

1 **Selective whole genome amplification as a tool to enrich specimens with low *Treponema pallidum* genomic DNA copies for whole genome**  
2 **sequencing**

3  
4 **Charles M. Thurlow,<sup>a#</sup> Sandeep J. Joseph,<sup>a</sup> Lilia Ganova-Raeva,<sup>b</sup> Samantha S. Katz,<sup>a</sup> Lara Pereira,<sup>a</sup> Cheng Chen,<sup>a</sup> Alyssa Debra,<sup>a</sup> Kendra**  
5 **Vilfort,<sup>a</sup> Kimberly Workowski,<sup>a,c</sup> Stephanie E. Cohen,<sup>d</sup> Hilary Reno,<sup>e,f</sup> Yongcheng Sun,<sup>a</sup> Mark Burroughs,<sup>g</sup> Mili Sheth,<sup>g</sup> Kai-Hua Chi,<sup>a</sup>**  
6 **Damien Danavall,<sup>a</sup> Susan S. Philip,<sup>c</sup> Weiping Cao,<sup>a</sup> Ellen N. Kersh,<sup>a</sup> and Allan Pillay<sup>a#</sup>**

7  
8 <sup>a</sup>Division of STD Prevention, Centers for Disease Control and Prevention, Atlanta, Georgia, USA

9 <sup>b</sup>Division of Viral Hepatitis, Centers for Disease Control and Prevention, Atlanta, Georgia, USA

10 <sup>c</sup>Department of Medicine, Emory University, Atlanta, Georgia, USA

11 <sup>d</sup>San Francisco Department of Public Health, San Francisco, California, USA

12 <sup>e</sup>St. Louis County Sexual Health Clinic, St. Louis, Missouri, USA

13 <sup>f</sup>Division of Infectious Diseases, Washington University, St. Louis, Missouri, USA

14 <sup>g</sup>Division of Scientific Resources, Centers for Disease Control and Prevention, Atlanta, Georgia, USA

15

16 Running title: Sequencing of *T. pallidum* from Clinical Specimens

17

18 #Address correspondence to Dr. Charles M. Thurlow, [cthurlow@cdc.gov](mailto:cthurlow@cdc.gov) and Dr. Allan Pillay, [apillay@cdc.gov](mailto:apillay@cdc.gov).

19

## 20 **Abstract.**

21 Downstream next generation sequencing of the syphilis spirochete *Treponema pallidum* subspecies *pallidum* (*T. pallidum*) is hindered by  
22 low bacterial loads and the overwhelming presence of background metagenomic DNA in clinical specimens. In this study, we investigated  
23 selective whole genome amplification (SWGA) utilizing Multiple Displacement Amplification (MDA) in conjunction with custom  
24 oligonucleotides with an increased specificity for the *T. pallidum* genome, and the capture and removal of CpG-methylated host DNA followed by  
25 MDA as enrichment methods to improve the yields of *T. pallidum* DNA in rabbit propagated isolates and lesion specimens from patients with  
26 primary and secondary syphilis. Sequencing was performed using the Illumina MiSeq v2 500 cycle or NovaSeq 6000 SP platform. These two  
27 enrichment methods led to 93-98% genome coverage at 5 reads/site in 5 clinical specimens from the United States and rabbit propagated isolates,  
28 containing >14 *T. pallidum* genomic copies/μl input for SWGA and >129 genomic copies/μl for CpG methylation capture with MDA. Variant  
29 analysis using sequencing data derived from SWGA-enriched specimens, showed that all 5 clinical strains had the A2058G mutation associated  
30 with azithromycin resistance. SWGA is a robust method that allows direct whole genome sequencing (WGS) of specimens containing very low  
31 numbers of *T. pallidum*, which have been challenging until now.

## 32 **Importance**

33 Syphilis is a sexually transmitted, disseminated acute and chronic infection caused by the bacterial pathogen *Treponema pallidum*  
34 subspecies *pallidum*. Primary syphilis typically presents as single or multiple mucocutaneous lesions, and if left untreated, can progress through  
35 multiple stages with varied clinical manifestations. Molecular studies rely on direct amplification of DNA sequences from clinical specimens;  
36 however, this can be impacted by inadequate samples due to disease progression or timing of patients seeking clinical care. While genotyping has  
37 provided important data on circulating strains over the past two decades, whole genome sequencing data is needed to better understand strain  
38 diversity, perform evolutionary tracing, and monitor antimicrobial resistance markers. The significance of our research is the development of a  
39 SWGA DNA enrichment method that expands the range of clinical specimens that can be directly sequenced to include samples with low numbers  
40 of *T. pallidum*.

## 41 **Introduction**

42 Syphilis rates have been steadily increasing in the United States with 38,992 cases (11.9 per 100,000 people) of primary and secondary  
43 syphilis and 1,870 cases (48.5 per 100,000 live births) of congenital syphilis reported to the CDC during 2019 (1). This represents a 167.2%  
44 increase in primary and secondary syphilis rates since 2010 and a 291.1% increase in congenital syphilis reported since 2015. While syphilis rates  
45 have been on the rise in the U.S., the genetic diversity of the bacterial pathogen *Treponema pallidum* subspecies *pallidum* (hereafter referred to as  
46 *T. pallidum*), in this setting, is not well understood due to the lack of recently sequenced whole genomes from clinical specimens. Strain diversity  
47 has been gleaned from molecular epidemiology studies, which are based on 3 to 4 genetic loci, but may not be representative of the entire *T.*  
48 *pallidum* genome (2-5).

49 Molecular studies have relied primarily on *T. pallidum* strains propagated in rabbits or DNA amplified directly from clinical specimens,  
50 because *T. pallidum* cannot be grown on routine laboratory media. However, advances have been made with *in vitro* tissue culture and the

51 propagation of *T. pallidum* in rabbits from cryopreserved genital lesion specimens, which may make routine culture directly from clinical  
52 specimens a possibility in the near future (6-7). Despite these advances, the methods are still time-consuming and impractical for laboratory  
53 diagnosis and molecular epidemiological studies of syphilis.

54 Metagenomic shotgun sequencing approaches have made significant advances in recent years with sequence data being used for pathogen  
55 detection, *in silico* or whole genome typing, and antimicrobial resistance marker detection, in addition to phylogenetic analyses (8-10). However,  
56 direct whole genome sequencing (WGS) of *T. pallidum* from clinical specimens and rabbit isolates can be problematic due to bacterial genomic  
57 DNA being outweighed by either human or rabbit DNA. Several DNA enrichment methods have been described for *T. pallidum* including RNA  
58 bait capture techniques, methyl-directed enrichment using the restriction nuclease DpnI, and pooled whole genome amplification, which have  
59 generated *T. pallidum* specific WGS data from over 700 metagenomic samples; however, specimens with low numbers of *T. pallidum* remains  
60 challenging (11-16). Therefore, additional approaches that would enable sequencing of samples with low bacterial loads are needed.

61 Azithromycin has been used as an alternative to penicillin for treating early syphilis in the US; however, macrolide-resistant *T. pallidum*  
62 strains, associated with two mutations (A2508G, A2509G) in the 23S rRNA genes, have been reported in many states (17-18). While macrolides  
63 are no longer recommended for treatment of syphilis in the US, periodic monitoring is useful to determine the prevalence of resistant strains (19).

64 In this study, we describe a robust DNA enrichment method based on selective whole genome amplification (SWGA) using multiple  
65 displacement amplification (MDA) and custom primers that enables WGS of clinical specimens with very low genomic copies of *T. pallidum* and  
66 use of the sequence data for macrolide mutation analysis. We also investigated an alternative method that uses CpG methylated capture of host  
67 DNA followed by MDA with random oligonucleotide primers.

## 68 Results

69 **Real-time qPCR on clinical specimens and spiked samples.** The *T. pallidum* PCR results for all clinical specimens are shown in Table 1. Out of  
70 the 11 Emory specimens processed using the standard extraction protocol, only one specimen exceeded 100 genomic copies/ $\mu$ l based on *polA*  
71 qPCR (Table 1). The remaining 10 specimens had an average copy number  $<1$  copy/ $\mu$ l of DNA extract. These 11 specimens had an average  
72 standardized RNP cycle threshold (RNP<sub>Ct</sub>) value of  $30.71 \pm 0.13$ , and the lowest C<sub>t</sub> value (highest concentration of RNP) was 25.22. Based on this  
73 data, an RNP<sub>Ct</sub> value of 25.22 was targeted as the cut-off for the spiked samples below.

74 **NEBNext microbiome enrichment with MDA.** The serially diluted spiked samples enriched with the NEB Microbiome Enrichment Kit with  
75 subsequent REPLIg Single Cell MDA (hereafter referred to as NEB+MDA) showed a marked increase in *polA* copy number by qPCR (Fig. 1,  
76 Table S1). The non-diluted samples indicated an average *polA* copy number of  $6.67 \times 10^6 \pm 2.74 \times 10^5$  per  $\mu$ l of enriched DNA, which was 603.02  
77 times greater than the input copy number. The 10-fold diluted samples also indicated increases in *polA* copy numbers, with an average of  $7.85 \times$   
78  $10^5 \pm 3.79 \times 10^4$ ,  $1.28 \times 10^5 \pm 1.27 \times 10^4$ ,  $8.66 \times 10^3 \pm 2.54 \times 10^3$ , and  $964 \pm 574.23$  copies/ $\mu$ l from 1:10 -1:10,000 dilution, respectively (Table S1).  
79 This was a 482 – 995.09 times enrichment when compared to the input copy number. Upon comparing the average RNP<sub>Ct</sub> of each dilution in the  
80 series, the enriched samples indicated  $29.28 \pm 1.07$ ,  $31.15 \pm 0.46$ ,  $30.25 \pm 0.56$ ,  $31.08 \pm 0.59$ ,  $31.42 \pm 0.45$  for the neat – 1:10,000 dilution,  
81 respectively (Table S1). The RNP<sub>Ct</sub> value of each enriched sample in the dilution series were insignificantly different from one another, with an  
82 average RNP<sub>Ct</sub> =  $30.64 \pm 0.33$  for all dilutions in the series (P = 0.22).

83 After enriching with NEB+MDA, the average DNA percent for the neat to 1:10,000 dilutions indicated a range of  $2.33\% \pm 0.10$ -  $3.91 \times$   
84  $10^{-4}\% \pm 2.50 \times 10^{-4}\%$  of the total DNA belonging to *T. pallidum*, respectively (Fig. 2). Further, this form of enrichment generated up to a 26.12-

85 fold increase in the percent of *T. pallidum* DNA, and an average of 16.27-fold  $\pm$  1.92-fold increase, amongst all enriched replicates when  
86 compared to the unenriched input. All samples enriched by NEB+MDA were significantly different in their percent *T. pallidum* DNA when  
87 compared to their respective inputs ( $P < 0.01$ ). Apart from enriched samples from the 1:100 and 1:1,000 diluted *polA* inputs, we observed that by  
88 increasing the *polA* input copy number 10-fold resulted in a significant increase in the total DNA belonging to *T. pallidum* post-enrichment ( $P =$   
89 0.06 and  $P < 0.05$ , respectively).

