1  **High-throughput nanopore sequencing of *Treponema pallidum* tandem repeat genes *arp* and**
2  **_tp0470_ reveals clade-specific patterns and recapitulates global whole genome phylogeny**
3

4  *Nicole AP Lieberman[1], Thaddeus D Armstrong[1†], Benjamin Chung[1†], Daniel Pfalmer[1], Christopher*
5  *M Hennelly[2], Austin Haynes[3], Emily Romeis[3], Qian-Qiu Wang[4,5], Rui-Li Zhang[6], Cai-Xia Kou[4,5],*
6  *Giulia Ciccarese[7], Ivano Dal Conte[8], Marco Cusini[9], Francesco Drago[7], Shu-ichi Nakayama[10],*
7  *Kenichi Lee[10], Makoto Ohnishi[10], Kelika A Konda[11,12], Silver K Vargas[11,13], Maria Eguiluz[11], Carlos*
8  *F Caceres[11], Jeffrey D Klausner[12], Oriol Mitja[14,15], Anne Rompalo[16], Fiona Mulcahy[17], Edward W*
9  *Hook 3r d[18,19,20], Irving F Hoffmann[2,21], Mitch M Matoga[2,21], Heping Zheng[22,23], Bin Yang[22,23],*
10  *Eduardo Lopez-Medina[24,25], Lady G Ramirez[24,26], Justin D Radolf[27,28,29,30,31], Kelly L*
11  *Hawley[27,28,30,32], Juan C Salazar[28,30,32], Sheila A Lukehart[3,33], Arlene C Seña[2], Jonathan B Parr[2],*
12  *Lorenzo Giacani[3,33], Alexander L Greninger[1,34,\*]*
13

14  [1] Department of Laboratory Medicine and Pathology, University of Washington School of
15  Medicine, Seattle, Washington, United States
16  [2] Division of Infectious Diseases and Institute for Global Health and Infectious Diseases,
17  University of North Carolina, Chapel Hill, NC, USA
18  [3] Department of Medicine, Division of Allergy and Infectious Diseases, University of Washington
19  School of Medicine, Seattle, Washington, United States of America
20  [4] Institute of Dermatology, Chinese Academy of Medical Science & Peking Union Medical
21  College, Beijing, China
22  [5] National Center for STD Control, China Centers for Disease Control and Prevention, Nanjing,
23  China
24  [6] Department of Dermatology, The Second Affiliated Hospital of Nanjing Medical University,
25  Nanjing, China
26  [7] Health Sciences Department, Section of Dermatology, San Martino University Hospital, Genoa,
27  Italy
28  [8] Sexual Health Center, Department of Prevention, ASL Città di Torino, Turin, Italy
29  [9] Fondazione IRCCS Ca' Granda, Ospedale Maggiore Policlinico, Milan, Italy
30  [10] Department of Bacteriology I, National Institute of Infectious Diseases, Tokyo, Japan
31  [11] Unit of Health, Sexuality and Human Development and Laboratory of Sexual Health,
32  Universidad Peruana Cayetano-Heredia, Lima, Peru
33  [12] Keck School of Medicine, University of Southern California, Los Angeles, California, United
34  States of America
35  [13] School of Public Health and Administration "Carlos Vidal Layseca", Universidad Peruana
36  Cayetano-Heredia, Lima, Peru
37  [14] Fight Aids and Infectious Diseases Foundation, Hospital Germans Trias i Pujol, Barcelona,
38  Spain
39  [15] Lihir Medical Centre-International SOS, Lihir Island, Papua New Guinea
40  [16] Department of Infectious Diseases, Johns Hopkins Medical Institutions, Baltimore, Maryland,
41  United States of America
42  [17] Department of Genito Urinary Medicine and Infectious Diseases, St James's Hospital, Dublin,
43  Ireland

44   [18] Department of Medicine, University of Alabama, Birmingham, Birmingham, Alabama, United
45   States of America
46   [19] Department of Epidemiology, University of Alabama, Birmingham, Birmingham, Alabama,
47   United States of America
48   [20] Department of Microbiology, University of Alabama, Birmingham, Birmingham, Alabama,
49   United States of America
50   [21] UNC-Project Malawi, Lilongwe, Malawi
51   [22] Dermatology Hospital of Southern Medical University, Guangzhou, P.R. China
52   [23] Institute for Global Health and Sexually Transmitted Infections, Guangzhou, P.R. China
53   [24] Centro Internacional de Entrenamiento e Investigaciones Medicas (CIDEIM), Cali, Colombia
54   [25] Centro de Estudios en Infectología Pediátrica (CEIP), Cali, Colombia
55   [26] Universidad ICESI, Cali, Colombia
56   [27] Department of Medicine, UConn Health, Farmington, CT, USA
57   [28] Department of Pediatrics, UConn Health, Farmington, CT, USA
58   [29] Department of Molecular Biology and Biophysics, UConn Health, Farmington, CT, USA
59   [30] Department of Immunology, UConn Health, Farmington, CT, USA
60   [31] Department of Genetics and Genome Sciences, UConn Health, Farmington, CT, USA
61   [32] Division of Infectious Diseases and Immunology, Connecticut Children's, Hartford, CT, USA
62   [33] Department of Global Health, University of Washington School of Medicine, Seattle,
63   Washington, United States of America
64   [34] Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle,
65   Washington, United States of America
66
67   * Corresponding Author, agrening@uw.edu
68   † Authors contributed equally
69
70   Short title: Nanopore sequencing *arp* and *tp0470*
71

