomal RNA gene of *Escherichia coli*), were investigated in separate analyses. Since
367    the operon contains two copies of the gene, mapping of reads with BWA was carried out
368    independently for each of the genes, including a flanking region of 200 bases on both the 5'
369    and 3' end of each genes. Following variant calling, the presence/absence of each of the two
370    mutations was recorded for each sample. The two operons could not, however, be
371    distinguished.
372    In addition, we used primers specific for each of the two operons to carry out PCR
373    amplifications as well as Sanger sequencing on the samples, following the protocol in
374    Matějková et al.[31]. Details on the samples sequenced, as well as resistance or sensitivity to
375    the macrolide as determined by the presence or absence of the associated mutations are
376    given in Supplementary Table 7.

377    **Data availability**
378    All samples sequenced in this study are available in an NCBI Bioproject under accession
379    number PRJNA313497. Raw sequencing reads in FASTQ format were uploaded to the Short
380    Read Archive (SRA). All accession codes are listed in Supplementary Table 2. Code for the in-
381    house scripts developed for some of the analyses are available upon request from the
382    authors.

383 **Reference List**

384 1. Gall, G. E. C., Lautenschlager, S. & Bagheri, H. C. Quarantine as a public health measure against an

385 emerging infectious disease: syphilis in Zurich at the dawn of the modern era (1496–1585). *GMS Hyg Infect*

386 *Control* **11,** (2016).

387 2. Rowley, J. *et al. Global incidence and prevalence of selected curable sexually transmitted infections, 2008.*

388 (World Health Organization, 2012).

389 3. Stamm, L. V. Global Challenge of Antibiotic-Resistant *Treponema pallidum. Antimicrob. Agents*

390 *Chemother.* **54,** 583–589 (2010).

391 4. Fraser, C. M. *et al.* Complete genome sequence of *Treponema pallidum*, the syphilis spirochete. *Science*

392 **281,** 375–88 (1998).

393 5. Quétel, C. *History of syphilis.* (Johns Hopkins University Press, 1990).

394 6. Fernandez de Oviedo y Valdes, G. *Sumario de la natural historia de las Indias.* (Fondo de Cultura

395 Economico, 1526).

396 7. Harper, K. N., Zuckerman, M. K., Harper, M. L., Kingston, J. D. & Armelagos, G. J. The origin and

397 antiquity of syphilis revisited: an appraisal of Old World pre-Columbian evidence for treponemal infection.

398 *Am. J. Phys. Anthropol.* **146 Suppl 53,** 99–133 (2011).

399 8. Šmajs, D., Norris, S. J. & Weinstock, G. M. Genetic diversity in *Treponema pallidum*: Implications for

400 pathogenesis, evolution and molecular diagnostics of syphilis and yaws. *Infect. Genet. Evol.* **12,** 191–202

401 (2012).

402 9. Giacani, L. *et al.* Complete Genome Sequence of the *Treponema pallidum* subsp. *pallidum* Sea81-4 Strain.

403 *Genome Announc.* **2,** e00333-14-e00333-14 (2014).

404 10. Štaudová, B. *et al.* Whole Genome Sequence of the Treponema pallidum subsp. endemicum Strain Bosnia

405 A: The Genome Is Related to Yaws Treponemes but Contains Few Loci Similar to Syphilis Treponemes.

406 *PLoS Negl. Trop. Dis.* **8,** e3261 (2014).

407 11. Marra, C. M. *et al.* Enhanced Molecular Typing of *Treponema pallidum*: Geographical Distribution of

408 Strain Types and Association with Neurosyphilis. *J. Infect. Dis.* **202,** 1380–1388 (2010).

409 12. Grillová, L. *et al.* Molecular Typing of *Treponema pallidum* in the Czech Republic during 2011 to 2013:

410 Increased Prevalence of Identified Genotypes and of Isolates with Macrolide Resistance. *J. Clin. Microbiol.*

411 **52,** 3693–3700 (2014).

412 13. Hodges, E. *et al.* Hybrid selection of discrete genomic intervals on custom-designed microarrays for

413 massively parallel sequencing. *Nat. Protoc.* **4,** 960–974 (2009).

414    14. Meyer, M. & Kircher, M. Illumina Sequencing Library Preparation for Highly Multiplexed Target Capture

415        and Sequencing. *Cold Spring Harb. Protoc.* **2010,** pdb.prot5448-prot5448 (2010).