90 Genome sequencing data derived from samples enriched by NEB+MDA showed 0.01 – 10.52% of the quality-controlled reads binned as  
91 *T. pallidum*, along with a mean mapping read depth to *T. pallidum* Nichols reference genome (*NC\_021490.2*) ranging from 0.05 – 501.75. An  
92 average percent coverage of 99.99%, 99.99%, and 97.29% across the Nichols reference genome with at least 5 reads mapped per site (5X) for the  
93 neat, 1:10, and 1:100 diluted samples, respectively, was observed among the NEB+MDA enriched samples (Fig.3A; Table S1 and Fig. S1). The  
94 coverage estimates indicated low deviations from this average in all replicates, with  $2.92 \times 10^{-4} \% - 1.38\%$  standard error between all replicates for  
95 the neat – 1:100 diluted samples. At the same time, for a higher coverage of at least 10 reads mapped per nucleotide (10X), the 1:100 diluted  
96 samples had an average percentage coverage of 84.14% while neat and 1:10 dilution samples were covered at 99.99% and 99.99% across the  
97 reference genome, respectively. A sharp decline in coverage was observed in the 1:1,000 diluted samples, with a break down in replication at an  
98 average coverage of  $27.08\% \pm 18.62\%$  for the 1:1,000 dilution and  $4.80\% \pm 0.63\%$  for the 1:10,000 diluted samples at 5X read depth. With the  
99 QC criteria for efficiency set at  $\geq 90\%$  at  $\geq 5X$  read depth, samples sequenced post NEB+MDA enrichment had a limit of detection (LoD) of 129  
100 *polA* copies/ $\mu$ l of extract (Fig. 3A; Table S1 and Fig. S1).

101 Post NEB+MDA enrichment of isolate CDC-SF003, we observed  $2.39 \times 10^6 \pm 1.35 \times 10^5$  *polA* copies/ $\mu$ l of DNA extract. Further, 1.06%  
102 of the total DNA belonged to *T. pallidum* post enrichment and 3.29% of the host removed quality-controlled sequencing reads were classified as *T.*

103 *pallidum*. Sequencing indicated a 98.60% coverage across the *T. pallidum* SS14 reference genome (*NC\_021508.1*) at 5X read depth with a mean  
104 mapping depth of 46.43 (Fig. 4; Table 2 and Fig. S2).

105 **SWGA Enrichment of *T. pallidum* Nichols.** A total of 12 primer sets were tested by SWGA using Equiphi29 MDA (Table S2-S3). The 1:100  
106 diluted Nichols DNA sample (~129 copies/ $\mu$ l) was used to evaluate each of the 12 primers since it was comparable to specimen EUHM-004,  
107 which had 106.7 *polA* copy/ $\mu$ l (Table 1; Table S4). Each of the primer sets indicated a 6.86 – 1.16 x 10<sup>5</sup> times enrichment when compared to the  
108 input Nichols copy number (Fig. 5A). Further, we observed a >10,000-fold increase in *polA* copy number in samples enriched with 7 of the 12  
109 primer sets (SWGA Pal 2, 4, 5, 9, 10, 11, and 12). SWGA Pal 9 and Pal 11 gave the highest enrichment at 1.13 x 10<sup>5</sup>, and 1.16 x 10<sup>5</sup> times,  
110 respectively (Table S4). The difference observed between Pal 9 and Pal 11 in the *T. pallidum polA* copy number and relative percent DNA  
111 belonging to *T. pallidum* was insignificant; however, Pal 11 was selected for testing the SWGA limit of detection ( $P > 0.1$ ; Fig. 5 and Table S4).

112 To determine the SWGA Pal 11 primer set's LoD and enrichment for *T. pallidum*, SWGA was performed in triplicate on the 10-fold  
113 dilution series. The ~1.11x10<sup>4</sup> copies/ $\mu$ l (neat) sample was eliminated from the dilution series, as this was ~100-fold increase in *T. pallidum* copy  
114 number when compared to the clinical specimens tested. We observed a marked increase in *polA* copy number in every dilution in the series post  
115 enrichment (Fig. 1; Table S1). The *polA* copy number ranged from 1.11 x 10<sup>6</sup>  $\pm$  6.68 x 10<sup>5</sup> for the 1:10,000 dilution to 2.04 x 10<sup>7</sup>  $\pm$  1.20 x 10<sup>7</sup> in  
116 the 1:10 dilution (Table S1). When compared to the input *polA* copy number, this was a 2.01 x 10<sup>4</sup>-fold, 1.19 x 10<sup>5</sup>-fold, 3.53 x 10<sup>5</sup>-fold, and 5.53  
117 x 10<sup>5</sup>-fold increase in the enriched samples, from 1:10 -1:10,000 dilution, respectively. Upon comparing the average RNP<sub>Ct</sub> of each dilution in the  
118 series, the SWGA enriched samples indicated a 29.36  $\pm$  0.37 - 28.65  $\pm$  0.16 for the 1:10 -1:10,000 dilution, respectively (Table S1). The average  
119 RNP<sub>Ct</sub> at each 10-fold increase in *polA* concentration were insignificantly different from one another ( $P > 0.1$ ); however, by increasing the *polA*  
120 input 100-fold, we observed a significant decrease in RNP concentration ( $P < 0.03$ ).

121 After enriching with SWGA, we observed that dilutions ranging from 1:10 to 1:10,000 held  $27.93\% \pm 1.57\%$  -  $3.29\% \pm 1.93\%$  of the total  
122 DNA belonging to *T. pallidum*, respectively (Fig. 2). This reflected up to a  $1.63 \times 10^5$ -fold increase in the relative *T. pallidum* and an average of  
123  $2.43 \times 10^4$ -fold  $\pm 1.05 \times 10^4$ -fold increase amongst all replicate SWGA enriched samples when compared to the unenriched samples. All samples  
124 were significantly increased in their relative *T. pallidum* DNA when compared to their respective inputs ( $P < 0.0001$ ). While there was observed  
125 deviations in the percent DNA between replicates, the 1:10,000 diluted replicates still yielded a 28.40-fold  $\pm 17.71$ -fold increase in DNA  
126 belonging to *T. pallidum* post SWGA when compared to the non-enriched neat dilution ( $P < 0.0001$ ).

127 Genome sequencing data derived from the SWGA enriched Nichols samples showed 0.98%-78.05% of the quality-controlled reads binned  
128 as *T. pallidum*, along with a mean mapping read depth to *T. pallidum* Nichols reference genome ranging from  $65.82 - 4.89 \times 10^3$ . An average  
129 percent coverage of  $98.67\% \pm 0.005\%$ ,  $98.62\% \pm 0.003\%$ , and  $96.15\% \pm 0.082\%$  across the Nichols genome at 5X read depth was observed  
130 among the SWGA enriched 10-fold dilution series samples for the 1:10, 1:100 and 1:1,000 diluted samples, respectively (Fig. 3B; Table S1 and  
131 Fig. S3). Further, coverage indicated low deviations from this average in all replicates, with a 0.0002% - 1.72% standard error between all  
132 replicates for the 1:10 - 1:1,000 diluted samples. We did observe a sharp decline in coverage from the 1:1,000 to 1:10,000 dilution with an average  
133 coverage of  $38.46\% \pm 2.50\%$  for the 1:10,000 diluted replicates a 5X read depth (Fig. 3B; Table S1 and Fig. S3).

134 Upon comparing the percent *T. pallidum* DNA derived from both enrichment methods, we observed that SWGA consistently produced  
135 higher relative *T. pallidum* DNA in all samples (Fig. 2). We observed that the 10-fold dilutions enriched with SWGA exhibited an average of  
136 94.08-fold -  $1.41 \times 10^4$ -fold increase in relative *T. pallidum* DNA in the 1:10-1:10,000 diluted samples when compared to the dilutions enriched by  
137 NEB+MDA. All dilutions of each enrichment were significantly different from one another ( $P < 0.01$ ), apart from the 1:10,000 and 1:1,000 diluted  
138 samples enriched by SWGA and the neat diluted samples enriched by NEB+MDA ( $P > 0.07$ ).



139 Comparing the sequencing data derived from the 1:10 and 1:100 diluted Nichols samples enriched using the NEB+MDA and SWGA, all  
140 samples exhibited >95% coverage at 5X read depth (Fig. 3; Table S1 and Fig. S1, S3). There was a decline in coverage observed in the 1:1,000  
141 diluted samples enriched by NEB+MDA, with an average coverage of  $27.08\% \pm 24.80\%$  at 5X read depth. This drop was not observed in the  
142 1:1,000 diluted samples enriched by SWGA, which still held >95% coverage at 5X read depth. The 1:10,000 diluted samples enriched NEB+MDA  
143 and SWGA exhibited <95% coverage at 5X read depth.

#### 144 **Enrichment of Clinical Strains.**

145 Due to the increased sequencing coverage derived from the SWGA enriched Nichols strain, SWGA was chosen for enriching clinical  
146 specimens with low numbers of *T. pallidum* (Fig. 3, Table S1). SWGA on clinical specimen EUHM-004 gave an average *polA* of  $6.37 \times 10^6 \pm 2.24$   
147  $\times 10^5$  copies/ $\mu$ l with 5.56% of the total DNA belonging to *T. pallidum* (Table 2). Next generation sequencing using the MiSeq v2 (500 cycle)  
148 platform revealed 95.13% coverage across the *T. pallidum* genome at 5X read depth (Fig. 4; Table 2 and Fig. S2). After large-scale DNA  
149 extraction, we observed  $31.5 \pm 0.5$ ,  $122 \pm 1.15$ , and  $103 \pm 6.55$  *polA* copies/ $\mu$ l for specimens EUHM-012 – EUHM-014, respectively (Table 1).  
150 For specimen EUHM-012, we observed an average *polA* of  $2.14 \times 10^6 \pm 2.82 \times 10^4$  copies/ $\mu$ l with 1.72% of the total DNA belonging to *T.*  
151 *pallidum* post-enrichment by SWGA (Table 2). Sequencing indicated a 93.98% coverage across the *T. pallidum* genome at 5X read depth (Fig. 4;  
152 Table 2 and Fig. S2).

