

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

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

**The abrupt onslaught of the syphilis pandemic that started in the late fifteenth century established this devastating infectious disease as one of the most feared in human history<sup>1</sup>. Surprisingly, despite the availability of effective antibiotic treatment since the mid-twentieth century, this bacterial infection, which is caused by *Treponema pallidum* subsp. *pallidum* (TPA), has been re-emerging globally in the last few decades with an estimated 10.6 million cases in 2008 (ref. 2). Although resistance to penicillin has not yet been identified, an increasing number of strains fail to respond to the second-line antibiotic azithromycin<sup>3</sup>. Little is known about the genetic patterns in current infections or the evolutionary origins of the disease due to the low quantities of treponemal DNA in clinical samples and difficulties in cultivating the pathogen<sup>4</sup>. Here, we used DNA capture and whole-genome sequencing to successfully interrogate genome-wide variation from syphilis patient specimens, combined with laboratory samples of TPA and two other subspecies. Phylogenetic comparisons based on the sequenced genomes indicate that the TPA strains examined share a common ancestor after the fifteenth century, within the early modern era. Moreover, most contemporary strains are azithromycin-resistant and are members of a globally dominant cluster, named here as SS14-Ω. The cluster diversified from a common ancestor in the mid-twentieth century subsequent to the discovery of antibiotics. Its recent phylogenetic divergence and global presence point to the emergence of a pandemic strain cluster.**

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

The paucity of molecular studies and the focus on typing of a few genes means that we have limited information regarding the

<sup>1</sup>Institute for Evolutionary Biology and Environmental Studies, University of Zurich, 8057 Zurich, Switzerland. <sup>2</sup>Zurich Institute of Forensic Medicine, University of Zurich, 8057 Zurich, Switzerland. <sup>3</sup>Institute for Archaeological Sciences, University of Tübingen, 72070 Tübingen, Germany. <sup>4</sup>Center for Bioinformatics, University of Tübingen, 72076 Tübingen, Germany. <sup>5</sup>Department of Biology, Faculty of Medicine, Masaryk University, 625 00 Brno, Czech Republic. <sup>6</sup>Unidad Mixta Infección y Salud Pública FISABIO/Universidad de Valencia; CIBER in Epidemiology and Public Health, 46020, Spain. <sup>7</sup>Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SA, UK. <sup>8</sup>Institute of Integrative Biology, Department of Environmental Systems Science, ETH Zürich, 8092 Zurich, Switzerland. <sup>9</sup>Department of Infectious Diseases, Public Health Laboratory, GGD Amsterdam, 1018 WT Amsterdam, the Netherlands. <sup>10</sup>Department of Dermatology, Medical University of Graz, A-8036 Graz, Austria. <sup>11</sup>The Mortimer Market Centre CNWL, Camden Provider Services, London NW1 2PL, UK. <sup>12</sup>Department of Clinical Microbiology and Virology, University College London Hospitals NHS Foundation Trust, London W1T 4EU, UK. <sup>13</sup>Instituto de Investigaciones Biomédicas en Retrovirus y SIDA (INBIRS), Universidad de Buenos Aires-CONICET, 1121 Buenos Aires, Argentina. <sup>14</sup>Facultad de Farmacia y Bioquímica, Departamento de Bioquímica Clínica, Microbiología Clínica, Universidad de Buenos Aires, 1113 Buenos Aires, Argentina. <sup>15</sup>Servicio de Dermatología, Hospital General Universitario de Valencia, 46014 Valencia, Spain. <sup>16</sup>Department of Medicine, Division of Allergy and Infectious Diseases, and Department of Global Health, University of Washington, Seattle, Washington 98105, USA. <sup>17</sup>Department of Pathology and Laboratory Medicine, UTHealth McGovern Medical School, Houston, Texas 77225, USA. <sup>18</sup>Department of Dermatology, University Hospital of Zurich, 8091 Zurich, Switzerland. <sup>†</sup>Present address: Department of Archaeogenetics, Max Planck Institute for the Science of Human History, D-07745 Jena, Germany (A.P., A.H., K.I.B., J.K.); Department of Infectious Disease Epidemiology, Imperial College London, London SW7 2AZ, UK (L.R.D.); Repsol Technology Center, 28935 Mostoles, Madrid, Spain (H.C.B.). \*e-mail: [natasha.arora@uzh.ch](mailto:natasha.arora@uzh.ch); [fernando.gonzalez@uv.es](mailto:fernando.gonzalez@uv.es); [kay.nieselt@uni-tuebingen.de](mailto:kay.nieselt@uni-tuebingen.de); [krause@shh.mpg.de](mailto:krause@shh.mpg.de); [homayoun.bagheri@repsol.com](mailto:homayoun.bagheri@repsol.com)

evolution and spread of epidemic TPA. In this study, we interrogated genome-wide variation across geographically widespread isolates. In total, we obtained 70 samples from 13 countries, including 52 syphilis swabs collected directly from patients between 2012 and 2013, and 18 syphilis, yaws and bejel samples collected from 1912 onwards and propagated in laboratory rabbits (Supplementary Table 1). Through comparative genome analyses and phylogenetic reconstruction, we shed light on the evolutionary history of TPA and identify epidemiologically relevant haplotypes.

Due to the large background of host DNA, samples were enriched for treponemal DNA prior to Illumina sequencing<sup>13,14</sup>. The resultant reads were mapped to the Nichols TPA reference genome (RefSeq NC\_021490; Supplementary Table 3)<sup>4,15</sup>. Genomic coverage ranged from 0.13-fold to over 1,000-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. We therefore 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 proportions of genome coverage with at least 3-fold or 10-fold depth were 97% and 82%, respectively (Supplementary Table 4).