416    15. Pětrošová, H. *et al.* Resequencing of *Treponema pallidum* ssp. *pallidum* Strains Nichols and SS14:

417        Correction of Sequencing Errors Resulted in Increased Separation of Syphilis Treponeme Subclusters. *PLoS*

418        *ONE* **8,** e74319 (2013).

419    16. Petrosova, H. *et al.* Whole genome sequence of *Treponema pallidum* ssp. *pallidum*, strain Mexico A,

420        suggests recombination between yaws and syphilis strains. *PLoS Negl Trop Dis* **6,** e1832 (2012).

421    17. Centurion-Lara, A. *et al.* Fine Analysis of Genetic Diversity of the tpr Gene Family among Treponemal

422        Species, Subspecies and Strains. *PLoS Negl. Trop. Dis.* **7,** e2222 (2013).

423    18. Mikalova, L. *et al.* Genome analysis of *Treponema pallidum* subsp. *pallidum* and subsp. pertenue strains:

424        most of the genetic differences are localized in six regions. *PLoS One* **5,** e15713 (2010).

425    19. Achtman, M. Evolution, Population Structure, and Phylogeography of Genetically Monomorphic Bacterial

426        Pathogens. *Annu. Rev. Microbiol.* **62,** 53–70 (2008).

427    20. Bouckaert, R. *et al.* BEAST 2: A Software Platform for Bayesian Evolutionary Analysis. *PLoS Comput.*

428        *Biol.* **10,** e1003537 (2014).

429    21. Lukehart, S. A. & Giacani, L. When Is Syphilis Not Syphilis? Or Is It?: *Sex. Transm. Dis.* **41,** 554–555

430        (2014).

431    22. Mikalova, L. *et al.* Genome analysis of Treponema pallidum subsp. pallidum and subsp. pertenue strains:

432        most of the genetic differences are localized in six regions. *PLoS One* **5,** e15713 (2010).

433    23. Stadler, T., Kuhnert, D., Bonhoeffer, S. & Drummond, A. J. Birth-death skyline plot reveals temporal

434        changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc. Natl. Acad. Sci.* **110,** 228–233

435        (2013).

436    24. Holt, K. E. *et al.* Shigella sonnei genome sequencing and phylogenetic analysis indicate recent global

437        dissemination from Europe. *Nat Genet* **44,** 1056–1059 (2012).

438    25. Mutreja, A. *et al.* Evidence for several waves of global transmission in the seventh cholera pandemic.

439        *Nature* **477,** 462–465 (2011).

440    26. Giacani, L. *et al.* Footprint of Positive Selection in T*reponema pallidum* subsp. *pallidum* Genome

441        Sequences Suggests Adaptive Microevolution of the Syphilis Pathogen. *PLoS Negl. Trop. Dis.* **6,** e1698

442        (2012).

443    27. Strouhal, M. *et al.* Genome Differences between *Treponema pallidum* subsp. *pallidum* Strain Nichols and T.

444        paraluiscuniculi Strain Cuniculi A. *Infect. Immun.* **75,** 5859–5866 (2007).

445 28. Marra, C. M. *et al.* Antibiotic selection may contribute to increases in macrolide-resistant Treponema

446       pallidum. *J. Infect. Dis.* **194,** 1771–1773 (2006).

447 29. Mitchell, S. J. *et al.* Azithromycin-resistant syphilis infection: San Francisco, California, 2000–2004. *Clin.*

448       *Infect. Dis.* **42,** 337–345 (2006).

449 30. Morshed, M. & Jones, H. Treponema pallidum macrolide resistance in BC. *CMAJ* 349 (2006).

450 31. Matejkova, P. *et al.* Macrolide treatment failure in a case of secondary syphilis: a novel A2059G mutation

451       in the 23S rRNA gene of Treponema pallidum subsp. pallidum. *J. Med. Microbiol.* **58,** 832–836 (2009).

452 32. Stamm, L. V. & Bergen, H. L. A Point Mutation Associated with Bacterial Macrolide Resistance Is Present

453       in Both 23S rRNA Genes of an Erythromycin-Resistant *Treponema pallidum* Clinical Isolate. *Antimicrob.*

454       *Agents Chemother.* **44,** 806–807 (2000).