153 When compared to EUHM-012, EUHM-013 had a higher *polA* copy number at  $5.16 \times 10^6 \pm 2.20 \times 10^5$  copies/ $\mu$ l with 15.48% of the total  
154 DNA belonging to *T. pallidum* (Table 2). The sequencing data correlated with the qPCR data, indicating a 98.56% coverage across the *T. pallidum*  
155 genome at 5X read depth (Fig. 4; Table 2 and Fig. S2). We also observed EUHM-014 held an increased *polA* copy number post-SWGA, with  $2.57$   
156  $\times 10^6 \pm 2.21 \times 10^5$  copies/ $\mu$ l and 4.72% of the total DNA belonging to *T. pallidum* (Table 2). Upon sequencing, we observed 98.49% coverage

157 across the *T. pallidum* genome at 5X read depth (Fig. 4; Table 2 and Fig. S2). The *polA* copy number for specimen STLC-001 was  $7.42 \times 10^6 \pm$   
158  $7.20 \times 10^5$  copies/ $\mu$ l with 8.34% of the total DNA belonging to *T. pallidum* (Table 2). The sequencing coverage was 95.94% at 5X read depth  
159 where 38.91% of the quality-controlled reads binned as *T. pallidum* along with a mean depth read coverage of 1,133.43X (Fig. 4; Table 2 and Fig.  
160 S2).

### 161 **Phylogenetic Analysis and Characterization of Genotypic Macrolide Resistance.**

162 To analyze whether genomes generated from the 7 clinical specimens or isolates clustered to any of the two deep-branching monophyletic  
163 *T. pallidum* lineages, Nichols-like and Street-14(SS14)-like, a whole genome phylogenetic tree was constructed using the genomes derived from  
164 the clinical specimens/isolates along with 126 high quality published *T. pallidum* genome sequences as of May 2021 (12-15, 20-22; see Table S5,  
165 methods in supplemental materials). Phylogenetic analysis revealed the presence of two dominant lineages, of which most strains belonged to the  
166 SS14-like lineage. We identified a total of four monophyletic clades within this phylogenetic tree with  $\geq 30$  bootstrap support (Fig. 6). Three of the  
167 clinically derived genomes from Atlanta, EUHM-004 (2019) EUHM-012 (2019), and EUHM-014 (2020), belonged to Nichols-like lineage (clade  
168 1; n=12; Fig. 6). Interestingly, the other nine Nichols-like genomes in clade 1 were recent clinically derived genomes from Cuba (n=2; 2015-  
169 2016), Australia (n=1; 2014), France (n=2; 2012-2013) and UK (n=3; 2016), and were distinct from the original Nichols strain isolated in 1912  
170 and sent to different North American labs as *in vivo* derived clones, suggesting that we might not yet fully understand the current diversity of this  
171 lineage. The three clinical specimens from Atlanta (EUHM-004, EUHM-012, and EUHM-014) and three clinically derived genomes from UK  
172 isolated in 2016 (NL14, NL19 and NL17) carried the 23S rRNA A2058G mutation that confers macrolide resistance, suggesting a recent  
173 acquisition of this antibiotic resistance variant in the Nichols-like lineage.

174 Even though previous phylogenomic analyses indicated that SS14-lineage showed a polyphyletic structure, our phylogenetic analysis with  
175 a greater number of genomes showed the presence of 3 monophyletic clades (Clades 2, 3 and 4)(12, 14; Fig. 6). Clades 2 and 4 contained genomes  
176 clustered within the previously reported SS14 $\Omega$ -A sub cluster, which also contained two clades corresponding to the clades 2 and 4 detected in this  
177 study, and contained genomes derived from Europe and North America; while clade 3 was similar to sub cluster SS14 $\Omega$ -B and composed of  
178 Chinese and North American derived *T. pallidum* genomes. The rabbit-derived clinical isolate, CDC-SF003 (San Francisco, U.S; 2017) sequenced  
179 in this study, clustered within clade 2; while EUHM-013 (Atlanta, U.S; 2020) and STLC-001 (St. Louis, U.S; 2020) genomes clustered within  
180 clade 4. Sequence analysis showed that all 3 strains carried the A2058G AMR variant for macrolide resistance. Macrolide resistance strains were  
181 widespread among the SS14-lineage with higher proportion among the genomes in clades 2 and 3 compared to clade 4 genomes. The A2058G  
182 point mutation identified in 4 patient specimens and isolate CDC-SF003 was verified by real-time PCR testing of genomic DNA and SWGA-  
183 enriched samples (data not shown). There was inadequate sample for the fifth specimen to confirm the mutation by real-time PCR testing.

184 All the Nichols-like genomes derived from the NEB+MDA and SWGA 10-fold dilution series that contained *T. pallidum* reads mapped to  
185  $\geq 90\%$  of the genome with at least 5X read depth formed a tight monophyletic clade (bootstrap support of 88/100) and clustered with the lab-  
186 derived Nichols-Houston-J genome (bootstrap support of 100/100), indicating that genomes generated from both methods are adequate to capture  
187 genetic variants required to perform a high resolution phylogenetic analysis (Fig. S4).

## 188 Discussion

189 WGS of *T. pallidum* is often challenging due to low bacterial loads or the difficulty of obtaining adequate samples for testing. In this  
190 study, we sought to develop a method for performing WGS from rabbit propagated isolates and clinical specimens containing lower *T. pallidum*  
191 numbers, leading us to investigate CpG capture and SWGA.

192 CpG capture has been successfully used for enriching bacterial genomic DNA in metagenomic samples (23-24), but this method has not  
193 been used for *T. pallidum*. During our testing, we observed increases in *polA* copy numbers and relative *T. pallidum* percent DNA in the neat to  
194 1:1,000 dilutions enriched by NEB+MDA when compared to the non-enriched inputs. Further, the results of the percent *T. pallidum* observed in  
195 the enriched 1:10,000 diluted samples correlated with the decrease in overall coverage across the Nichols genome. Even though we observed an  
196 increase in both *polA* copy number and relative percent *T. pallidum* DNA for the enriched diluted 1:1,000 diluted samples, we still only gained  
197 ~50% genomic coverage. This could be due to the remnant human DNA that was not initially captured prior to MDA, or the loss of *T. pallidum*  
198 DNA during the enrichment. While there was no significant difference in the relative human RNP copy number from dilution to dilution, there is a  
199 minimum *T. pallidum* copy number input required to outweigh the remnant human DNA during the metagenomic shotgun sequencing. Taking the  
200 above into consideration, we observed that >129 *polA* copies/ $\mu$ l can generate >95% coverage at 5X read depth from the Nichols strain post  
201 NEB+MDA. The results observed post NEB+MDA enrichment of clinical isolate CDC-SF003 correlated with the Nichols limit of detection  
202 validation, with >98% coverage at 5X read depth across the *T. pallidum* genome. *In silico* variant analysis correlated with real-time PCR detection  
203 of the mutations associated with macrolide resistance in clinical isolate CDC-SF003. Further, phylogenetics revealed that this strain belonged to  
204 the SS14 lineage, which correlated with its enhanced CDC typing method (ECDCT) strain type, 4d9f, as previously reported (7). While this  
205 enrichment method yielded good results with isolates, most clinical specimens collected in this study had lower than 100 *polA* DNA copies/ $\mu$ l of  
206 *T. pallidum* leading us to consider an alternative method.

207 SWGA has been shown to be successful with other bacterial pathogens in metagenomic samples; however, it has not been investigated  
208 with *T. pallidum* (25-27). We observed that samples enriched by SWGA using multiple primer sets exhibited a 10,000-fold increase in *polA* copy  
209 number, with Pal 9 and 11 producing the highest relative percent *T. pallidum* DNA at 29% and 31%, respectively. While we chose to work with  
210 Pal 11 as the optimal set, Pal 9 could also be a good alternative for enriching syphilis specimens. Further testing using Pal 11 showed that the limit  
211 of detection was increased when compared to the *T. pallidum* enrichment obtained with NEB+MDA, with significant increases in both *polA* copy  
212 number and percent *T. pallidum* across the 10-fold dilution series. Coverage across the *T. pallidum* genome exceeded 95% at 5X read depth for all  
213 diluted samples, apart from the 1:10,000 diluted samples. Interestingly, we observed that increasing the input 100-fold resulted in a significant  
214 decrease in the presence of RNP post-enrichment. Our data shows that >14 *T. pallidum polA* copies/ $\mu$ l can generate at least 95% coverage at 5X  
215 read depth with the Nichols strain, which translated well to the clinical specimens tested. While there was a decrease in coverage in one of the  
216 clinical specimens at 94.44% with 5X read depth when compared to the 98.62% coverage at 5X read depth observed in the 1:100 diluted Nichols  
217 isolates, this could be primarily due to the improved capabilities of the NovaSeq 6000 when compared to the MiSeq v2 (500 cycle) platform used  
218 to sequence this clinical specimen. Another possible reason for the variation in coverage could be due to the lower *T. pallidum* input copy number  
219 in the clinical specimens.

220 The genomes derived directly from the 5 clinical specimens using SWGA were phylogenetically associated with the representative  
221 lineages (either Nichols-like or SS14-like) and also provided high levels of within lineage strain resolution, which is ideal for effective tracking of  
222 various strains circulating within a geographical area and outbreak investigations. In addition, the NGS methods described here can be used for  
223 macrolide resistance marker detection. As observed with NEB+MDA enrichment, *in silico* azithromycin mutation detection performed on the  
224 SWGA enriched specimens matched the results obtained with a real-time PCR, indicating that all clinical specimens contained the A2058G

225 mutation. SWGA-based enrichment also enabled sequencing of specimens within the range of detection limits for real-time PCR assays,  
226 suggesting that our NGS workflow can be adapted for *T. pallidum* detection in metagenomic samples.

227 In terms of expense, both methods are cost-effective for enriching *T. pallidum* genomic DNA, and while SWGA is cheaper than  
228 NEB+MDA, sequencing reagents are the true limiting factor for WGS. With the recent advancements in large-scale sequencing platforms, overall  
229 sequencing costs can be further reduced. While NovaSeq 6000 has a much higher potential for multiplex sequencing, our data shows compatibility  
230 of these enrichments for both NovaSeq 6000 and MiSeq platforms.

231 While we successfully enriched *T. pallidum* whole genomes in clinical specimens, the success of SWGA is limited by the constraint on  
232 primer size, which may reduce the selectivity for the target genome. Phi29 functions best between 30-35°C, and ramp-down incubations have been  
233 shown as an effective means of utilizing larger primers with increased melting temperatures (26-29). To help alleviate the constraints on primer  
234 size, we utilized a thermostable phi29 mutant which has a much higher optimal temperature at 45°C (30) compared to the 30-35°C functional  
235 range of the phi29 polymerase (26-27). This higher optimal temperature permits the use of longer oligonucleotides to be used in the SWGA  
236 reaction, potentially increasing the selectivity for the *T. pallidum* genome. The phi29 mutant has also shown to be more efficient, with a 3-hour  
237 exhaustion time when compared to the 8-16 hours required for the wild-type phi29 (30).