*De novo* assemblies for the four highest covered syphilis swab samples (NE17, NE20, CZ27 and 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<sup>16</sup>. The deletion was shared across all the yaws infection isolates (Supplementary Methods), consistent with other studies<sup>17</sup>.

Before phylogenetic reconstruction we checked for signatures of recombination. *T. pallidum* is considered to be a clonal species<sup>18</sup>, but previous studies suggest recombinant genes in a Mexican syphilis and a Bosnian bejel strain<sup>10,19</sup>. 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 single nucleotide polymorphism (SNP) densities, incongruent topologies with the genome-wide tree and more than four homoplasies in a pair of branches (Supplementary Methods). We identified four genes coding for outer membrane proteins (Supplementary Table 6), one of which (TP0136) is used in typing studies<sup>8</sup>.

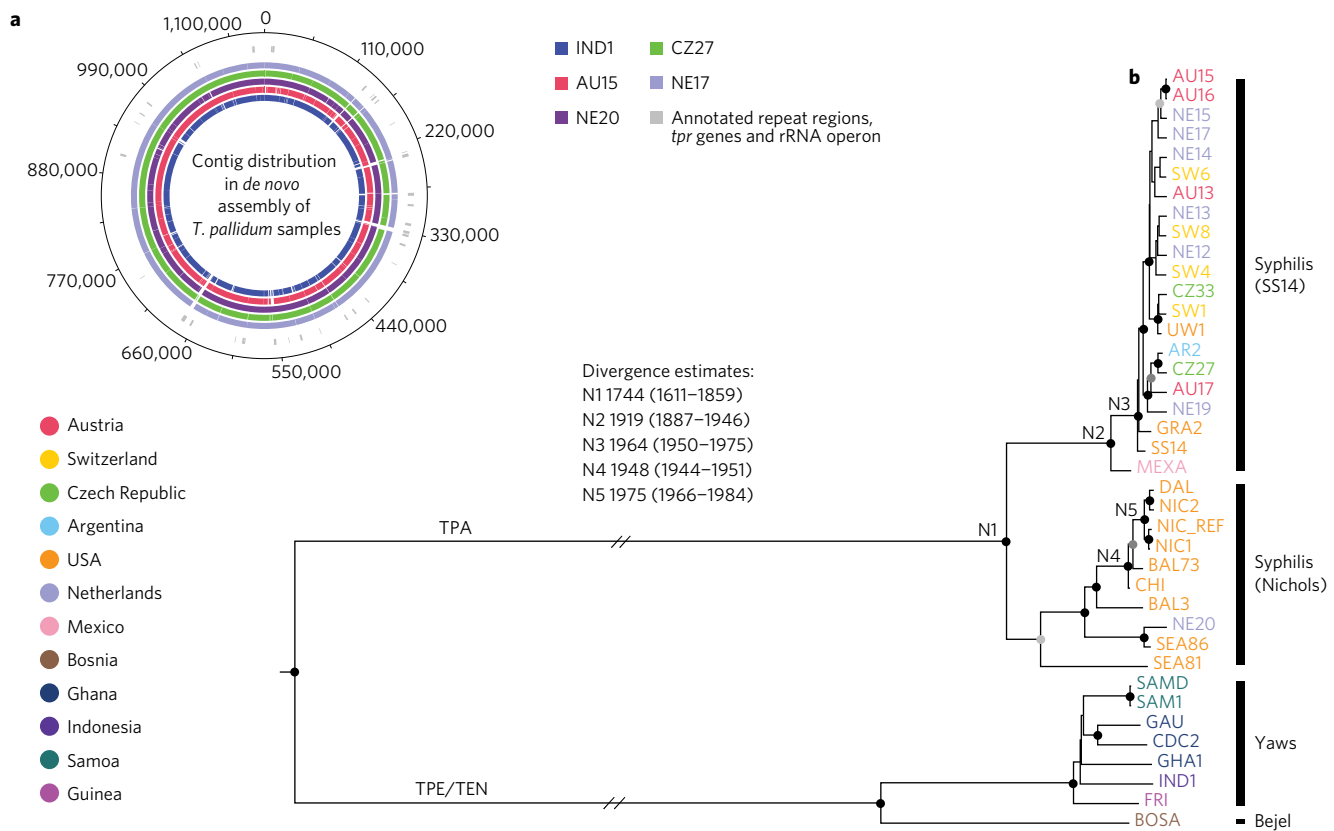
After excluding the four putative recombinant genes, the genome alignment for all 39 genomes contained 2,235 variable positions. We used the Bayesian framework implemented in BEAST<sup>20</sup> to reconstruct a phylogenetic tree (Fig. 1b). The tree topology revealed a marked separation between TPA and 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 1,225 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 1,703 to TPE/TEN genomes. Previous studies have found cross-subspecies groupings when relying on a limited set of markers<sup>21</sup>. 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<sup>10,17</sup>, but also illustrate

the need for a careful choice of taxonomic markers when genome-wide data are not available.

Using the sample isolation dates as tip calibration and applying the Birth Death Serial Skyline model<sup>22</sup>, we obtained a mean evolutionary rate of  $3.6 \times 10^{-4}$  (rate variance  $3.8 \times 10^{-8}$ ; 95% highest posterior density (HPD) of  $1.86 \times 10^{-4}$  to  $5.73 \times 10^{-4}$ ). This estimate is equivalent to a 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;  $8.0 \times 10^{-7}$ )<sup>23,24</sup>. Our divergence analyses for TPA samples provide a time to the most recent common ancestor (TMRCA) of less than 500 years ago (mean calendar year 1744, 95% HPD 1611–1859; Fig. 1b).