455 33. Šmajs, D., Paštěková, L. & Grillová, L. Macrolide Resistance in the Syphilis Spirochete, Treponema

456       pallidum ssp. pallidum: Can We Also Expect Macrolide-Resistant Yaws Strains? *Am. J. Trop. Med. Hyg.*

457       **93,** 678–683 (2015).

458 34. Geisler, W. M. *et al.* Azithromycin versus Doxycycline for Urogenital *Chlamydia trachomatis* Infection. *N.*

459       *Engl. J. Med.* **373,** 2512–2521 (2015).

460 35. Centurion-Lara, A. in *Pathogenic Treponema: Molecular and Cellular Biology* (eds. Radolf, J. D. &

461       Lukehart, S. A.) 267–283 (Caister Academic Press, 2006).

462 36. Knell, R. J. Syphilis in renaissance Europe: rapid evolution of an introduced sexually transmitted disease?

463       *Proc Biol Sci* **271 Suppl 4,** S174-6 (2004).

464 37. Kircher, M., Sawyer, S. & Meyer, M. Double indexing overcomes inaccuracies in multiplex sequencing on

465       the Illumina platform. *Nucleic Acids Res.* **40,** e3–e3 (2012).

466 38. Peltzer, A. *et al.* EAGER: Efficient Ancient Genome Reconstruction. *Genome Biol.* **17,** 1 (2016).

467 39. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *ArXiv Prepr.*

468       *ArXiv13033997* (2013).

469 40. McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation

470       DNA sequencing data. *Genome Res.* **20,** 1297–1303 (2010).

471 41. Simpson, J. T. & Durbin, R. Efficient de novo assembly of large genomes using compressed data structures.

472       *Genome Res.* **22,** 549–556 (2012).

473 42. Schmidt, H. A., Strimmer, K., Vingron, M. & von Haeseler, A. TREE-PUZZLE: maximum likelihood

474       phylogenetic analysis using quartets and parallel computing. *Bioinformatics* **18,** 502–504 (2002).

475    43.   Strimmer, K. & von Haeseler, A. Quartet Puzzling: A Quartet Maximum-Likelihood Method for

476        Reconstructing Tree Topologies. *Mol. Biol. Evol.* **13,** 964 (1996).

477    44.   Strimmer, K. & Rambaut, A. Inferring confidence sets of possibly misspecified gene trees. *Proc. R. Soc. B*

478        *Biol. Sci.* **269,** 137–142 (2002).

479    45.   Shimodaira, H. & Hasegawa, M. Multiple Comparisons of Log-Likelihoods with Applications to

480        Phylogenetic Inference. *Mol. Biol. Evol.* **16,** 1114 (1999).

481    46.   Tamura, K., Stecher, G., Peterson, D., Filipski, A. & Kumar, S. MEGA6: Molecular Evolutionary Genetics

482        Analysis Version 6.0. *Mol. Biol. Evol.* **30,** 2725–2729 (2013).

483    47.   Bouckaert, R. *et al.* BEAST 2: A Software Platform for Bayesian Evolutionary Analysis. *PLoS Comput.*

484        *Biol.* **10,** e1003537 (2014).

485    48.   Baele, G. *et al.* Improving the Accuracy of Demographic and Molecular Clock Model Comparison While

486        Accommodating Phylogenetic Uncertainty. *Mol. Biol. Evol.* **29,** 2157–2167 (2012).

487    49.   Rambaut, A., Suchard, M., Xie, D. & Drummond, A. *Tracer v1.6.* (2014).

488    50.   Rambaut, A. *FigTree v.1.4.2.* (2014).

489    51.   Bandelt, H.-J., Forster, P. & Röhl, A. Median-joining networks for inferring intraspecific phylogenies. *Mol.*

490        *Biol. Evol.* **16,** 37–48 (1999).

491    52.   *www.fluxux-engineering.com.*

492    53.   Dai, T. *et al.* Molecular Typing of Treponema pallidum: a 5-Year Surveillance in Shanghai, China. *J. Clin.*

493        *Microbiol.* **50,** 3674–3677 (2012).

494    54.   Flasarová, M. *et al.* Sequencing-based Molecular Typing of *Treponema pallidum* Strains in the Czech

495        Republic: All Identified Genotypes are Related to the Sequence of the SS14 Strain. *Acta Derm. Venereol.*

496        **92,** 669–674 (2012).

497    55.   Grange, P. A. *et al.* Molecular Subtyping of Treponema pallidum in Paris, France: *Sex. Transm. Dis.* **40,**

498        641–644 (2013).