72 **Abstract**

73  Sequencing of most *Treponema pallidum* (*T. pallidum*) genomes excludes repeat regions in
74  *tp0470* and the *tp0433* gene, encoding the acidic repeat protein (*arp*). As a first step to
75  understanding the evolution and function of these genes and the proteins they encode, we
76  developed a protocol to nanopore sequence *tp0470* and *arp* genes from 212 clinical samples
77  collected from ten countries on six continents. Both *tp0470* and *arp* repeat structures
78  recapitulate the whole genome phylogeny, with subclade-specific patterns emerging. The
79  number of *tp0470* repeats is on average appears to be higher in Nichols-like clade strains than
80  in SS14-like clade strains. Consistent with previous studies, we found that 14-repeat *arp*
81  sequences predominate across both major clades, but the combination and order of repeat
82  type varies among subclades, with many *arp* sequence variants limited to a single subclade.
83  Although strains that were closely related by whole genome sequencing frequently had the
84  same *arp* repeat length, this was not always the case. Structural modelling of TP0470 suggested
85  that the eight residue repeats form an extended $\alpha$-helix, predicted to be periplasmic. Modeling
86  of the ARP revealed a C-terminal sporulation-related repeat (SPOR) domain, predicted to bind
87  denuded peptidoglycan, with repeat regions possibly incorporated into a highly charged $\beta$-
88  sheet. Outside of the repeats, all TP0470 and ARP amino acid sequences were identical.
89  Together, our data, along with functional considerations, suggests that both TP0470 and ARP
90  proteins may be involved in *T. pallidum* cell envelope remodeling and homeostasis, with their
91  highly plastic repeat regions playing as-yet-undetermined roles.

92

93

**Introduction**

94  **Introduction**
95  In recent years, efforts to catalog genomic diversity and phylodynamics of the syphilis
96  spirochete, *Treponema pallidum* subsp. *pallidum* (*T. pallidum*), have resulted in a rapid increase
97  in the amount of sequencing data and number of near-complete genome assemblies available
98  in public databases (Arora et al., 2017; Beale et al., 2019, 2021; Chen et al., 2021; Grillová et al.,
99  2019; Lieberman et al., 2021; Pinto et al., 2016; Taouk et al., 2022; Thurlow et al., 2022).
100 Insights gained from these efforts, including the spread of azithromycin resistance (Beale et al.,
101 2019) and high-resolution information on antigenic diversity (Lieberman et al., 2021), have
102 aided our understanding of *T. pallidum* evolution and are invaluable to vaccine design. Although
103 whole genome sequencing of low abundance *T. pallidum* DNA directly from clinical specimens is
104 technically challenging due to the necessity of enrichment protocols such as hybrid capture
105 with RNA or DNA baits (Pinto et al., 2016; Arora et al., 2017; Beale et al., 2019, 2021; Lieberman
106 et al., 2021; Taouk et al., 2022), *Dpn1* enrichment (Grillová et al., 2019), whole genome
107 amplification (Chen et al., 2021; Thurlow et al., 2022), and/or the traditional technique of
108 passage of clinical strains through rabbits, sufficient progress has been made in development of
109 these techniques that sequencing throughput of samples on a scale appropriate for monitoring
110 of vaccine trials is feasible.
111
112 Most genomic analyses of *T. pallidum* have excluded portions of the genome difficult to resolve
113 by short-read sequencing, including the <u>T. pallidum r</u>epeat (*tpr*) family of paralogous genes, the
114 number of 60 bp tandem near-perfect repeats in the gene encoding the acidic repeat protein
115 (*arp; tp0433*), and the number of 24 bp tandem repeats in the gene encoding the
116 tetratricopeptide repeat protein TP0470 (*tp0470*). *Arp* repeat length has been interrogated
117 extensively using the CDC (Pillay et al., 1998) and enhanced CDC (Marra et al., 2010) typing
118 schemes. This has allowed monitoring of local strain composition over time (Marra et al., 2010;
119 Flasarová et al., 2012; Grimes et al., 2012) and of strains circulating worldwide (Sutton et al.,
120 2001; Pillay et al., 2002; Pope et al., 2005; Liu et al., 2020), and identification of subtypes
121 enriched in neurosyphilis (Molepo et al., 2007; Marra et al., 2010). However, these approaches
122 focus on obtaining the number of *arp* repeats, rather than providing sequence information on
123 the repeats, which is necessary for more robust genotyping. Our sequencing-based approach
124 starts to fill this important knowledge gap.
125
126 Little is known about the role of *tp0470* in syphilis pathogenesis, though tetratricopeptide
127 repeat proteins in other bacterial species act as scaffolds for protein-protein interactions and
128 are often critical to functionality of virulence factors (Cerveny et al., 2013). The highly charged
129 motif "EAEEARRK", encoded by the 24 bp repeat in *tp0470*, occurs C-terminal to the predicted
130 tetratricopeptide repeat domain. This motif is repeated between 4-29 times in publicly
131 available *T. pallidum* subsp. *pallidum* complete genomes, and up to 37 times in the *T. pallidum*
132 subsp. *pertenue* strain CDC-2. The role of the *arp* gene is similarly understudied. It encodes an
133 antigenic (Liu et al., 2007) protein containing at least four types of the 60 bp repeat that have
134 been previously identified in *T. pallidum* subsp. *pallidum*, with differences confined to six
135 positions, all of which result in amino acid substitutions (Liu et al., 2007; Harper et al., 2008).
136 Repeat lengths between 2-22 have been previously observed in *T. pallidum* subsp. *pallidum*
137 clinical specimens (Harper et al., 2008). Although the role of the acidic repeat protein in