Within the TPA lineage the samples group in two clades named after the SS14 and Nichols reference genomes (with 100% and 82% posterior probability values, respectively). The Nichols clade consists almost exclusively of samples collected from patients in North America from 1912 to 1986 and passed in rabbits before sequencing, with the exception of one patient sample from 2013 (NE20). In contrast, the SS14 clade has a geographically widespread distribution, encompassing European, North American and South American samples collected from infections between 1951 and 2013. We investigated the TPA clades further by generating a median-joining (MJ) network to illustrate the mutational differences among the TPA samples (Fig. 2a). As underscored by distances in the network, greater nucleotide diversity is found within the Nichols clade ( $\pi = 0.05$ ) than in the SS14 clade ( $\pi = 0.01$ ). Three closely related sequences derive from the original Nichols sample isolated from the cerebrospinal fluid of a patient in 1912 and propagated in laboratories in subsequent decades: NIC\_REF, the reference genome re-sequenced by Pětrošová *et al.*<sup>15</sup>, and NIC-1 and NIC-2, which we sequenced following independent propagation of the strains in Houston and Seattle, respectively, during different time periods (Supplementary Table 1 and Supplementary Table 3). These three group together with another three laboratory-propagated strains in a cluster labelled Nichols- $\alpha$  (Fig. 2a), with a TMRCA in the mid-twentieth century (Fig. 1a). The less diversified SS14 clade contains a dominant central haplotype (labelled SS14- $\Omega$ ) from which the other sequences radiate (Fig. 2a). Critically, the cluster associated with the SS14- $\Omega$  haplotype contains all but one of the recent patient samples from 2012 to 2013 ( $n = 17$ ) that were captured and sequenced directly, in addition to samples from 1977 ( $n = 1$ ) and 2004 ( $n = 2$ ). The genetic variation within the SS14- $\Omega$  cluster is found primarily as singleton mutations (95.5%), with no evidence here for geographical structuring. Bayesian analyses estimate coalescence for the SS14- $\Omega$  cluster in 1964 (mean calendar year; 95% HPD 1950–1975; Fig. 1b), at a time when incidence was reduced due to the introduction of antibiotics. The star-like topology of this cluster observed in both the tree and the network is suggestive of a recent and rapid clonal expansion.

To determine whether the dominance of SS14 clade sequences applies across other countries for which genetic data are available, we examined sequences from the widely typed TP0548 gene in worldwide epidemiological studies<sup>11</sup>. Phylogenies for the TP0548 typing regions separate the SS14 from the Nichols clade for the TPA samples, but do not distinguish the TPA and TPE/TEN lineages (Supplementary Methods and Supplementary Fig. 4). Across 1,354 worldwide TP0548 sequences from clinical samples, including the 78 from patients in this study, we found that 94% of them grouped in the SS14 clade (Supplementary Tables 8 and 9 and Supplementary Fig. 5), consistent with a probable recent spread of the epidemic cluster. The wide geographical distribution of the SS14 clade establishes it as representative of the present worldwide epidemic. Studies so far have focused on the Nichols strain<sup>25,26</sup>, but our results indicate that further work on the SS14 clade is warranted.