238 Our results show that SWGA is more sensitive, less cumbersome, and a faster method for enriching clinical specimens when compared to  
239 NEB+MDA, allowing for WGS of metagenomic samples with very low numbers of *T. pallidum*. In addition, the sequencing data generated is of  
240 sufficient quality to enable phylogenetic analyses and detection of mutations associated with azithromycin resistance. While the NEB+MDA was  
241 unsuitable for the clinical specimens in this study, our data suggests that it can be used for samples exceeding 129 genomic copies/ $\mu$ l.

242 **Materials and Methods.**

243 **Specimen collection, *T. pallidum* strains used for WGS, and real-time qPCR.** Specimens used in this study were collected from men  
244 presenting with lesions of primary or secondary syphilis to the Emory Infectious Diseases Clinic, Emory University Hospital Midtown (EUHM) in  
245 Atlanta, GA and St Louis County STD Clinic (STLC) in St. Louis, MO (Table 1). Patients were diagnosed with syphilis based on clinical  
246 presentation and serology testing. Fourteen swab specimens were collected in Aptima Multitest storage medium (Hologic, Inc., Marlborough, MA)  
247 at Emory Infectious Diseases Clinic and 1 specimen at St. Louis County STD Clinic (Table 1). All specimens were stored at -80°C until shipment  
248 on dry ice to the CDC. The *T. pallidum* Nichols reference strain was used for initial optimization and validation of the two enrichment methods. A  
249 recent rabbit propagated isolate, CDC-SF003, was also included for testing (Table 1; 7). Prior to study commencement, local IRB approvals were  
250 obtained from, Emory University, and St. Louis County Department of Public Health, and the project was approved at CDC

251 DNA was extracted from specimens and rabbit testis extracts using the QIAamp DNA Mini Kit (Qiagen, Germantown, MD) following the  
252 manufacturer's recommendations. Large-scale DNA extraction of three specimens was carried out on 1.5 ml of the Aptima stored specimen using  
253 the QIAamp DNA Mini Kit following the manufacturer's recommendations for upscaling with slight modifications (Table 1). Proteinase K was  
254 added at 0.1X total sample volume, and AL Buffer and absolute ethanol were added at 1X total sample volume. Each sample was processed  
255 through a single column, washed following the manufacturer's recommendations, and eluted in 100 µl AE Buffer (Qiagen). Following DNA  
256 extraction, each sample was tested by a real-time quantitative duplex PCR (qPCR) targeting the *po1A* gene of *T. pallidum* and human RNase P  
257 gene (RNP) using a Rotor-Gene 6000 instrument (Qiagen) as previously described with modifications (7; see additional methods in supplemental  
258 materials).

259 **Enrichment of *T. pallidum* by capture of CpG methylated host DNA and multiple displacement amplification (MDA).** Initially, DNA  
260 concentration of extracts from clinical specimens and rabbit propagated strains were measured using the Qubit dsDNA HS assay (Thermo Fisher  
261 Scientific, Waltham, MA). Capture and removal of CpG methylated host DNA from samples were carried out using the NEBNext Microbiome  
262 DNA Enrichment Kit following the manufacturer's recommendations with modifications (New England Biolabs, Ipswich, MA). For all samples  
263 tested, 250 ng of DNA was subjected to two rounds of bead capture using the NEBNext Microbiome DNA Enrichment Kit and enriched  
264 treponemal genomic DNA was purified using AMPure XP beads (Beckman Coulter, Indianapolis, IN). Enriched DNA samples were stored at -  
265 20°C until MDA was performed. MDA was carried out using the REPLI-g Single Cell Kit following the manufacturer's recommendations with  
266 slight modifications (Qiagen). Each MDA reaction was incubated at 30°C for 16 hr. Following amplification, the polymerase was inactivated at  
267 65°C for 10 min, samples were purified with AMPure XP beads, and eluted with 100 µl 1X AE Buffer (Qiagen). For each enrichment using the  
268 REPLI-g Single Cell Kit, non-template controls were included to confirm the absence of *T. pallidum*.

269 A 10-fold dilution series on the Nichols strain was used to determine the limit of detection (LoD) for enrichment (see supplemental  
270 materials) with NEB+MDA followed by sequencing on an Illumina NovaSeq 6000. After DNA extraction, each dilution in the series was enriched  
271 by NEB+MDA, genomic copy numbers estimated by *polA* qPCR, and sequencing performed in triplicate. Enriched samples were diluted 1:10  
272 prior to measuring RNP amplification. The LoD was set at the minimal genome copy number required to generate a  $\geq 5X$  read depth with  $\geq 95\%$   
273 genome coverage compared to the reference genome.

274 **Selective whole genome amplification (SWGA) primer design, validation, and enrichment.** Primers with an increased affinity to *T. pallidum*  
275 were identified using the *swga* Toolkit as previously described with slight modifications (<https://www.github.com/eclarke/swga>; 26; see  
276 supplemental materials). Eight primer sets (SWGA Pal 1-8), including 4 additional primer sets (SWGA Pal 9-12) generated by combining primers



277 in the initial set (Table S1), were chosen for SWGA using the EquiPhi29 DNA Polymerase (Thermo Fisher Scientific, Waltham, MA). To account  
278 for the 3'-5' exonuclease activity of the phi29 polymerase, all SWGA primers were generated with phosphorothioate bonds between the last two  
279 nucleotides at the 3' end (Table S1). Each of the 12 primer sets were tested in triplicate against the spiked sample diluted to an estimated 100 *T.*  
280 *pallidum* *polA* copies/ $\mu$ l (see supplemental materials).

281 Prior to SWGA enrichment, samples were denatured for 5 min at 95°C by adding 2.5  $\mu$ l of DNA to 2.5  $\mu$ l denaturing solution, containing  
282 custom primers, then placed immediately on ice until the EquiPhi29 master mix, prepared as per manufacturer's recommendations, was added  
283 (Thermo Fisher Scientific, Waltham, MA). MDA was carried out following the manufacturer's recommendations with modifications (Thermo  
284 Fisher Scientific; 30). The reaction contained EquiPhi29 master mix, with EquiPhi29 Reaction Buffer at a final concentration of 1X, each primer  
285 at a final concentration of 4  $\mu$ M, and nuclease-free H<sub>2</sub>O was added to a final reaction volume of 20  $\mu$ l. Reaction tubes were gently mixed by pulse  
286 vortexing and incubated at 45°C for 3 hr. MDA was stopped by inactivating the DNA polymerase at 65°C for 15 min. All reactions were purified  
287 using AMPure XP beads and eluted in 100  $\mu$ l AE buffer (Qiagen). Non-template controls were included to confirm the absence of contaminate *T.*  
288 *pallidum* DNA.

289 Relative percent *T. pallidum* in each sample was calculated as shown in Figure S1. SWGA Pal 11 was chosen for testing the LoD for  
290 downstream genome sequencing post-SWGA enrichment using the 10-fold dilution series, excluding the undiluted (neat) spiked sample. All  
291 enriched samples were validated by *polA* real-time qPCR in triplicate.

292 **Sequencing and genome analysis of *T. pallidum* strains.** Libraries were prepared using the NEBNext Ultra DNA Library Preparation Kit for  
293 NovaSeq and NEBNext Ultra II FS DNA Library Preparation Kit for MiSeq sequencing following the manufacturer's recommendations (New  
294 England Biolabs, Ipswich, MA). For the validation experiments, sequencing was carried out on the Nichols reference strain using the Illumina

295 NovaSeq 6000 platform following the manufacturer's recommendations (Illumina, San Diego, CA). Sequencing of isolate CDC-SF003 and swab  
296 specimens were carried out using the MiSeq v2 (500 cycle) platform following the manufacturer's recommendations (Illumina, San Diego, CA).

297 Post sequencing, reads were deduplicated, trimmed, and down selected for *T. pallidum* (supplemental materials). All down selected *T.*  
298 *pallidum* reads were mapped to the *T. pallidum* reference genomes, and *de novo* assembled. Phylogenetic analyses were performed as described in  
299 the supplemental materials. Apart from the genomes sequenced in this study, 122 high quality (with at least 5x read depth covering > 90% of the  
300 genome) *T. pallidum* genomes deposited in the NCBI's Sequencing Read Archive (SRA) under the BioProject number PRJEB20795 and  
301 PRJNA508872 were also included (12, 14). The publicly available raw sequencing data were re-analyzed to determine the quality as described in  
302 the supplemental materials. A second phylogenetic tree was also reconstructed by including all the genomes sequenced from the 10-fold dilution  
303 series for both NEB+MDA and SWGA enriched samples. Genomic sequencing data from samples included in the phylogenetic analyses covered  
304 at least 90% of the reference genome with 5X read depth. Variant calls for the A2058G and A2059G macrolide resistance mutations were  
305 validated using a real-time PCR assay as previously described (31).

306 **Statistical analyses.** Statistical analyses were performed in R (R Foundation for Statistical Computing, Vienna, Austria) using the R companion  
307 software RStudio (Rstudio, Boston, MA). Statistical significance was determined by analysis of variance (ANOVA) and Tukey post hoc multiple  
308 comparisons tests. *T. pallidum* percent DNA were normalized through Log<sub>10</sub> conversions. Quantitative data are presented as means ± standard  
309 error. Differences were considered statistically significant if a P < 0.05.

310 **Data availability.** All sequencing data associated with this study were submitted to the National Center for Biotechnology Information's sequence  
311 read archive (SRA) under the BioProject accession ID PRJNA744275.

## 312 **Acknowledgments**

313 We thank Teresa Burns at the Emory University Hospital Midtown; Tamara Jones from the St Louis County STD Clinic; Yetty Fakile,  
314 Kevin Pettus, and Jack Cartee at CDC's Division of STD Prevention; The veterinary staff at the CDC's Comparative Medicine Branch; Mark Itsko  
315 at CDC's Division of Bacterial Diseases; and Nikhat Sulaiman and Justin Lee at CDC's Division of Scientific Resources for their assistance,  
316 consults, and support throughout this study. This work was made possible through CDC's Division of STD Prevention with support from the  
317 Advanced Molecular Detection (AMD) program.

## 318 **Author Contributions**

319 Allan Pillay and Ellen N. Kersh conceived the study. Allan Pillay, Charles M. Thurlow, Cheng Chen, and Lilia Ganova-Raeva designed  
320 the study. Charles M. Thurlow and Allan Pillay designed the enrichment protocols. Charles M. Thurlow designed the SWGA specific custom  
321 primer sets used during this study and performed all enrichment experiments. Charles M. Thurlow, Allan Pillay, Samantha S. Katz, Lara Pereira,  
322 Alyssa Debra, Kendra Vilfort, Yongcheng Sun, Kai-Hua Chi, and Damien Danavall performed the laboratory experiments and assisted with  
323 specimen collection. Kimberly Workowski, Stephanie E. Cohen, Hilary Reno, and Susan S. Philip collected clinical specimens and patient data.  
324 Mark Burroughs, Mili Sheth, and Charles M. Thurlow performed Illumina sequencing. Sandeep J. Joseph performed the bioinformatic analyses of  
325 the genomic data, phylogenetic analysis and contributed to the generation of tables and figures. Charles M. Thurlow and Sandeep J. Joseph  
326 performed data analysis. Charles M. Thurlow wrote and prepared the manuscript with oversight by Allan Pillay and contributions from Sandeep J.  
327 Joseph and Weiping Cao, which was reviewed by all authors for revisions.