138 pathogenesis and colonization of various anatomic sites is unknown, some samples collected
139 from whole blood had a lower number of *arp* tandem repeats than in patient-matched lesion
140 swabs (Mikalová et al., 2013), and late stage syphilis samples are reported to have fewer *arp*
141 repeats on average than early stage (Harper et al., 2008). Importantly, *arp* and *tp0470* are
142 thought to evolve via intra-strain recombination (Grillová et al., 2019; Noda et al., 2022).

143

144 To date, most closed *T. pallidum* genomes have relied on Sanger sequencing of *arp* and *tp0470*
145 (Cejková et al., 2012; Pětrošová et al., 2012; Zobaníková et al., 2013; Grillová et al., 2019),
146 requiring extensive manual curation of data as well as ample starting material. Herein, we
147 describe new bench and bioinformatic protocols for highly multiplexed nanopore sequencing of
148 *arp* and *tp0470*, reducing the quantity of sample needed, as well as per-sample cost and hands-
149 on time. These methodological improvements allowed us to gain insight into the evolution of
150 both *tp0470* and *arp* and formulate hypotheses that pave the way for functional studies to
151 understand the role of these putative virulence factors in syphilis pathogenesis.

152

153 **Methods**
154 **Ethics Statement**
155 All human samples were collected and deidentified following protocols established at each
156 institution. All IRB information from samples collected in Japan; Italy; Ireland; Maryland, USA;
157 Madagascar; Peru; Papua New Guinea, and Nanjing, China has been previously published
158 (Hopkins et al., 2004; Lukehart et al., 2004; Van Damme et al., 2009; Hook III et al., 2010; Marra
159 et al., 2010; Lieberman et al., 2021). Collection of additional samples was covered by the
160 following IRBs: Malawi: National Health Sciences Research Committee Ministry of Health and
161 Population (IRB Approval Number 2252); Colombia: Centro Internacional de Entrenamiento e
162 Investigaciones Medicas (CIDEIM) Institutional Human Research Ethics Committee (CIEIH) (IRB
163 protocol number 1289); Guangzhou, China: Dermatology Hospital of Southern Medical
164 University (SMU) Medical Ethic Committee (IRB protocol number GDDHLS-20181202(R3));
165 Chapel Hill, USA: University of North Carolina IRB Protocol Number 19-0311. Sequencing of
166 deidentified strains was covered by the University of Washington Institutional Review Board
167 (IRB) protocol number STUDY00000885 and University of North Carolina IRB protocol number
168 19-0311.

169

170 **Whole genome sequencing and maximum likelihood phylogeny**
171 Samples were sequenced and genomes assembled as previously described (Lieberman et al.,
172 2021) or, for samples from Colombia, Malawi, Guangzhou, China, and Chapel Hill, USA,
173 consensus sequences were assembled as previously described with minor modifications (Chen
174 et al., 2021; Parr, 2021). Alignment, recombination masking and generation of the maximum
175 likelihood phylogeny were performed as previously described (Lieberman et al., 2021).

176

177 **Barcoding PCR of *arp* and *tp0470*:**
178 Repeat regions of *arp* and *tp0470* were amplified using 96 combinations of forward and reverse
179 barcodes barcoded primers containing a 24bp index (Supporting Information). Input volume
180 ranged between 0.5-8 µL depending on genome copy number and amount of sample available.
181 Samples were amplified by one of two methods, both using the Takara PrimeSTAR GXL