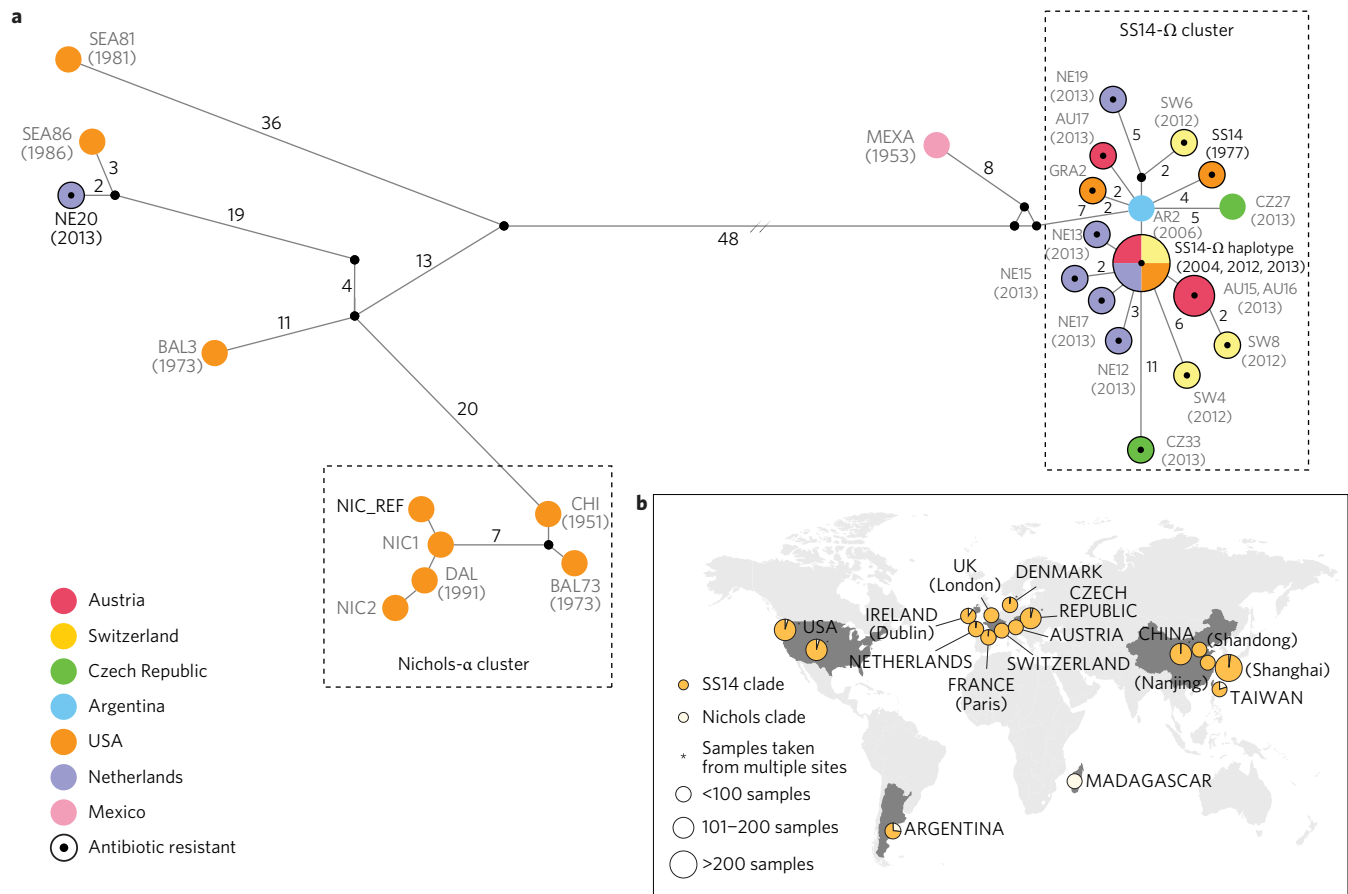
**Figure 1 | De novo genome assemblies and phylogenetic reconstruction.** **a**, *De novo* genome assembly for four syphilis patient samples and one yaws strain, with colour-coded geographic origin. Blank spaces correspond to gaps, overlapping with gene regions that are difficult to assemble from short reads such as the *trp* subfamilies and rRNA operons (regions shown in the outermost ring in grey). **b**, BEAST tree for the 39 genomes (excluding putative recombinant genes), with black circles for nodes with  $\geq 96\%$  posterior probabilities (PP), dark grey circles for nodes with 91–95% PP and light grey circles for nodes with 80–90% PP. Divergence date estimates (mean and 95% highest posterior density) for major well-supported TPA nodes are given in the legend.

Critically, typing of samples over multiple years in the Czech Republic, San Francisco, British Columbia and Seattle indicate that macrolide antibiotic resistance has increased over time<sup>3,12,27–29</sup>. We queried the presence of the two mutations (A2058G and A2059G) in the 23S ribosomal RNA (rRNA) genes associated with azithromycin resistance<sup>3,30,31</sup>. As observed in the MJ network, the resistance marker is a dominant characteristic of the SS14-Ω cluster (Fig. 2a), although it is also found in a recent patient sample (NE20) of the Nichols clade. Extending our analyses of the 23S rRNA gene to all sequenced samples from our study, including the 42 with lower coverage, revealed the mutations in 90% of the SS14 ( $n = 51$ ) and 25% of the Nichols ( $n = 12$ ) samples, indicating that neither resistance nor sensitivity is clade-specific (Supplementary Table 8). A likely scenario is that the extensive usage of azithromycin to treat syphilis and a wide range of bacterial infections, including co-infections with other sexually transmitted diseases (STDs) such as chlamydia, has played an important role in the selection and subsequent spread of resistance<sup>32,33</sup>.

The results here represent the first reported set of whole-genome sequences successfully obtained directly from syphilis patients, enabling us to disentangle evolutionary relationships at high resolution and paving the way for further clinical sequencing from current epidemics. Given our identification of putative recombinant genes in *Treponema* and previous reports on genes involved in homologous recombination<sup>4,34</sup>, further detailed analyses on the potential mechanisms of recombination will be necessary. Our phylogenetic reconstruction indicates that all TPA samples examined to date share a common ancestor that was infecting populations in the 1700s, within the early centuries of the modern era, and that

was successful in leaving descendants until today. This date is posterior to the colonization of the Americas and therefore potentially compatible with the post-Columbian model for the emergence of syphilis in Europe. Nonetheless, our work does not exclude the possibility that older TPA lineages had previously existed in Europe but went extinct. Obtaining more patient sample genomes with high coverage could refine our detection of putative recombinants and our phylogenetic inferences. In addition, sequencing from ancient skeletal material would help to further ascertain the history of syphilis. Interestingly, we observed a time difference between the first reported syphilis outbreak in 1495 and the last common ancestor of modern strains dated to the 1700s. Although this difference could stem from imprecision in the divergence estimates, an alternative scenario is the eventual establishment of a specific lineage due to selection. For instance, it has been hypothesized that the symptoms of syphilis became less severe after the first reported outbreaks in Europe because of the evolution of strains with lower virulence and higher transmission rates<sup>35</sup>. In this scenario, the eighteenth century provided the context for the origin and propagation of a lineage that successfully outcompeted other lineages.