## 328 **Disclaimer**

329           The findings and conclusions in this report are those of the authors and do not necessarily represent the official position of the Centers for  
330 Disease Control and Prevention. We declare that there are no competing interests.

## 331 **References**

- 332 1. CDC. 2020. Sexually Transmitted Disease Surveillance 2019. US Department of Health and Human Services, Atlanta
- 333 2. Marra C, Sahi S, Tantaló L, Godornes C, Reid T, Behets F, Rompalo A, Klausner JD, Yin Y, Mulcahy F, Golden MR, Centurion-Lara A,  
334 Lukehart SA. 2010. Enhanced molecular typing of *Treponema pallidum*: geographical distribution of strain types and association with  
335 neurosyphilis. *J Infect Dis* 202:1380-8.
- 336 3. Pillay A, Lee M-K, Slezak T, Katz SS, Sun Y, Chi K-H, Morshed M, Philip S, Ballard RC, Chen CY. 2019. Increased Discrimination of  
337 *Treponema pallidum* Strains by Subtyping With a 4-Component System Incorporating a Mononucleotide Tandem Repeat in *rpsA*.  
338 *Sexually Transmitted Diseases* 46:e42-e45.
- 339 4. Katz K, Pillay A, Ahrens K, Kohn R, Hermanstyné K, Bernstein K, Ballard R, Klausner J. 2010. Molecular Epidemiology of Syphilis—  
340 San Francisco, 2004–2007. *Sexually Transmitted Diseases* 37:660-3.
- 341 5. Grillová L, Bawa T, Mikalová L, Gayet-Ageron A, Nieselt K, Strouhal M, Sednaoui P, Ferry T, Cavassini M, Lautenschlager S, Dutly F,  
342 Pla-Díaz M, Krützen M, González-Candelas F, Bagheri HC, Šmajš D, Arora N, Bosshard PP. 2018. Molecular characterization of  
343 *Treponema pallidum* subsp. *pallidum* in Switzerland and France with a new multilocus sequence typing scheme. *PLoS One* 13:e0200773.
- 344 6. Edmondson DG, Wormser GP, Norris SJ. 2020. In Vitro Susceptibility of *Treponema pallidum pallidum* to Doxycycline. *Antimicrobial*  
345 *Agents and Chemotherapy* 64:e00979-20.

- 346 7. Pereira LE, Katz SS, Sun Y, Mills P, Taylor W, Atkins P, Thurlow CM, Chi K-H, Danavall D, Cook N, Ahmed T, Debra A, Philip S,  
347 Cohen S, Workowski KA, Kersh E, Fakile Y, Chen CY, Pillay A. 2020. Successful isolation of *Treponema pallidum* strains from patients'  
348 cryopreserved ulcer exudate using the rabbit model. PLOS ONE 15:e0227769.
- 349 8. Bachmann NL, Rockett RJ, Timms VJ, Sintchenko V. 2018. Advances in Clinical Sample Preparation for Identification and  
350 Characterization of Bacterial Pathogens Using Metagenomics. Frontiers in Public Health 6.
- 351 9. Thorburn F, Bennett S, Modha S, Murdoch D, Gunson R, Murcia PR. 2015. The use of next generation sequencing in the diagnosis and  
352 typing of respiratory infections. Journal of Clinical Virology 69:96-100.
- 353 10. Lefterova MI, Suarez CJ, Banaei N, Pinsky BA. 2015. Next-Generation Sequencing for Infectious Disease Diagnosis and Management: A  
354 Report of the Association for Molecular Pathology. The Journal of Molecular Diagnostics 17:623-634.
- 355 11. Beale MA, Marks M, Cole MJ, Lee M-K, Pitt R, Ruis C, Balla E, Crucitti T, Ewens M, Fernández-Naval C, Grankvist A, Guiver M,  
356 Kenyon CR, Khairulin R, Kularatne R, Arando M, Molini BJ, Obukhov A, Page EE, Petrovay F, Rietmeijer C, Rowley D, Shokoples S,  
357 Smit E, Sweeney EL, Taiaroa G, Vera JH, Wennerås C, Whiley DM, Williamson DA, Hughes G, Naidu P, Unemo M, Krajden M,  
358 Lukehart SA, Morshed MG, Fifer H, Thomson NR. 2021. Contemporary syphilis is characterised by rapid global spread of pandemic  
359 *Treponema pallidum* lineages. medRxiv doi:10.1101/2021.03.25.21250180:2021.03.25.21250180.
- 360 12. Beale MA, Marks M, Sahi SK, Tantalo LC, Nori AV, French P, Lukehart SA, Marra CM, Thomson NR. 2019. Genomic epidemiology of  
361 syphilis reveals independent emergence of macrolide resistance across multiple circulating lineages. Nature Communications 10:3255.

- 362 13. Pinto M, Borges V, Antelo M, Pinheiro M, Nunes A, Azevedo J, Borrego MJ, Mendonça J, Carpinteiro D, Vieira L, Gomes JP. 2016.  
363 Genome-scale analysis of the non-cultivable *Treponema pallidum* reveals extensive within-patient genetic variation. Nature Microbiology  
364 2:16190.
- 365 14. Grillová L, Oppelt J, Mikalová L, Nováková M, Giacani L, Niesnerová A, Noda AA, Mechaly AE, Pospíšilová P, Čejková D, Grange PA,  
366 Dupin N, Strnadel R, Chen M, Denham I, Arora N, Picardeau M, Weston C, Forsyth RA, Šmajš D. 2019. Directly Sequenced Genomes of  
367 Contemporary Strains of Syphilis Reveal Recombination-Driven Diversity in Genes Encoding Predicted Surface-Exposed Antigens.  
368 Frontiers in microbiology 10:1691-1691.
- 369 15. Arora N, Schuenemann VJ, Jäger G, Peltzer A, Seitz A, Herbig A, Strouhal M, Grillová L, Sánchez-Busó L, Kühnert D, Bos KI, Davis  
370 LR, Mikalová L, Bruisten S, Komericki P, French P, Grant PR, Pando MA, Vaulet LG, Fermepin MR, Martinez A, Centurion Lara A,  
371 Giacani L, Norris SJ, Šmajš D, Bosshard PP, González-Candelas F, Nieselt K, Krause J, Bagheri HC. 2016. Origin of modern syphilis and  
372 emergence of a pandemic *Treponema pallidum* cluster. Nature Microbiology 2:16245.
- 373 16. Chen W, Šmajš D, Hu Y, Ke W, Pospíšilová P, Hawley KL, Caimano MJ, Radolf JD, Sena A, Tucker JD, Yang B, Juliano JJ, Zheng H,  
374 Parr JB. 2021. Analysis of *Treponema pallidum* Strains From China Using Improved Methods for Whole-Genome Sequencing From  
375 Primary Syphilis Chancres. The Journal of Infectious Diseases 223:848-853.
- 376 17. Lukehart SA, Godornes C, Molini BJ, Sonnett P, Hopkins S, Mulcahy F, Engelman J, Mitchell SJ, Rompalo AM, Marra CM, Klausner  
377 JD. 2004. Macrolide resistance in *Treponema pallidum* in the United States and Ireland. N Engl J Med 351:154-8.

- 378 18. Chi KH, Danavall D, Taleo F, Pillay A, Ye T, Nachamkin E, Kool JL, Fegan D, Asiedu K, Vestergaard LS, Ballard RC, Chen CY. 2015.  
379 Molecular differentiation of *Treponema pallidum* subspecies in skin ulceration clinically suspected as yaws in Vanuatu using real-time  
380 multiplex PCR and serological methods. *Am J Trop Med Hyg* 92:134-8.
- 381 19. Workowski KA, Bolan GA. 2015. Sexually transmitted diseases treatment guidelines, 2015. *MMWR Recomm Rep* 64:1-137.
- 382 20. Čejková D, Zbaníková M, Chen L, Pospíšilová P, Strouhal M, Qin X, Mikalová L, Norris SJ, Muzny DM, Gibbs RA, Fulton LL,  
383 Sodergren E, Weinstock GM, Šmajš D. 2012. Whole Genome Sequences of Three *Treponema pallidum* ssp. *pertenue* Strains: Yaws and  
384 Syphilis *Treponemes* Differ in Less than 0.2% of the Genome Sequence. *PLOS Neglected Tropical Diseases* 6:e1471.
- 385 21. Pětrošová H, Pospíšilová P, Strouhal M, Čejková D, Zbaníková M, Mikalová L, Sodergren E, Weinstock GM, Šmajš D. 2013.  
386 Resequencing of *Treponema pallidum* ssp. *pallidum* Strains Nichols and SS14: Correction of Sequencing Errors Resulted in Increased  
387 Separation of Syphilis *Treponeme* Subclusters. *PLOS ONE* 8:e74319.
- 388 22. Sun J, Meng Z, Wu K, Liu B, Zhang S, Liu Y, Wang Y, Zheng H, Huang J, Zhou P. 2016. Tracing the origin of *Treponema pallidum* in  
389 China using next-generation sequencing. *Oncotarget* 7.
- 390 23. Feehery GR, Yigit E, Oyola SO, Langhorst BW, Schmidt VT, Stewart FJ, Dimalanta ET, Amaral-Zettler LA, Davis T, Quail MA, Pradhan  
391 S. 2013. A method for selectively enriching microbial DNA from contaminating vertebrate host DNA. *PLoS One* 8:e76096.
- 392 24. Thoendel M, Jeraldo PR, Greenwood-Quaintance KE, Yao JZ, Chia N, Hanssen AD, Abdel MP, Patel R. 2016. Comparison of microbial  
393 DNA enrichment tools for metagenomic whole genome sequencing. *Journal of Microbiological Methods* 127:141-145.
- 394 25. Leichty AR, Brisson D. 2014. Selective Whole Genome Amplification for Resequencing Target Microbial Species from Complex Natural  
395 Samples. *Genetics* 198:473-481.