182  polymerase in a 25 µL reaction: Those with low volume were first amplified with non-barcoded
183  primers (98˚C for 2 minutes, 35 cycles of 98˚C for 10 seconds, 62˚C for 15 seconds, 68˚C for 2
184  minutes, then held at 68˚C for 10 minutes before storing at 4˚C), the amplicon cleaned with
185  0.8x Ampure XP beads, diluted 1:100, and then 1 µL template barcoded with 14 additional
186  cycles of PCR using the indexed primers, using a 65˚C annealing temperature. Alternatively,
187  samples were amplified directly from the genomic DNA using barcoded primers. Both methods
188  produced equivalent results and technical replicates for each method agreed with each other
189  (Supporting Information, "Methods Equivalency")).
190  PCR products were electrophoresed using 1% or 2% TAE agarose gels and purified with 0.6x or
191  0.8x volumes of AMPure XP beads for *arp* and *tp0470* amplicons, respectively. Clean PCR
192  products were quantified using Qubit 1X dsDNA HS buffer (Invitrogen).
193
**Nanopore library preparation and sequencing**
195  Following barcoding PCR and purification/quantification, samples were pooled to meet
196  recommended input amounts for Oxford Nanopore (ONT) Adapter ligation kit (SQK-LSK109).
197  Because most samples fell in the range of 2-10 ng/uL, we chose to maximize efficiency by
198  pooling an equivalent volume (0.5 µL) per amplicon.
199
200  DNA end repair of the pool was performed as outlined in the ONT protocol for SQK-LSK109
201  using both ONT and NEBNext reagents, then purified with Ampure XP beads at a 1:1 ratio.
202  Adapter ligation proceeded according to the manufacturer's protocol, with room temperature
203  incubation for 10 minutes, followed by purification with a ratio of 0.8x Ampure XP beads,
204  washing with kit short fragment buffer (SFB), and elution in 47 µL water. The resulting pool was
205  then quantified on Qubit (1X high sensitivity dsDNA kit) to calculate the total molecular weight
206  of DNA. Molarity was calculated assuming an average fragment size of 800bp.
207
208   Amplicons were sequenced on Flongle flow cells via the MinION mk1B platform. ONT
209  MinKNOW software interfaced with the MinION to perform pre-run flow cell checks and
210  initiate/monitor the sequencing experiment. Each Flongle was primed and loaded with 20 fmol
211  DNA per SQK-LSK109 and Flongle sequencing expansion (EXP-FSE001).
212
213  Sequencing was run for 24 hours, selecting "SQK-LSK109" for the DNA amplicon kit with high
214  accuracy basecalling and other default parameters. Basecalling was performed in MinKNOW
215  (v21.11.7) running Guppy (v.5.1.12) (Wick et al., 2019). Reads with an average phred score
216  greater than 9 passed the quality filter.
217
**Demultiplexing**
219  Fastq files that passed the quality filter were processed through Porechop (Wick et al., 2017)
220  twice: First, using a customized adapters.py script that contained forward barcodes, using the
221  default stringency of up to 5 mismatches in the 24 bp barcodes. Barcodes were modeled after
222  those used previously in a dual-indexing protocol (Currin et al., 2019). Second round
223  demultiplexing proceeded using reverse primers with the same parameters in the adapters2.py
224  script, but ensuring no additional bases were trimmed past the barcodes.
225

226 ***arp* and *tp0470* consensus generation**
227 Demultiplexed reads were aligned with bwa mem 0.7.17 (Li and Durbin, 2009) using the
228 nanopore preset to reference files containing various numbers of tandem *arp* or *tp0470* repeats
229 and flanking regions of several hundred bases. Per-reference statistics were extracted with the
230 BBtools v38.18 (Bushnell) utility pileup.sh, and the reference assigned the most reads
231 determined for each sample using a custom R script. Consensus sequences were extracted
232 using the sam2consensus utility script (https://github.com/edgardomortiz/sam2consensus).
233 Predicted number of repeats and resultant band size was then cross referenced with the
234 agarose gel images, and the automated call either confirmed or overridden.
235
236 **Structural modeling**
237 Locations of predicted signal peptides and lipidation sites were determined in "slow" mode in
238 SignalP 6.0 (Teufel et al., 2022). Default settings were used for PSORT analysis (Yu et al., 2010).
239 Conserved domains were determined using the NCBI Conserved Domain Database (Lu et al.,
240 2020). Following removal of predicted secretion sequences, full length protein sequences for
241 TP0470 and ARP were modeled in trRosetta (Yang et al., 2020; Du et al., 2021) with default
242 settings and AlphaFold (Jumper et al., 2021; Varadi et al., 2022) using the pdb70 database for
243 homology modeling and otherwise default settings. Protein models were visualized in PyMOL,
244 and electrostatic surface potential shown with the Adaptive Poisson-Boltzmann Solver (APBS)
245 plugin.
246
247 **Statistics and code availability**
248 All statistical analysis was performed in R v4.0.3. Phylogenetic trees and *tp0470* and *arp*
249 variants were visualized with the R packages ggtree (Yu et al., 2017), treeio (Wang et al., 2020),
250 and ggplot (Wickham, 2016), and multiple sequence alignments by R package ggmsa (Zhou et
251 al., 2022). Bash, R, and python scripts for all data processing are available at
252 https://github.com/greninger-lab/TP_genome_finishing.
253
254 **Results**
255 **Maximum likelihood phylogeny of 207 *T. pallidum* subsp. *pallidum* strains**
256 We have previously reported the near-complete genomes of 196 *T. pallidum* strains, of which
257 191 were *T. pallidum* subsp. *pallidum* (Lieberman et al., 2021). Herein, we have updated the *T.*
258 *pallidum* subsp. *pallidum* whole genome phylogeny to include 16 additional strains collected in
259 Colombia, China, Malawi, and the United States. With the exception of one strain collected
260 from Malawi (TPVMW082H), all newly-added strains fit within the subclades defined previously
261 (Lieberman et al., 2021) (Fig 1A; Supporting Information). Malawi strain TPVMW082H appears
262 to have diverged from the lineage that gave rise to the Nichols B and Nichols C subclades, but
263 the maximum likelihood support values were below 0.95; therefore the phylogenetic
264 relationship between TPVMW082H and Nichols B and C could not be clearly delineated. To
265 demonstrate the fact that strain TPVMW082H is only distantly related to other samples
266 included in this phylogeny, we have assigned this strain its own subclade, Nichols X.
267
268 **Nanopore sequencing can resolve *tp0470* and *arp* repeat sequences in a high throughput**
269 **fashion**