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



**Figure 2 | Median-joining (MJ) network analysis and geographic distribution of the SS14 and Nichols clades.** **a**, MJ network for genome-wide variable positions after excluding sites with missing data ( $n = 682$ ). Coloured circles represent haplotypes, with colours representing the geographical origin. The number of mutations, when above one, is shown next to the lines. Inferred haplotypes (median vectors) are shown as black connecting circles. Central black circles within haplotypes indicate mutations associated with azithromycin resistance. **b**, Relative frequencies of SS14 versus Nichols clade isolates across the globe are shown by the pie charts, with sizes proportional to sampling efforts. The SS14 and Nichols clade classifications are based on the *TP0548* gene.

instrumental in understanding the genetic and epidemiological basis for the spread of SS14- $\Omega$  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) 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<sup>14,36</sup>, at the University of Tübingen (Supplementary 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 and co-authors<sup>13</sup>. 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<sup>37</sup>, 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 Methods). All reads (merged and unmerged) were treated as single-end reads and mapping was performed using the BWA-MEM algorithm<sup>38</sup> with default parameters, using the Nichols genome as reference. Subsequently, we selected the samples that had at least 80% coverage of the Nichols genome and a minimum of three reads ( $n = 28$  samples, Supplementary Table 2). For each of these samples, we used the Genome Analysis Toolkit (GATK)<sup>39</sup> 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 position was covered by at least three reads, the Phred-scaled genotype quality score of the call was at least 30, and at least 90% of the reads agreed with the reference. A variant position (SNP) was called if the position was covered by at least three reads, the genotype quality of the call was at least 30, and the minimum SNP allele frequency was 90% (detailed in the Supplementary Methods). 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 that produced by the GATK UnifiedGenotyper and incorporates—for each row and thus for each call—one nucleotide into the new draft sequence.

To apply our analysis pipeline also to those samples for which complete genomic sequences are available in GenBank (Supplementary Table 3), we produced artificial reads in these cases using an in-house tool (Genome2Reads) and then applied the same mapping, SNP calling and genome reconstruction procedure as for the sequenced samples to obtain consistent and comparable results.

To investigate conservation of structure and gene order in the genomes, in addition to the mapping assembly we also performed a *de novo* assembly for the five samples with highest coverage (Supplementary Table 5). Our *de novo* assembly pipeline started with the merged reads, and in a first step used the short read assembler software SOAPdenovo2 using ten different  $k$ -mer sizes ( $k = 37 + 10i$ , where  $i = \{0, \dots, 9\}$ ). Different  $k$ -mer sizes were used because merging of read pairs into one single read results in very different lengths (between 30 and 190 bases). Next, all input reads were mapped back against the resulting contigs using

BWA-MEM (ref. 38). Contigs that were not supported by any reads (no read mapped against these contigs) were removed. To assemble the contigs resulting from the different *k*-mers, the remaining contigs were subject to the overlap-based String Graph Assembler (SGA)<sup>40</sup>. Finally, contigs smaller than 1,000 bp were removed before these contigs were mapped against the Nichols reference genome for comparison of genome architectures.

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

**Recombination detection.** Tests for the non-vertical transmission of genes were carried out on the TPA, TPE and TEN genomes ( $n = 39$ ) by identifying those genes that (1) had an unexpectedly high number of SNPs and (2) displayed patterns of transmission (that is, phylogenies) incongruent with most other genes. First, an expected substitution rate was computed by dividing the total number of observed SNPs in the 978 annotated genes ( $n = 2,098$ ) by the total length of these genes (1,046,421 bp). This rate was then used to calculate the expected number of polymorphisms per gene according to its length. A total of 87 genes displayed at least twice the expected number of polymorphisms. Second, for each of these 87 genes the gene sequence alignment and the gene tree topology were tested against the maximum likelihood tree topology of the draft genome in TREE-PUZZLE v5.2 (refs 41,42). Genes for which both the expected likelihood weight<sup>43</sup> and the Shimodaira–Hasegawa<sup>44</sup> test rejected the genome tree ( $P < 0.05$ ) were examined more closely. Third, genes within which we identified a minimum of five homoplasies (identical mutations in separate lineages) in at least two branches of the tree were marked as putative recombinants (Supplementary Table 6).

**Genome-wide variation and phylogenetic analyses.** We investigated genome-wide patterns of polymorphism and divergence using MEGA 6.0 (ref. 45) and DnaSP v.5.10 (ref. 46) to compute various measures of diversity including the average pairwise nucleotide differences, Nei's  $\pi$  ( $\pi$ ) and the number of singletons in each group. We also estimated the number of SNPs private to particular groups. A comparison of the TPA and TPE/TEN genomes revealed between 1 (NIC1) and 339 (AR2) SNPs observed in the TPA samples and between 1,091 (GHA1) and 1,443 (Bosnia A) SNPs in the TPE/TEN strains (Supplementary Table 4). Furthermore, we produced a heat map to display the number of SNPs (with respect to the Nichols reference genome) that any two genomes share (Supplementary Fig. 3).

The molecular clock hypothesis was tested with maximum likelihood analysis in MEGA 6.0 (ref. 45). Tests were conducted for all TPA, TPE and TEN genomes (39 samples) using (1) multiple whole genome alignments and (2) alignments with only the variable positions, in both cases excluding the four putative recombinant genes. The molecular clock hypothesis was rejected at the 5% significance level.

Bayesian phylogenetic trees were produced in BEAST 2.3 (ref. 20) 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 positions) and the tree generated with the set of variable positions after excluding the four putative recombinant genes (2,235 positions). Additionally, rooted trees were generated with maximum parsimony by including *Treponema paraluisuncinuli* (NC\_015714) as the outgroup.