- 396 26. Clarke EL, Sundararaman SA, Seifert SN, Bushman FD, Hahn BH, Brisson D. 2017. swga: a primer design toolkit for selective whole  
397 genome amplification. *Bioinformatics* 33:2071-2077.
- 398 27. Itsko M, Retchless AC, Joseph SJ, Norris Turner A, Bazan JA, Sadjji AY, Ouédraogo-Traoré R, Wang X. 2020. Full Molecular Typing of  
399 *Neisseria meningitidis* Directly from Clinical Specimens for Outbreak Investigation. *Journal of Clinical Microbiology* 58:e01780-20.
- 400 28. Sundararaman SA, Plenderleith LJ, Liu W, Loy DE, Learn GH, Li Y, Shaw KS, Ayoub A, Peeters M, Speede S, Shaw GM, Bushman  
401 FD, Brisson D, Rayner JC, Sharp PM, Hahn BH. 2016. Genomes of cryptic chimpanzee *Plasmodium* species reveal key evolutionary  
402 events leading to human malaria. *Nature Communications* 7:11078.
- 403 29. Cowell AA-O, Loy DE, Sundararaman SA, Valdivia H, Fisch K, Lescano AG, Baldeviano GC, Durand S, Gerbasi V, Sutherland CJ,  
404 Nolder D, Vinetz JM, Hahn BH, Winzeler EA. 2017. Selective Whole-Genome Amplification Is a Robust Method That Enables Scalable  
405 Whole-Genome Sequencing of *Plasmodium vivax* from Unprocessed Clinical Samples. *mBio* 8:e02257-16.
- 406 30. Povilaitis T, Alzbutas G, Sukackaite R, Siurkus J, Skirgaila R. 2016. In vitro evolution of phi29 DNA polymerase using isothermal  
407 compartmentalized self replication technique. *Protein Eng Des Sel* 29:617-628.
- 408 31. Chen C-Y, Chi K-H, Pillay A, Nachamkin E, Su JR, Ballard RC. 2013. Detection of the A2058G and A2059G 23S rRNA Gene Point  
409 Mutations Associated with Azithromycin Resistance in *Treponema pallidum* by Use of a TaqMan Real-Time Multiplex PCR Assay.  
410 *Journal of Clinical Microbiology* 51:908-913.

411

412 **Tables**



**Table 1.** Clinical and laboratory data for specimens and clinical isolate CDC-SF003.

Sample/ Isolate ID	Collection Year	Source	Gender	Sexual Status	Syphilis Stage	Site of Lesion	Antibody Titer (Assay)	qPCR ( <i>T. pallidum</i> <i>polA</i> in DNA Extract)	Extraction Method	Reference
CDC-SF003	2017	San Francisco	Male	MSM	Primary	Penis	1:4 (VDRL)	9,680 copies/μl	Standard	Pereira <i>et al.</i> , 2020
EUHM-001	2019	Atlanta	Male	MSM	Secondary	Neck	1:128 (RPR)	< 1 copy/μl	Standard	This study
EUHM-002	2019	Atlanta	Male	MSM	Secondary	Perianal	1:256 (RPR)	< 1 copy/μl	Standard	This study
EUHM-003	2019	Atlanta	Male	MSM	Secondary	Penis	1:32 (RPR)	< 1 copy/μl	Standard	This study
EUHM-004	2019	Atlanta	Male	MSM	Primary	Penis	1:4 (RPR)	106.7 ± 6.5 copies/μl	Standard	This study
EUHM-005	2019	Atlanta	Male	MSM	Secondary	Penis	1:64 (RPR)	< 1 copy/μl	Standard	This study
EUHM-006	2019	Atlanta	Male	MSM	Primary	Penis	1:16 (RPR)	< 1 copy/μl	Standard	This study
EUHM-007	2019	Atlanta	Male	MSM	Secondary	Hand	1:64 (RPR)	< 1 copy/μl	Standard	This study
EUHM-008	2019	Atlanta	Male	MSM	Secondary	Scrotum	1:64 (RPR)	0.9 ± 0.1 copy/μl	Standard	This study
EUHM-009	2019	Atlanta	Male	MSM	Secondary	Scrotum	1:64 (RPR)	< 1 copy/μl	Standard	This study
EUHM-010	2019	Atlanta	Male	MSM	Secondary	Scrotum	1:128 (RPR)	< 1 copy/μl	Standard	This study
EUHM-011	2019	Atlanta	Male	MSM	Primary	Penis	1:32 (RPR)	< 1 copy/μl	Standard	This study
EUHM-012	2019	Atlanta	Male	MSM	Primary	Penis	1:8 (RPR)	31.5 ± 0.5 copies/μl	Large Scale	This study
EUHM-013	2020	Atlanta	Male	MSM	Secondary	Penis	1:64 (RPR)	122 ± 1.2 copies/μl	Large Scale	This study
EUHM-014	2020	Atlanta	Male	MSM	Secondary	NA*	1:16 (RPR)	103 ± 6.7 copies/μl	Large Scale	This study
STLC-001	2020	St. Louis	Male	MSW	Primary	Penis	NR** (RPR)	28.8 ± 3.1 copies/μl	Standard	This study

\* Not available

\*\* Non-reactive

**Table 2.** Sequencing percent coverage for the Nichols isolates, clinical isolate CDC-SF003, and clinical specimens across the *T. pallidum* reference genome.

Sample	Enrichment method*	Clonal complex	<i>T.pallidum</i> <i>polA</i> post enrichment genome copies/μl	Raw read pairs	Non-host read pairs	Total read pairs after QC	Read pairs classified as <i>T. pallidum</i>	Percent of total read pairs classified as <i>T. pallidum</i>	Mean read depth	Percent genome covered ≥1X	Percent genome covered ≥5X	Percent genome covered ≥10X
Nichols_CDC	non-enriched	Nichols-like	NA***	4,053,500	3,645,649	3,588,414	70,299	1.96	6.33	86.26	60.30	22.28
Nichols_CDC**	SWGA	Nichols-like	11,565,333 ± 1,294,672	3,701,303	3,692,932	3,648,044	3,414,111	93.59	751.17	98.39	98.24	98.16
CDC-SF003	NEB + MDA	SS14-like	2,394,930 ± 135,210	5,798,777	3,988,173	3,949,036	129,998	3.29	46.44	98.87	98.60	98.01
EUHM-004	SWGA	Nichols-like	6,367,089.5 ± 240,811.5	6,102,826	4,440,618	4,280,401	1,403,645	32.79	370.39	96.99	95.13	92.67
EUHM-012	SWGA	Nichols-like	2,140,753 ± 28,192	10,350,274	5,870,287	5,716,082	2,793,693	48.87	639.86	96.34	93.98	91.89
EUHM-013	SWGA	SS14-like	5,159,716 ± 220,318.5	11,975,324	11,966,460	11,838,431	8,308,234	70.18	2,503.96	98.72	98.56	98.37
EUHM-014	SWGA	Nichols-like	2,573,508 ± 221,900.5	11,250,518	9,266,926	9,059,022	2,355,426	26.00	930.87	98.79	98.49	98.04
STLC-001	SWGA	SS14-like	7,420,534 ± 719,765	11,293,960	7,770,834	7,721,767	3,004,631	38.91	1,133.43	98.32	95.94	94.10

\*All sequencing was performed using Illumina's MiSeq v2 (500 cycle) platform

\*\* Enrichment performed on 1,000 copies/μl *T. pallidum* *polA* input

\*\*\* Not available

414 **Figures**

415 **Fig 1.** *T. pallidum* *polA* copies/ $\mu$ l for the 10-fold dilution series spiked samples enriched by the NEBNext  
416 Microbiome Enrichment Kit with REPLIg Single Cell MDA (NEB+MDA) or SWGA. The input *T.*  
417 *pallidum* *polA* copies/ $\mu$ l for each dilution is displayed as Non-Enriched. The y-axis has been  $\log_{10}$  scaled  
418 for depiction of the Non-Enriched dilution series. Error bars represent standard error among three  
419 replicate enriched *T. pallidum* samples.

420 **Fig 2.** Relative percent *T. pallidum* Nichols DNA for Non-Enriched, NEBNext Microbiome Enrichment  
421 Kit with REPLIg Single Cell MDA (NEB+MDA), and SWGA enriched samples. Percent *T. pallidum*  
422 DNA was calculated based on the input DNA concentration and *polA* copies/ $\mu$ l (Non-Enriched), and the  
423 DNA concentration and *polA* copies/ $\mu$ l for the Nichols -spiked samples post-enrichment (NEB+MDA or  
424 SWGA). The y-axis has been  $\log_{10}$  scaled for depiction of the Non-Enriched dilution series. Error bars  
425 represent standard error among three replicate samples.

426 **Fig 3.** Percent coverage of sequencing reads of enriched *T. pallidum* Nichols spiked samples. (A)  
427 Sequencing reads of samples enriched using the NEB Microbiome Enrichment Kit and REPLIg Single  
428 Cell MDA (NEB+MDA). (B) Sequencing reads of samples enriched using SWGA. All samples were  
429 sequenced using the Illumina NovaSeq 6000 platform. Error bars represent standard error between the  
430 mapped reads derived from three replicate enriched Nichols samples.

431 **Fig 4.** Percent coverage of isolates and clinical specimens. All samples were sequenced using the Illumina  
432 MiSeq v2 (500 cycle) platform. Percent of *T. pallidum* reads are derived from down selected *T. pallidum*  
433 reads. Prefiltered reads for Nichols-CDC were mapped to the Nichols reference genome (*NC\_000919.1*).  
434 The prefiltered reads in all clinical isolates and specimens were mapped against the SS14 reference  
435 genome (*NC\_021508.1*).