270   To date, no systematic analysis of full-length *T. pallidum* tandem repeat genes has been
271   performed, in part due to the fact that the short-read Illumina data used to generate most
272   whole genomes cannot resolve the many 24- and 60-bp repeats found in *tp0470* and *arp*,
273   respectively. Although Sanger sequencing can be employed to examine repeats from both
274   genes, this approach is labor intensive and relies on the *a priori* assumption that the tandem
275   repeats comprise fewer than ~800 bp, the length limit for Sanger sequencing; lengths of up to
276   22 repeats (1320 bases) have been reported for *arp* (Harper et al., 2008). For *tp0470*, the DAL1
277   strain has 29 repeats of 24 bases (696 bases total) (Cejková et al., 2012), and the *T. pallidum*
278   subsp. *pertenue* strain CDC-2 has 37 24 base repeats (888 bases total) with no upper limit
279   known. Therefore, we developed bench and bioinformatic protocols for highly multiplexed,
280   long read sequencing of the *arp* and *tp0470* loci. We amplified portions of the *arp* and *tp0470*
281   genes using dual-indexed primers, allowing up to 192 amplicons on a single nanopore Flongle,
282   followed by demultiplexing, reference mapping, and consensus calling. Across six Flongle runs
283   to generate *arp* and *tp0470* data on 212 strains, an average of 408,286 reads passing filter were
284   generated. Following demultiplexing of forward and reverse barcodes, an average of 3327
285   reads were assigned to each sample, comprising both *arp* and *tp0470* reads. Reads were then
286   mapped to *tp0470* and *arp* reference files containing between 1-60 and 1-24 tandem repeats,
287   respectively, using bwa mem (Li and Durbin, 2009) with the Oxford Nanopore preset (-ont2d) to
288   account for low fidelity reads.
289
290   To validate our method, we cross referenced the band size seen on gel electrophoresis with the
291   automated call from our pipeline. S1 and S2 Figs show the gel electrophoresis bands from a
292   subset of samples and histograms of the distribution of mapping to the number of repeats for
293   *tp0470* and *arp*, respectively. We found 89% (183/206) concordance between the *tp0470*
294   pipeline call and band size, with discordance likely due to considerable low molecular weight
295   byproducts produced during amplification of *tp0470* repeats, which are 75% GC-rich. However,
296   as is clear for sample China11, the correct repeat number was readily apparent in these
297   samples upon inspection of the mapping distribution (boxed region and arrows, S1B Fig).
298   Technical replicates of select samples agreed with each other in 100% of cases, however,
299   amplicons produced using a two-stage amplification (see Methods) required more manual
300   comparison with band size to eliminate the low molecular weight byproducts, which were
301   unsurprisingly more likely to appear during two stage barcoding. Replicates generated using
302   both methods gave the same result (Supporting Information, "Methods Equivalency").
303
304   At 98.5% (200/203), concordance between the number of *arp* repeats determined
305   automatically by our pipeline and the electrophoresis band size was extremely high (S2 Fig).
306   Sanger sequencing was used to validate *tp0470* repeat lengths of select samples, with 100%
307   concordance seen. Confirmation of novel *arp* repeat types, as well as linkage between adjacent
308   60 bp repeats, was performed by manual inspection of Illumina WGS reads for all unique *arp*
309   variants. Technical replicates and replicates using each barcoding method gave the same result.
310
311   **The number of *tp0470* repeats is higher in samples from the Nichols-like clade**
312   We first examined the consensus sequence and number of tandem 24 bp repeats in the *tp0470*
313   gene in the context of the recombination-masked whole genome phylogeny we had previously