As a calibration for the BEAST (Bayesian evolutionary analysis sampling trees) 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<sup>20,47</sup>. We report estimates for the five combined BDSS model runs, each with the following specifications: uncorrelated lognormal relaxed clock model, generalized time reversible plus gamma substitution model, and 100 million generations with parameter sampling every 10,000 generations. We used Tracer 1.6 (ref. 48) to assess convergence and suitable burn-in periods (see Supplementary Methods for full details). The annotated maximum clade credibility tree was visualized and edited using Figtree v1.4.2 (ref. 49). 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 an MJ network analysis in Network 4.6 and Network Publisher (<http://www.fluxus-engineering.com>)<sup>50</sup> 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 (for full details, see Supplementary Methods). 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. 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 ref. 30. 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 was congruent with the *TP0548* classification scheme wherever prediction strength was above 0.4, with the exception of one TEN sample (detailed in the Supplementary Methods).

**Samples from typing studies.** We put together all publicly available *TP0548* sequences obtained in typing studies of syphilis infections around the world<sup>12,51–58</sup>. We also incorporated *TP0548* sequences obtained for 36 Argentinian clinical samples by LGV at the University of Buenos Aires, Argentina (Supplementary Table 9). All *TP0548* sequences were classified as part of the SS14 clade or part of the Nichols clade based on an ML tree (Supplementary Fig. 5). Subtypes were distinguished through visual inspection (Supplementary Table 9).

**Antibiotic resistance.** The two mutations associated with resistance to the macrolide azithromycin, A2058G and A2059G on the 23S ribosomal RNA operon (with positions referring to coordinates in the 23S ribosomal RNA gene of *Escherichia coli*), were investigated in separate analyses. Because the operon contains two copies of the gene, mapping of reads with BWA was carried out independently for each of the genes, including a flanking region of 200 bases on both the 5' and 3' end of each gene. Following variant calling, the presence/absence of each of the two mutations was recorded for each sample. The two operons could not, however, be distinguished.

In addition, we used primers specific for each of the two operons to carry out PCR amplifications as well as Sanger sequencing on the samples, following the protocol in ref. 30. Details on the samples sequenced, as well as resistance or sensitivity to the macrolide as determined by the presence or absence of the associated mutations, are provided in Supplementary Table 8.

**Data availability.** All samples sequenced in this study are available in an NCBI Bioproject under accession code [PRJNA313497](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA313497). Raw sequencing reads in FASTQ format were uploaded to the Short Read Archive (SRA). All accession codes are listed in Supplementary Table 2. Codes for the in-house scripts developed for some of the analyses are available upon request from the authors. All raw read files have been deposited in the trace archive of the NCBI Sequence Read Archive under accession code [SRR072086](https://www.ncbi.nlm.nih.gov/sra/SRR072086).

Received 24 June 2016; accepted 3 November 2016;  
published 5 December 2016

#### References

- Gall, G. E. C., Lautenschlager, S. & Bagheri, H. C. Quarantine as a public health measure against an emerging infectious disease: syphilis in Zurich at the dawn of the modern era (1496–1585). *GMS Hyg. Infect. Control* **11**, 13 (2016).
- Rowley, J. et al. *Global Incidence and Prevalence of Selected Curable Sexually Transmitted Infections, 2008* (World Health Organization, 2012).
- Stamm, L. V. Global challenge of antibiotic-resistant *Treponema pallidum*. *Antimicrob. Agents Chemother.* **54**, 583–589 (2010).
- Fraser, C. M. et al. Complete genome sequence of *Treponema pallidum*, the syphilis spirochete. *Science* **281**, 375–388 (1998).
- Quétel, C. *History of Syphilis* (Johns Hopkins Univ. Press, 1990).
- Fernandez de Oviedo y Valdes, G. *Sumario de la natural historia de las Indias* (Fondo de Cultura Económico, 1526).
- Harper, K. N., Zuckerman, M. K., Harper, M. L., Kingston, J. D. & Armelagos, G. J. The origin and antiquity of syphilis revisited: an appraisal of Old World pre-Columbian evidence for treponemal infection. *Am. J. Phys. Anthropol.* **146**(Suppl 53), 99–133 (2011).
- Šmajš, D., Norris, S. J. & Weinstock, G. M. Genetic diversity in *Treponema pallidum*: implications for pathogenesis, evolution and molecular diagnostics of syphilis and yaws. *Infect. Genet. Evol.* **12**, 191–202 (2012).
- Giacani, L. et al. Complete genome sequence of the *Treponema pallidum* subsp. *pallidum* Sea81-4 strain. *Genome Announc.* **2**, e00333-14 (2014).
- Štaudová, B. et al. Whole genome sequence of the *Treponema pallidum* subsp. *endemicum* strain Bosnia A: the genome is related to yaws treponemes but contains few loci similar to syphilis treponemes. *PLoS Negl. Trop. Dis.* **8**, e3261 (2014).
- Marra, C. M. et al. Enhanced molecular typing of *Treponema pallidum*: geographical distribution of strain types and association with neurosyphilis. *J. Infect. Dis.* **202**, 1380–1388 (2010).
- Grillová, L. et al. Molecular typing of *Treponema pallidum* in the Czech Republic during 2011 to 2013: increased prevalence of identified genotypes and of isolates with macrolide resistance. *J. Clin. Microbiol.* **52**, 3693–3700 (2014).