436 **Fig 5.** SWGA primer set validation. (A) *T. pallidum* *polA* copies/ $\mu$ l for the Nichols mock sample (1:100  
437 diluted) enriched with each SWGA primer set. (B) Relative percent *T. pallidum* DNA for the Nichols

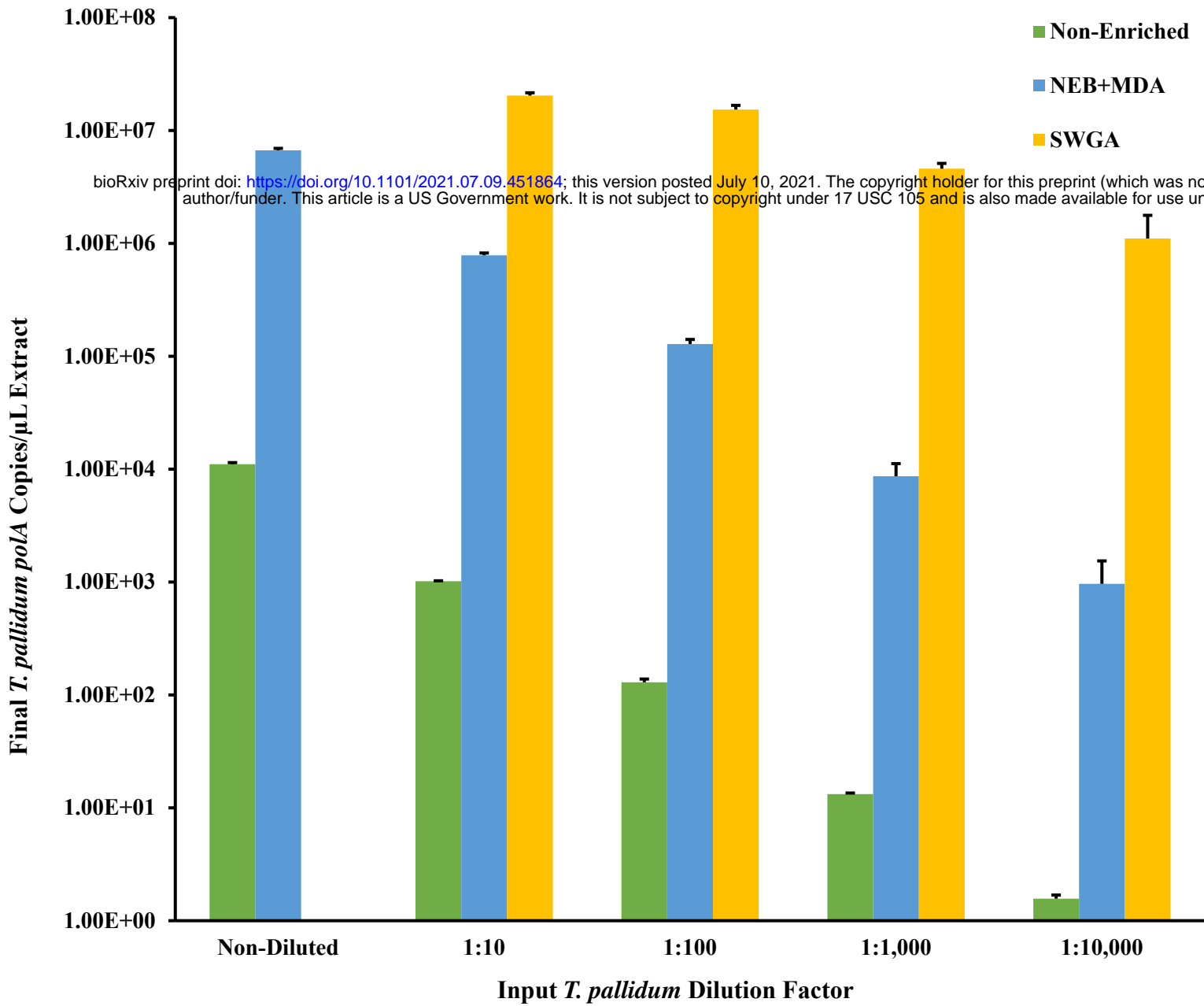
438 spiked sample (1:100 dilution) enriched with each SWGA primer set. Percent *T. pallidum* DNA was  
439 calculated based on the input DNA concentration and *polA* copies/ $\mu$ l for the Nichols mock samples post-  
440 SWGA enrichment. The y-axis has been  $\log_{10}$  scaled for depiction of the relative percent *T. pallidum*  
441 post-enrichment with each primer set. Error bars represent standard error among three replicate Nichols  
442 samples.

bioRxiv preprint doi: <https://doi.org/10.1101/2021.07.09.451864>; this version posted July 10, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. This article is a US Government work. It is not subject to copyright under 17 USC 105 and is also made available for use under aCC-BY 4.0 International license.

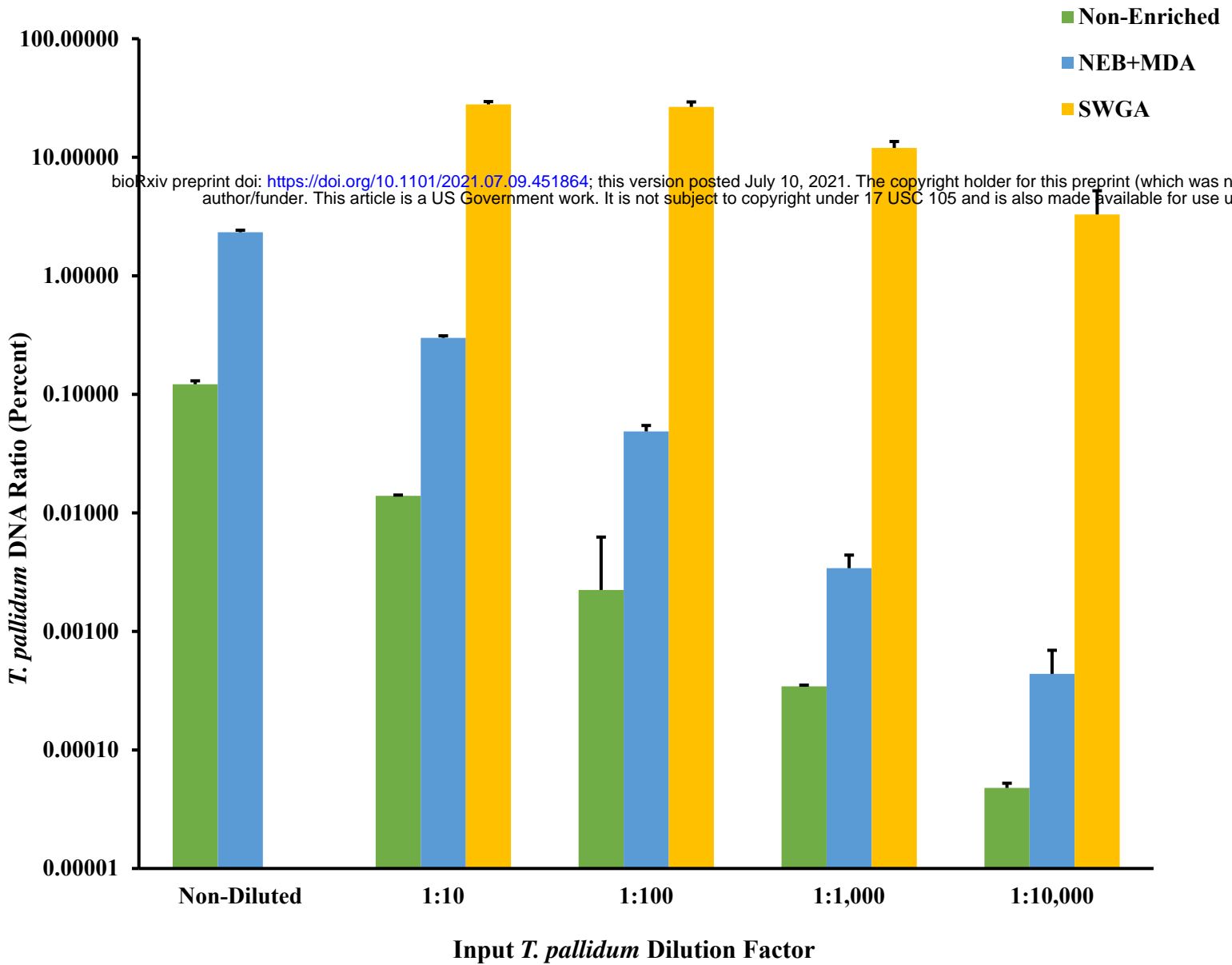
443 **Fig 6.** Maximum likelihood global phylogenetic tree of the clinical isolate/specimen genome sequenced in  
444 this study along with publicly available *T. pallidum* genomes. The two major lineages, Nichols-like and  
445 SS14-like are highlighted along with presence of genotypic mutation responsible for macrolide resistance  
446 and country of origin.

447

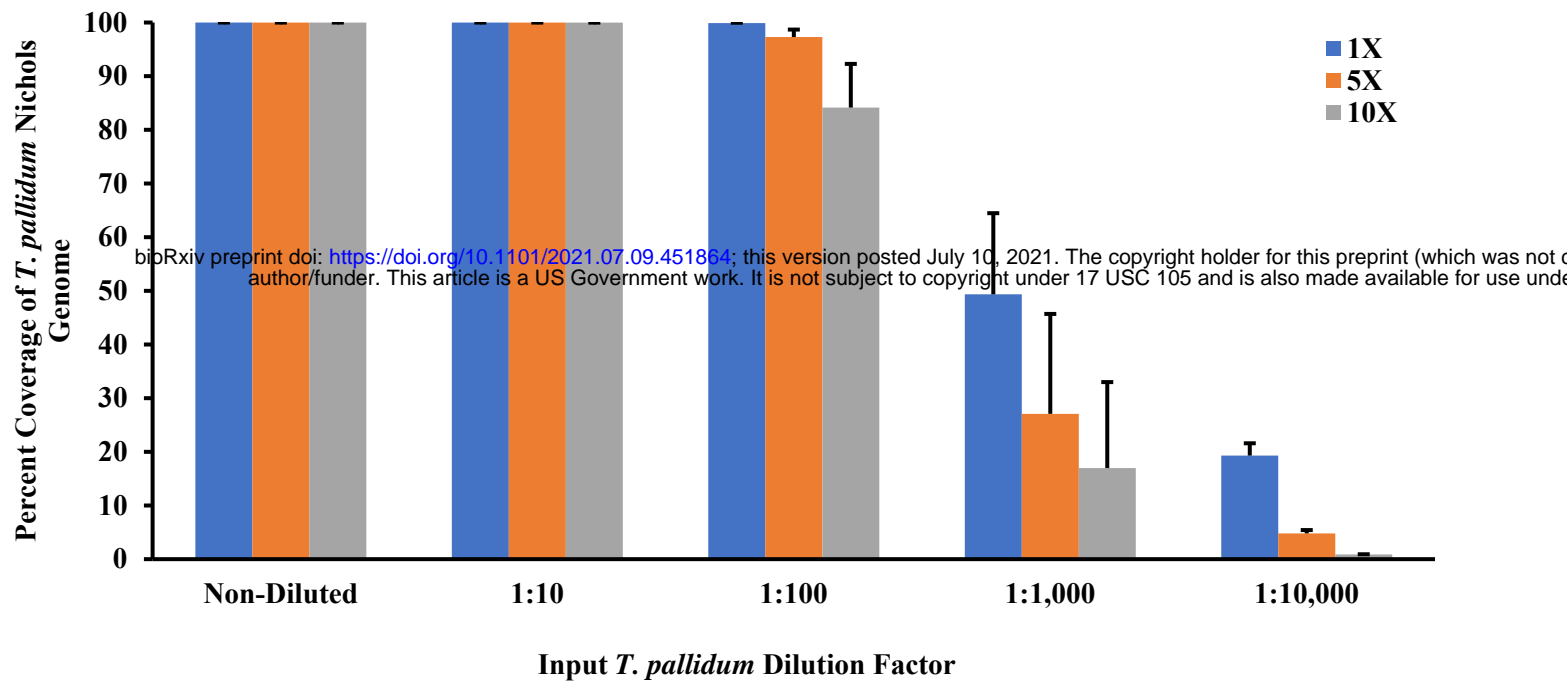
bioRxiv preprint doi: <https://doi.org/10.1101/2021.07.09.451864>; this version posted July 10, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. This article is a US Government work. It is not subject to copyright under 17 USC 105 and is also made available for use under aCC-BY 4.0 International license.