314    determined (Lieberman et al., 2021). Unlike *arp*, which has some sequence variation at six
315    positions per 60 bp repeat (see below; (Liu et al., 2007; Harper et al., 2008)), we did not find
316    any *tp0470* sequence changes in any of the 233 strains we examined, which included 211
317    strains we successfully sequenced by nanopore and the remainder from public databases.
318    However, we did note wide variability in the number of repeats (Fig 1A), ranging between 4-51,
319    with Nichols-like strains generally having more repeats than SS14-like strains (mean (median)
320    repeats for Nichols and SS14 25.7 (27) and 15.3 (14), respectively; $p<2.2 \times 10^{-16}$, Welch's t-test)
321    (Fig 1B), although this effect was difficult to disentangle from the strong location-specific
322    effects we observed (Fig 1C, $p<2 \times 10^{-16}$, ANOVA).
323    
324    **The number of *arp* repeats varies per subclade**
325    We then examined the number of *arp* repeats in each sample in the whole genome phylogeny.
326    Among the 226 TPA strains with *arp* sequence information, including 203 strains we sequenced
327    and the remainder from public databases, the majority (170, 83.7%) had 14 tandem repeats, 28
328    (13.8%) had ten repeats, and the remainder had between 4-24 repeats (Figure 2A, Supporting
329    Information), consistent with previous reports of the prevalence of repeat lengths (Pillay et al.,
330    1998; Harper et al., 2008; Marra et al., 2010). No relationship between *tp0470* repeat length
331    and *arp* repeat length was observed (SFig3A, Pearson coefficient = 0.008), as seen in prior
332    observations (Šmajs et al., 2018), even when strains containing 14 repeat *arp* sequences were
333    removed (Pearson coefficient = 0.196). With the exception of the Nichols-like clade Laboratory
334    Strain "Chicago" (NC_017268), which was excluded due to a known sequencing error in the
335    reference sequence, all variants were in frame and SNVs limited to the six positions per 60 bp
336    repeat previously recognized to be variable, resulting in four amino acid substitutions.
337    Subclade-specific repeat length variation was noted ($p < 2 \times 10^{-16}$, ANOVA), with all but five
338    samples in Nichols subclade B having ten repeats, three of the four strains in Nichols C with 19
339    repeats, and 15 or 16 repeats in the SS14 Mexico subclade (Figure 2A-B). When samples from
340    Madagascar, which had a bimodal distribution of *arp* repeats, were removed from analysis, no
341    repeat length variation was found among different countries (Figure 2C, $p > 0.05$, ANOVA).
342    Among the 170 samples with 14 ARP repeats, nine different gene sequences, which we have
343    named A-I in order of decreasing prevalence, were represented (Figure 2D), with different
344    usage patterns in different subclades. For example, ARP14 variant A was found in 103 samples
345    in SS14 subclades exclusively. ARP14 variants B and C, which differ by only three nucleotides
346    within a single 60 bp repeat, are found exclusively within the Nichols A subclade, variant E
347    found exclusively in Nichols B, and variant D found in Nichols subclades C, E, and the subclade
348    containing the laboratory strains, Nichols D.
349    
350    **The sequence of *arp* repeats varies per subclade**
351    We also characterized the pattern of the modular 60 bp near-identical repeats. Three types of
352    *arp* repeat (Type I, II, and III) were originally identified by Liu et al (Liu et al., 2007). A fourth
353    type, Type II/III, which likely formed by recombination of Types II and III between the two sets
354    of three variable positions to form a chimera, was discovered in a larger analysis of laboratory
355    and clinical strains (Harper et al., 2008). In addition to these "canonical" types of *arp* repeat, we
356    found three additional repeat Types that had not been previously described (Fig 3A): Type I/III,
357    which appears to be a chimera of Types I and III and found only in a single Peruvian strain with

358  seven *arp* repeats; Type III/I, which is a chimera of Type III and either Type I or II and found in a
359  14-repeat *arp* variant found in 17 samples from Madagascar as well as the Cuban strain CW83
360  (Grillová et al., 2019); and Type IIIG, found in a single sample from the United States, which
361  shares the Type III sequence at the first four variable positions and likely recombination to
362  match Types I or II at the final two variable positions. All repeat Types generated unique amino
363  acid sequences (Fig. 3B). Consistent with previous reports that non-venereal *T. pallidum*
364  subspecies use only Type II repeats (Harper et al., 2008; Cejková et al., 2012; Staudová et al.,
365  2014), the single Lihir Island *T. pallidum* subsp. *pertenue* strain and four Japan *T. pallidum*
366  subsp. *endemicum* strains included in our previous study (Lieberman et al., 2021) had only four
367  Type II repeats, or ten or eleven Type II repeats, respectively (Supporting Information).
368
369  To better visualize the relationships and pattern of repeat Type use between different *arp*
370  variants, we plotted them in the context of the whole genome phylogeny (Fig. 3C). In this
371  context, several general patterns emerge. Most strikingly, out of 121 strains in the SS14-like
372  clade with ARP sequence information, 115 include least two Type II repeats, a penultimate Type
373  II/III repeat, and a 3' Type III repeat, in contrast to the Nichols-like clade, which instead contains
374  four or more Type II repeats followed by one (Nichols subclade A) or two (Nichols subclades B-
375  E) Type III repeats. Furthermore, with the exception of *arp* sequences in the Nichols B subclade,
376  which start with a Type I repeat followed by Type II, *arp* variants in 99% of strains across both
377  major clades start with two Type I repeats.
378
379  Within each subclade, we defined the dominant *arp* sequence as the one most commonly
380  found. Out of the 226 *T. pallidum* subsp. *pallidum* strains with *arp* sequence information, we
381  identified 170 strains with the dominant *arp* sequence in the subclade to which it belongs, 25
382  strains where the *arp* repeat sequence clearly did not match the dominant sequence in its
383  subclade, and 26 strains that comprised three clusters of closely-related strains all exhibiting
384  the same *arp* variant but diverged from the predominant variant in the subclade to which they
385  belonged. ("Dominant *arp* Variant in Subclade", Supporting Information). The dominant strain
386  in the SS14 Mexico or Nichols X subclades could not be determined, and Nichols D Lab Strain
387  Chicago was excluded due to a known sequencing error. We did not find any SNVs enriched
388  among strains with an *arp* sequence altered from the dominant strain (p>0.05, Fisher's exact
389  test). We also examined the relationship between non-dominant *arp* sequence and branch
390  length on the whole genome maximum likelihood tree. We found that the strains with the non-
391  dominant *arp* sequence had terminal branch lengths (defined as the number of SNPs that
392  separate a tip from its most recent ancestral node) that were on average 3.6 times longer than
393  those with the dominant sequence (SFig4, mean (median) 4.11 (3) SNPs vs 1.14 (0) SNPs,
394  *p*=0.0046, Welch's t-test). While these observations do not account for sources of selective
395  pressure such as host immune response or anatomic site of infection, they do suggest that
396  events that result in a novel *arp* sequence are likely stochastic and less frequent than SNP
397  fixation, which occurs approximately once per genome every five years (Beale et al., 2019;
398  Lieberman et al., 2021; Taouk et al., 2022).
399
400  **Novel *arp* sequences likely arise through intra- or inter-strain recombination**