13. Hodges, E. *et al.* Hybrid selection of discrete genomic intervals on custom-designed microarrays for massively parallel sequencing. *Nat. Protoc.* **4**, 960–974 (2009).
14. Meyer, M. & Kircher, M. Illumina sequencing library preparation for highly multiplexed target capture and sequencing. *Cold Spring Harb. Protoc.* <http://dx.doi.org/10.1101/pdb.prot5448> (2010).
15. Pětrošová, H. *et al.* Resequencing of *Treponema pallidum* ssp. *pallidum* strains Nichols and SS14: correction of sequencing errors resulted in increased separation of syphilis *treponeme* subclusters. *PLoS ONE* **8**, e74319 (2013).
16. Centurion-Lara, A. *et al.* Fine analysis of genetic diversity of the *tpr* gene family among *treponemal* species, subspecies and strains. *PLoS Negl. Trop. Dis.* **7**, e2222 (2013).
17. Mikalova, L. *et al.* Genome analysis of *Treponema pallidum* subsp. *pallidum* and subsp. *pertenue* strains: most of the genetic differences are localized in six regions. *PLoS ONE* **5**, e15713 (2010).
18. Achtman, M. Evolution, population structure, and phylogeography of genetically monomorphic bacterial pathogens. *Annu. Rev. Microbiol.* **62**, 53–70 (2008).
19. Pětrošová, H. *et al.* Whole genome sequence of *Treponema pallidum* ssp. *pallidum*, strain Mexico A, suggests recombination between yaws and syphilis strains. *PLoS Negl. Trop. Dis.* **6**, e1832 (2012).
20. Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Comput. Biol.* **10**, e1003537 (2014).
21. Lukehart, S. A. & Giacani, L. When is syphilis not syphilis? Or is it? *Sex. Transm. Dis.* **41**, 554–555 (2014).
22. Stadler, T., Kuhner, D., Bonhoeffer, S. & Drummond, A. J. Birth–death skyline plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc. Natl Acad. Sci. USA* **110**, 228–233 (2013).
23. Holt, K. E. *et al.* *Shigella sonnei* genome sequencing and phylogenetic analysis indicate recent global dissemination from Europe. *Nat. Genet.* **44**, 1056–1059 (2012).
24. Mutreja, A. *et al.* Evidence for several waves of global transmission in the seventh cholera pandemic. *Nature* **477**, 462–465 (2011).
25. Giacani, L. *et al.* Footprint of positive selection in *Treponema pallidum* subsp. *pallidum* genome sequences suggests adaptive microevolution of the syphilis pathogen. *PLoS Negl. Trop. Dis.* **6**, e1698 (2012).
26. Strouhal, M. *et al.* Genome differences between *Treponema pallidum* subsp. *pallidum* strain Nichols and *T. paraluisaniculiculi* strain Cuniculi A. *Infect. Immun.* **75**, 5859–5866 (2007).
27. Marra, C. M. *et al.* Antibiotic selection may contribute to increases in macrolide-resistant *Treponema pallidum*. *J. Infect. Dis.* **194**, 1771–1773 (2006).
28. Mitchell, S. J. *et al.* Azithromycin-resistant syphilis infection: San Francisco, California, 2000–2004. *Clin. Infect. Dis.* **42**, 337–345 (2006).
29. Morshed, M. & Jones, H. *Treponema pallidum* macrolide resistance in BC. *Can. Med. Assoc. J.* **174**, 349 (2006).
30. Matejkova, P. *et al.* Macrolide treatment failure in a case of secondary syphilis: a novel A2059G mutation in the 23S rRNA gene of *Treponema pallidum* subsp. *pallidum*. *J. Med. Microbiol.* **58**, 832–836 (2009).
31. Stamm, L. V. & Bergen, H. L. A point mutation associated with bacterial macrolide resistance is present in both 23S rRNA genes of an erythromycin-resistant *Treponema pallidum* clinical isolate. *Antimicrob. Agents Chemother.* **44**, 806–807 (2000).
32. Šmajs, D., Paštěková, L. & Grillová, L. Macrolide resistance in the syphilis spirochete, *Treponema pallidum* ssp. *pallidum*: can we also expect macrolide-resistant yaws strains? *Am. J. Trop. Med. Hyg.* **93**, 678–683 (2015).
33. Geisler, W. M. *et al.* Azithromycin versus doxycycline for urogenital *Chlamydia trachomatis* infection. *N. Engl. J. Med.* **373**, 2512–2521 (2015).
34. Centurion-Lara, A. in *Pathogenic Treponema: Molecular and Cellular Biology* (eds Radolf, J. D. & Lukehart, S. A.) 267–283 (Caister Academic, 2006).
35. Knell, R. J. Syphilis in renaissance Europe: rapid evolution of an introduced sexually transmitted disease? *Proc. Biol. Sci.* **271**(Suppl 4), S174–S176 (2004).
36. Kircher, M., Sawyer, S. & Meyer, M. Double indexing overcomes inaccuracies in multiplex sequencing on the Illumina platform. *Nucleic Acids Res.* **40**, e3 (2012).
37. Peltzer, A. *et al.* EAGER: Efficient ancient genome reconstruction. *Genome Biol.* **17**, 60 (2016).
38. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Preprint at <https://arxiv.org/abs/1303.3997> (2013).
39. McKenna, A. *et al.* The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
40. Simpson, J. T. & Durbin, R. Efficient *de novo* assembly of large genomes using compressed data structures. *Genome Res.* **22**, 549–556 (2012).
41. Schmidt, H. A., Strimmer, K., Vingron, M. & von Haeseler, A. TREE-PUZZLE: maximum likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics* **18**, 502–504 (2002).
42. Strimmer, K. & von Haeseler, A. Quartet puzzling: a quartet maximum-likelihood method for reconstructing tree topologies. *Mol. Biol. Evol.* **13**, 964–969 (1996).
43. Strimmer, K. & Rambaut, A. Inferring confidence sets of possibly misspecified gene trees. *Proc. R. Soc. B* **269**, 137–142 (2002).
44. Shimodaira, H. & Hasegawa, M. Multiple comparisons of log-likelihoods with applications to phylogenetic inference. *Mol. Biol. Evol.* **16**, 1114–1116 (1999).
45. Tamura, K., Stecher, G., Peterson, D., Filipski, A. & Kumar, S. MEGA6: molecular evolutionary genetics analysis version 6.0. *Mol. Biol. Evol.* **30**, 2725–2729 (2013).
46. Librado, P. & Rozas, J. DnaSP v5: A software for comprehensive analysis of DNA polymorphism data. *Bioinformatics* **25**, 1451–1452 (2009).
47. Baele, G. *et al.* Improving the accuracy of demographic and molecular clock model comparison while accommodating phylogenetic uncertainty. *Mol. Biol. Evol.* **29**, 2157–2167 (2012).
48. Rambaut, A., Suchard, M., Xie, D. & Drummond, A. *Tracer v1.6*. (2014); <http://beast.bio.ed.ac.uk/Tracer>
49. Rambaut, A. *FigTree v.1.4.2*. (2014); <http://tree.bio.ed.ac.uk/software/figtree/>
50. Bandelt, H.-J., Forster, P. & Röhl, A. Median-joining networks for inferring intraspecific phylogenies. *Mol. Biol. Evol.* **16**, 37–48 (1999).
51. Dai, T. *et al.* Molecular typing of *Treponema pallidum*: a 5-year surveillance in Shanghai, China. *J. Clin. Microbiol.* **50**, 3674–3677 (2012).
52. Flasarová, M. *et al.* Sequencing-based molecular typing of *Treponema pallidum* strains in the Czech Republic: all identified genotypes are related to the sequence of the SS14 strain. *Acta Derm. Venereol.* **92**, 669–674 (2012).
53. Grange, P. A. *et al.* Molecular subtyping of *Treponema pallidum* in Paris, France. *Sex. Transm. Dis.* **40**, 641–644 (2013).
54. Grimes, M. *et al.* Two mutations associated with macrolide resistance in *Treponema pallidum*: increasing prevalence and correlation with molecular strain type in Seattle, Washington. *Sex. Transm. Dis.* **39**, 954–958 (2012).
55. Peng, R.-R. *et al.* Molecular typing of *Treponema pallidum* causing early syphilis in China: a cross-sectional study. *Sex. Transm. Dis.* **39**, 42–45 (2012).
56. Tian, H. *et al.* Molecular typing of *Treponema pallidum*: identification of a new sequence of *tp0548* gene in Shandong, China. *Sex. Transm. Dis.* **41**, 551 (2014).
57. Tipple, C., McClure, M. O. & Taylor, G. P. High prevalence of macrolide resistant *Treponema pallidum* strains in a London centre. *Sex. Transm. Infect.* **87**, 486–488 (2011).
58. Wu, B.-R. *et al.* Multicentre surveillance of prevalence of the 23S rRNA A2058G and A2059G point mutations and molecular subtypes of *Treponema pallidum* in Taiwan, 2009–2013. *Clin. Microbiol. Infect.* **20**, 802–807 (2014).