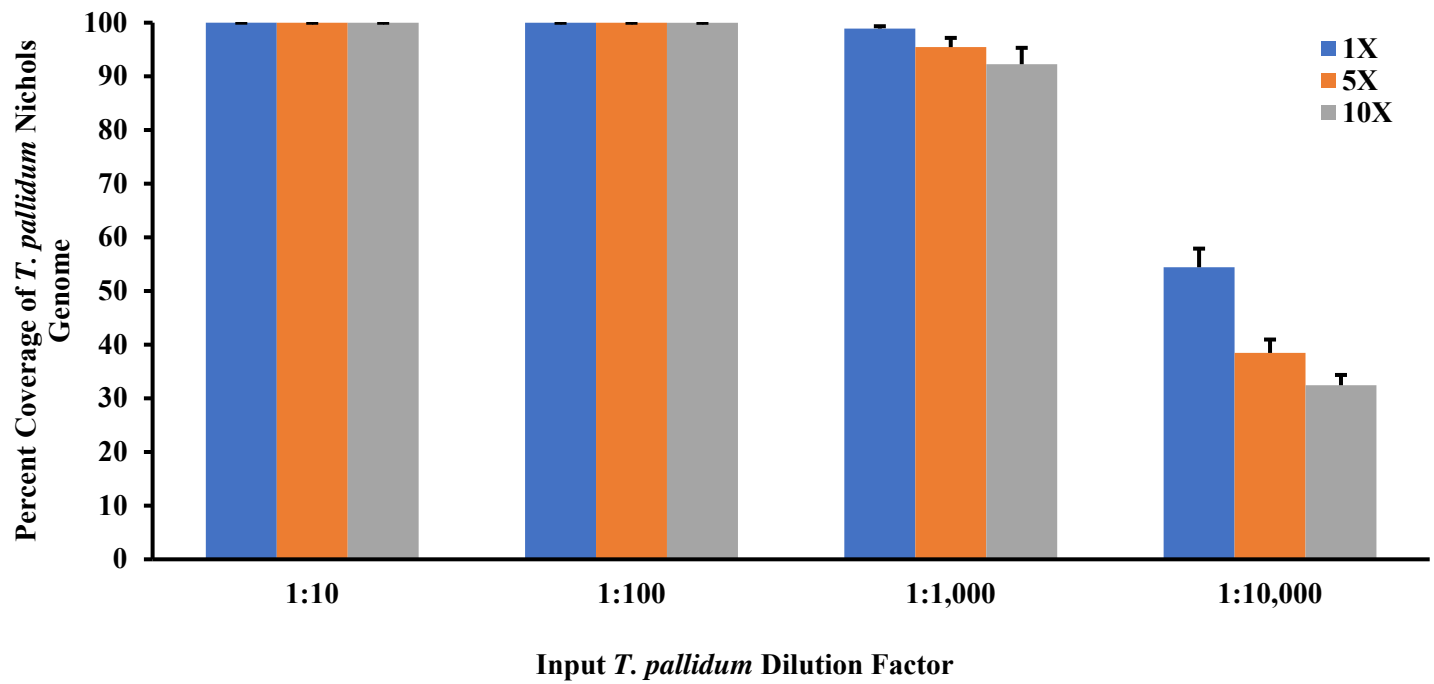
bioRxiv preprint doi: <https://doi.org/10.1101/2021.07.09.451864>; this version posted July 10, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. This article is a US Government work. It is not subject to copyright under 17 USC 105 and is also made available under aCC-BY 4.0 International license.

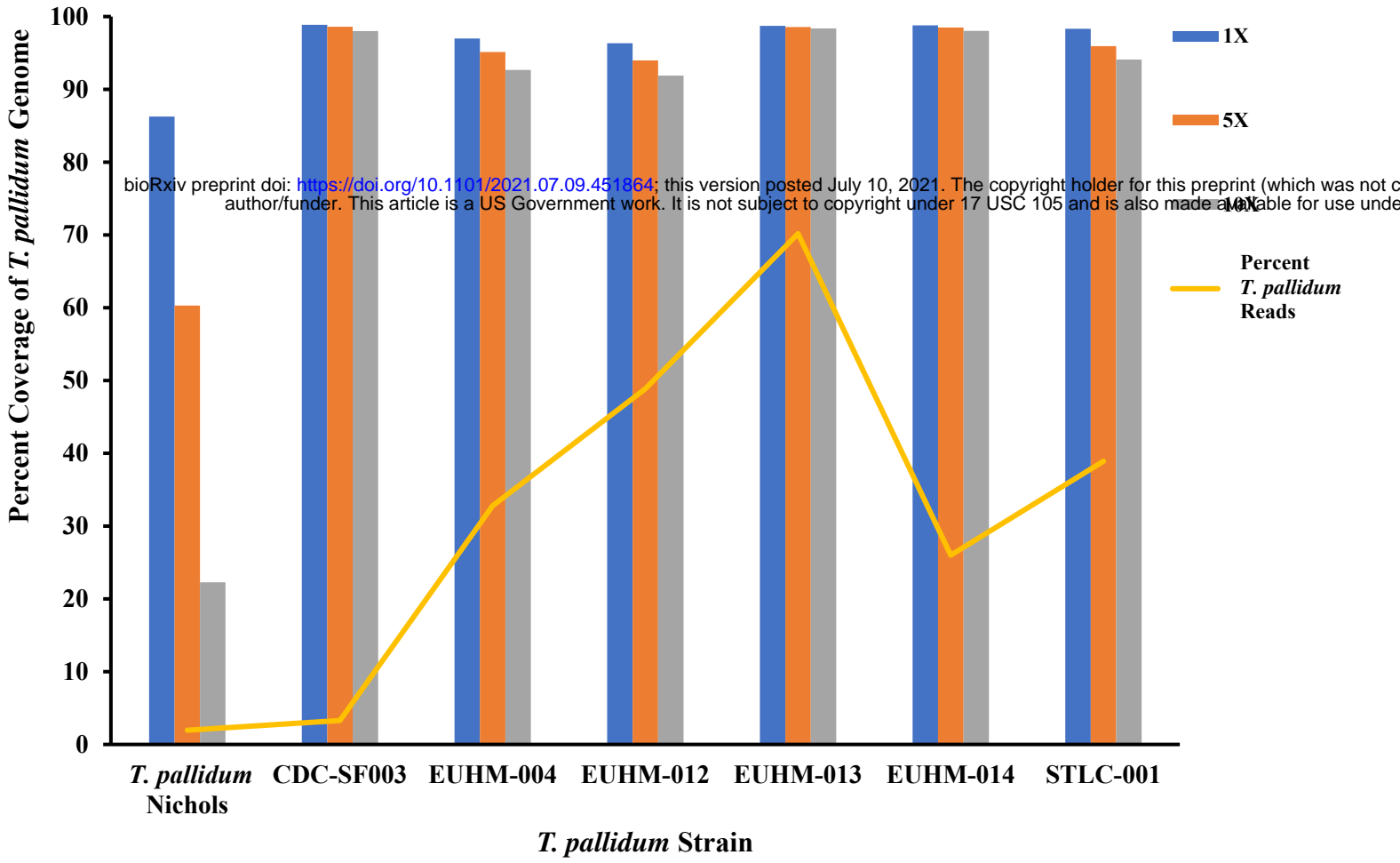


**A.**



**B.**

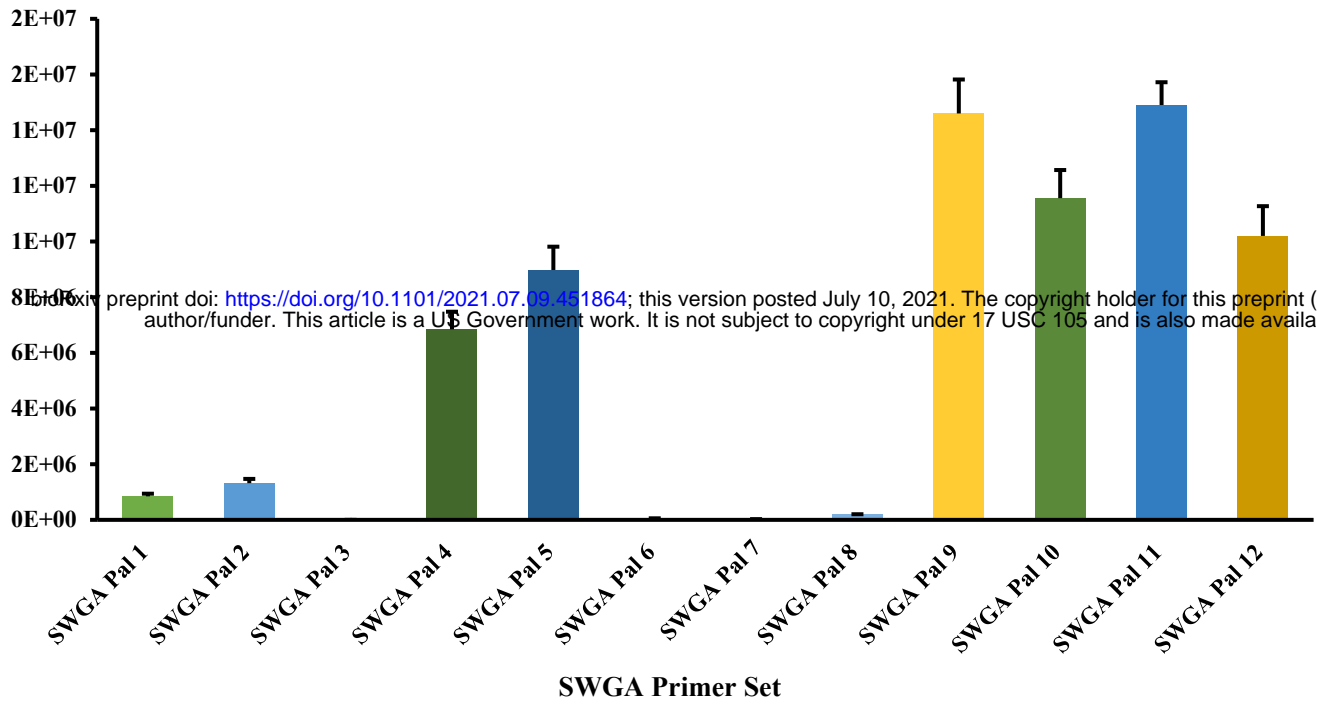






A.

Final *T. pallidum* *poLa* Copies/ $\mu$ L Extract



B.

*T. pallidum* DNA Ratio (Percent)