401    In cases where a strain's *arp* sequence did not match the dominant sequence in its subclade,
402    we confirmed the sequence by examination of WGS short read linkage and/or Sanger
403    sequencing. We then attempted to determine the simplest mechanism to generate the novel
404    sequence. In most cases, a single intra-strain recombination event, resulting in insertion,
405    deletion, or substitution of one or more *arp* modules, is most parsimonious (SFig5; most
406    variants could theoretically be generated by addition or removal of modules from the dominant
407    sequence in a different pattern than shown). For example, in the SS14 Omega´ and East Asia
408    subclades, ARP14.A is most commonly modified presumptively through loss of repeat modules,
409    resulting in sequence lengths between 7 and 12 modules. Notably, though, we do not know
410    how long the recombination junctions must be, therefore sequences such as ARP11.B, which
411    could have been formed by a recombination event between the 4th and 5th variable nucleotides
412    of type III and type II modules as shown, or could also have been generated via deletion of the
413    three tandem type II repeats in conjunction with mutation of the 5th and 6th variable positions
414    in the type III repeat from A to G, to generate the novel Type IIIG module.

416    In addition to novel *arp* variants that can be generated via a single intrastrain recombination,
417    there are other variants whose presumptive lineage is less clear. For example, the strain
418    UAB46xei is very unusual, both starting and ending with type II modules, unlike any other *arp*
419    variants. Interestingly, though, the UAB46xei *arp* sequence is 10 repeats, like the ARP10.A
420    sequence that predominates the Nichols B subclade. Strains TPVMW082H and Dublin57B,
421    belonging to Nichols X and SS14 Omega´ subclades, respectively, also contain unique sequences
422    generated via complex mechanisms, though it is plausible that several individual recombination
423    events generated the 20-repeat variant found in TPVMW082H, which is only distantly related to
424    any other strains in this dataset.

426    Other strains show clear evidence of inter-strain recombination. For example, ARP10.A is the
427    dominant sequence in the Madagascar strains that comprise Nichols subclade B and found in
428    no other strains except for the single Madagascar sample in the SS14 East Asia subclade.
429    Similarly, although the dominant *arp* variant in the SS14 Mexico subclade cannot be
430    determined since each strain has a unique sequence, the *arp* variant in SS14 Mexico strain
431    MD06B is only shared with the SS14 Omega´ subclade strain MD51x, both of which were
432    collected from Maryland, USA. Finally, strain Japan317x in Nichols subclade C harbors the same
433    ARP14.D variant as is found in strains in Nichols subclades D and E, including in one Japan
434    sample; however, deletion of 5 modules from the ARP19.A variant private to the Nichols C
435    subclade is also a possible mechanism for generation of the 14.D variant in Japan317x.
436    Together, these data suggest that both intra- and inter-strain recombination is employed by *T.*
437    *pallidum* to generate diversity at the *arp* locus.

439    We also attempted to determine if *tp0470* repeat length or *arp* repeat length and sequence
440    were associated with syphilis stage. Among the 79 strains with stage information available, 49
441    were primary and 30 were secondary; although longer *tp0470* variants were seen on average in
442    secondary syphilis samples (SFig6A, mean (median) 23.6 (25) for secondary vs 18.8 (15) for
443    primary, *p*=0.03624, Welch's t-test), no significant differences in *tp0470* repeat length by
444    disease stage were observed when samples were further split by SS14- or Nichols-like clades,