499    56.   Grimes, M. *et al.* Two Mutations Associated With Macrolide Resistance in Treponema pallidum: Increasing

500        Prevalence and Correlation With Molecular Strain Type in Seattle, Washington. *Sex. Transm. Dis.* **39,** 954–

501        958 (2012).

502    57.   Peng, R.-R. *et al.* Molecular Typing of *Treponema pallidum* Causing Early Syphilis in China: A Cross-

503        Sectional Study: *Sex. Transm. Dis.* **39,** 42–45 (2012).

504    58.   Tian, H. *et al.* Molecular typing of *Treponema pallidum*: identification of a new sequence of tp0548 gene in

505        Shandong, China. *Sex. Transm. Dis.* **41,** 551 (2014).

506    59. Tipple, C., McClure, M. O. & Taylor, G. P. High prevalence of macrolide resistant Treponema pallidum

507          strains in a London centre. *Sex. Transm. Infect.* **87,** 486–488 (2011).

508    60. Wu, B.-R. *et al.* Multicentre surveillance of prevalence of the 23S rRNA A2058G and A2059G point

509          mutations and molecular subtypes of *Treponema pallidum* in Taiwan, 2009–2013. *Clin. Microbiol. Infect.*

510          **20,** 802–807 (2014).

511

**Additional information**

Supplementary information is available for this paper. Correspondence and requests for materials should be addressed to N.A. (natasha.arora@uzh.ch), F.G.C. (fernando.gonzalez@uv.es),K.N. (kay.nieselt@uni-tuebingen.de), J.K. (johannes.krause@uni-tuebingen.de) or H.C.B (homayoun.bagheri@repsol.com).

**Accession codes**

All raw read files have been deposited in the trace archive of the NCBI Sequence Read Archive under accession number SRP072086.

**Competing interests**

The authors declare no competing financial interests.

**Figure Legends**

**Figure 1 | De novo genome assemblies and phylogenetic reconstruction. a**, De novo genome assembly for four syphilis patient samples and one yaws strain, with color coded geographic origin (inset legend). Blank spaces correspond to gaps, overlapping with gene regions that are difficult to assemble from short reads such as the tpr subfamilies and rRNA operons (regions shown in the outermost ring in gray). **b**, BEAST tree for the 39 genomes (excluding putative recombinant genes), with black circles for nodes with ≥96% posterior probabilities (PP); dark gray circles for nodes with 91-95% PP; and white circles for nodes with 81-85% PP. Divergence date estimates (mean and 95% highest posterior density) for major well-supported TPA nodes are given in the legend.

**Figure 2│Median-joining (MJ) network analysis and geographic distribution of the SS14 and Nichols clades**. **a**, Median-joining network for genome-wide variable positions after excluding sites with missing data (n=682). Circles represent haplotypes, with geographical origin color-coded. Number of mutations, when above one, is shown next to the lines. Inferred haplotypes (median vectors) are shown as black connecting circles. Central black