## Acknowledgements

Research in Zurich by N.A. and H.C.B. was funded by the Forschungskredit and the University of Zurich. A.H. was funded by an ERC Starting Grant. F.G.C. and L.S.B. were funded by MINECO (Spanish Government) and PROMETEO (Generalitat Valenciana). K.I.B. was funded by the Social Sciences and Humanities Research Council of Canada. L.M. was funded by the Faculty of Medicine of Masaryk University. The authors thank S. Lautenschlager for guidance, A. Drummond for input on BEAST, S. Lukehart for providing HaitiB, Sea86-1, Bal3, Bal9, Bal73-1 and Grady1 strain DNA, and C. Marra for providing UW249B and UW231B strain DNA. The authors also thank A. Messina and the S3IT at the University of Zurich for providing computational resources and services, and I. Schoechli and L. Keller's group for their valued support.

## Author contributions

N.A. and H.C.B. conceived the investigation. N.A., L.G., S.J.N., D.S., P.P.B., F.G.-C., K.N., J.K. and H.C.B. devised research and analyses. N.A., G.J., A.P., A.S., A.H., M.S., L.G., L.S.-B., D.K., L.R.D., L.M., F.G.-C. and K.N. analysed data. N.A., V.J.S., M.S., L.G., K.I.B., L.R.D., L.G.V. and P.P.B. contributed to or performed experiments. M.S., L.G., S.B., P.K., P.F., P.R.G., M.A.P., L.G.V., M.R.F., A.M., D.S., P.P.B. and F.G.-C. provided clinical samples and A.C.L., L.G., S.J.N. and D.S. provided laboratory samples. N.A. and H.C.B. wrote the manuscript with significant contributions from M.S., L.G., L.S.-B., D.K., K.I.B., L.R.D., L.M., S.B., L.G., S.J.N., D.S., P.P.B., F.G.-C., K.N. and J.K. and with comments from all co-authors.

## Additional information

Supplementary information is available for this paper. Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). Correspondence and requests for materials should be addressed to N.A., F.G.C., K.N., J.K. and H.C.B.

## Competing interests

The authors declare no competing financial interests.