557 circles within haplotypes indicate mutations associated with azithromycin resistance. **b**,
558 Relative frequencies of SS14 versus Nichols clade isolates across the globe shown in the pie
559 charts, with sizes proportional to sampling efforts. SS14 clade and Nichols classification are
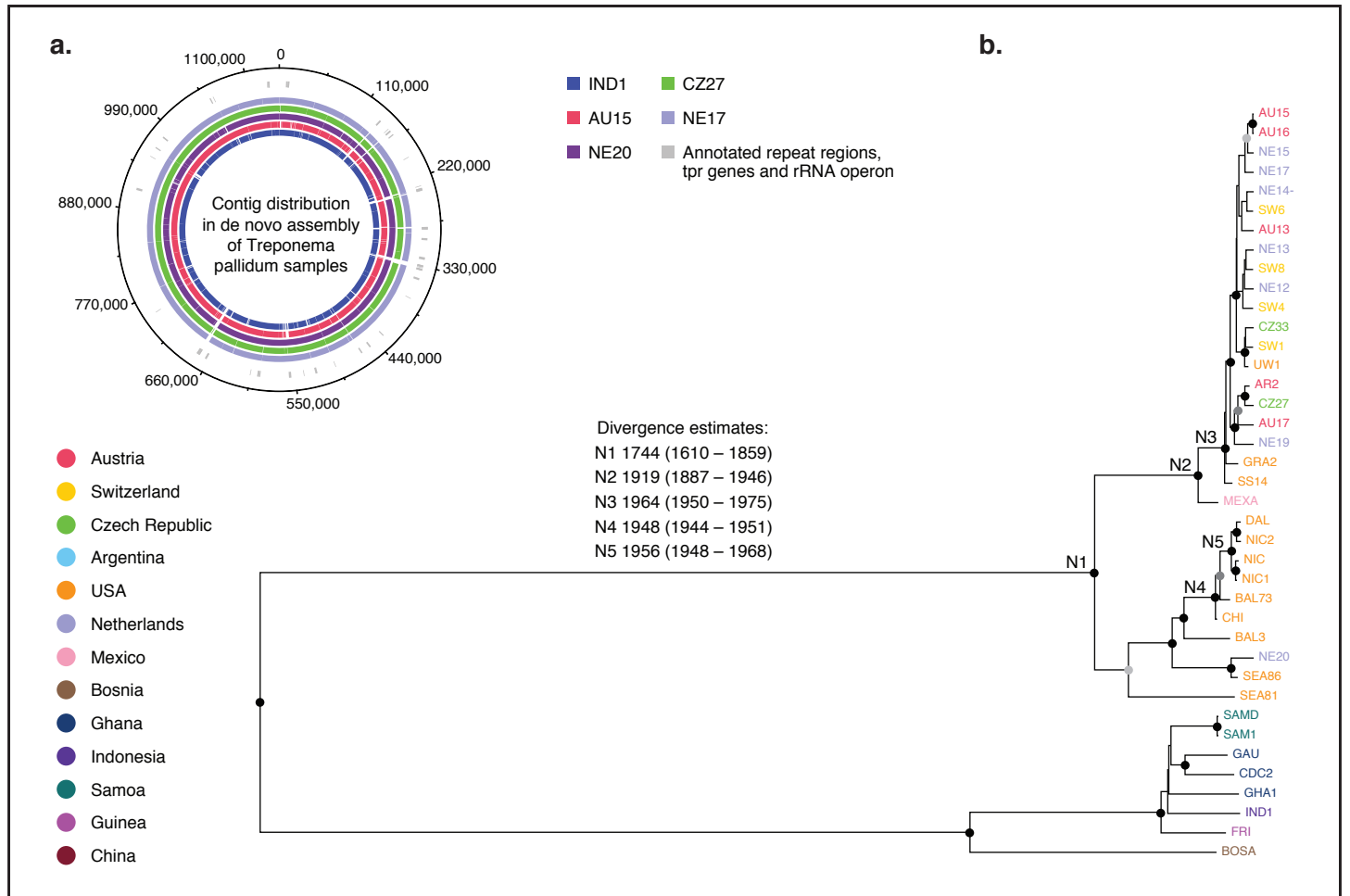560 based on the TP0548 gene.

FIGURE 1

# FIGURE 2



**a.**

SEA81
(1981)

SEA86
(1986)

36

NE20
(2013)

2

3

19

13

4

11

BAL3
(1973)

MEXA
(1953)

8

48

**SS14 - Ω cluster**

NE19
(2013)

SW6
(2012)

AU17
(2013)

GRA2

7

2

5

SS14
(1977)

2

4

AR2
(2006)

5

CZ27
(2013)

NE13
(2013)

NE15
(2013)

2

SS14 - Ω haplotype
(2004, 2012, 2013)

AU15, AU16
(2013)

NE17
(2013)

3

6

2

SW8
(2012)

NE12
(2013)

11

SW4
(2012)

CZ33
(2013)

20

**Nichols - α cluster**

NIC_REF

NIC1

7

CHI
(1951)

DAL
(1991)

NIC2

BAL73
(1973)

- ● Austria
- ● Switzerland
- ● Czech Republic
- ● Argentina
- ● USA
- ● Netherlands
- ● Mexico
- ◉ Antibiotic Resistant

**b.**

UK
(London)

DENMARK

IRELAND
(Dublin)

CZECH
REPUBLIC

USA

NETHERLANDS

AUSTRIA

CHINA

(Shandong)

FRANCE
(Paris)

SWITZERLAND

(Nanjing)

(Shanghai)

TAIWAN

- ● SS14 clade
- ○ Nichols clade
- * Samples taken
  from multiple sites
- ○ < 100 samples
- ○ 101 - 200 samples
- ○ > 200 samples

MADAGASCAR

ARGENTINA