445    suggesting that sampling bias may be confounding interpretation of the association of *tp0470*
446    repeat length with disease stage. There was no association between non-dominant *arp*
447    sequence and primary or secondary syphilis (p>0.05, Fisher's exact test), nor did we find a
448    significant difference in the number of *arp* repeats among primary vs secondary syphilis
449    samples (Sfig6B, *p*=0.358, Welch's t-test). However, secondary syphilis was overrepresented
450    among strains with the ARP10.A sequence (Fisher's Exact Test, *p*=0.0171), while primary syphilis
451    was overrepresented among strains with the ARP14.A sequence (Fisher's Exact Test, *p*=0.0188).
452    No other *arp* sequence variants had enough samples with stage data to determine
453    overrepresentation.
454
455    **Structural modeling of ARP and TP0470 to localize 3D repeat structure**
456    There is ample evidence the proteins encoded by *tp0470* and *arp* genes are present during
457    infection: Previous studies have shown that the *tp0470* transcript is expressed (Smajs et al.,
458    2005; De Lay et al., 2021), and sera from infected rabbits (McKevitt et al., 2005) and patients
459    (Brinkman et al., 2006) are reactive to TP0470 protein. The *arp* transcript is expressed (Smajs et
460    al., 2005; De Lay et al., 2021), and the ARP protein was found to be one of the top 10% most
461    abundant proteins by mass spectrometry (Osbak et al., 2016). *T. pallidum*-infected rabbit sera
462    are reactive to ARP (McKevitt et al., 2005; Liu et al., 2007), while sera from infected human
463    patients are weakly reactive to ARP during primary infection (Brinkman et al., 2006). Therefore,
464    we attemped to model select variants of the full-length proteins encoded by *tp0470* and *arp*.
465    The *tp0470* gene is identical in all strains included in this study outside of the repeat length
466    variation, and TP0470 is confidently predicted by both SignalP 6.0 (Teufel et al., 2022) and
467    PSORTb V3.0 (Yu et al., 2010) to contain a signal sequence with no lipid anchor, which suggests
468    it resides in the periplasm. A conserved domain search (Lu et al., 2020) reveals the
469    tetratricopeptide repeat protein domain (e-value: $1.52e^{-8}$) at the N-term of the protein, with no
470    predicted conserved domains otherwise. This is consistent in structures predicted by both
471    trRosetta (Yang et al., 2020; Du et al., 2021) and AlphaFold (Jumper et al., 2021; Varadi et al.,
472    2022), which contain four pairs of antiparallel $\alpha$-helices that comprise the conserved
473    tetratricopeptide motif in the N terminus, followed by an extended $\alpha$-helix largely composed of
474    the highly charged eight amino acid repeat motif "EAEEARRK" (Fig 4A-B; tetratricopeptide
475    repeat motifs in green, pre-repeat linker in grey, 15 repeats of eight amino acids in purple, post-
476    repeat C terminus in gold). Confidence metrics for trRosetta are high for the overall structure
477    (TM-score = 0.756), while for AlphaFold the local Difference Distance Test score is >80 (high)
478    throughout the tetratricopeptide repeat domain, and drops throughout the length of the
479    extended helix. Variants with longer tandem repeats are predicted to have a helix that folds
480    back on itself by AlphaFold, while trRosetta predicts an elongated structure; although at low
481    confidence (SFig7A-B). The length of the repeat portion of the helix ranges between 4-51
482    repeats, or a total of 32-408 residues in repeats, with a modal number of repeats of 15 (SFig7C).
483    Assuming 0.54nm in length per helical turn of 3.6 residues, the elongated length of the
484    predicted helix may range between 11.5 nm and 68.1 nm including the non-repetitive 23 amino
485    acids N terminal and 23 amino acids C terminal to the repeats, with 90% of lengths between 7-
486    36 repeats (15.3-50.1 nm), and the helix of the 15 repeat variant measuring approximately 24.9
487    nm. Within the helix, a single eight-residue repeat makes just over two helical turns. Fig 4C
488    zooms in on four repeats, comprising approximately 9 helical turns. Although the orientation of

489    amino acid residues in a structural model does not reflect the precise native conformation, the
490    stick representation of sidechains (Fig 4C, top) and smoothed surface charge (Fig 4C, bottom)
491    demonstrates the highly polar nature of the TP0470 repeats.

492

493    We predicted domains and structures for select variants of the acidic repeat protein. By SignalP
494    6.0, it is predicted to have a signal sequence (probability=0.67) and possibly lipidation site at
495    cysteine-29 (probability=0.33), however, PSORT predicts neither of these elements. A
496    conserved domain search reveals a C-terminal SPOR domain (e value: $7.5e^{-3}$), which in other
497    proteins is a peptidoglycan binding domain (Yahashiri et al., 2017). Together, these results
498    suggest that the acidic repeat protein is localized to the periplasm.

499

500    For model generation, we first examined the ARP14.A variant, by far the most common variant
501    in our phylogeny, harbored by 103 strains. While both trRosetta and AlphaFold predicted the
502    expected twisted beta strand structure of the C-terminal SPOR domain (Fig 5A-B, magenta),
503    AlphaFold's low confidence prediction of the repeat regions (local Distance Difference Test ~40)
504    is entirely unstructured (Fig 5B), whereas the trRosetta prediction is for the acidic repeats to
505    form a disordered linker comprising the first five repeats, followed by a parallel β-sheet
506    structure that contains nine strands composed of the last nine acidic repeats (Fig 5A), although
507    the confidence in the prediction is quite low (TM-score=0.288). This structure would contain an
508    extremely acidic face of the β-sheet (Fig 5C) with a periodicity of 20 residues, the same as the
509    repeat. Interestingly, the repeat region in most other variants was not predicted by trRosetta to
510    fold into a beta sheet, rather, they were highly disordered (SFig8); only variants ARP14.H and
511    ARP15.A were also predicted to form a β-sheet from the repeats. Overall, despite a plausible
512    structure for some variants, structural modeling of the acidic repeat protein remains
513    challenging with only low confidence models returned by two methods and any attempt to
514    infer function based on these results should be made cautiously.

515

516    **Discussion**
517    Despite the relatively low rate of SNP fixation, with a mean rate of approximately $1-3x10^{-7}$
518    substitutions per site per year in putative non-recombinogenic loci (Beale et al., 2019;
519    Lieberman et al., 2021; Taouk et al., 2022)), *T. pallidum* uses additional mechanisms to increase
520    its genetic diversity and antigenic repertoire. These include inter-strain and inter-species
521    recombination in genes encoding the Tpr family of antigens (Gray et al., 2006; Kumar et al.,
522    2018; Grillová et al., 2019), gene conversion in the variable regions of *tprK* (Centurion-Lara et
523    al., 2004; Giacani et al., 2010; Reid et al., 2014), and homopolymer expansion and contraction
524    to alter promoter activity and hence expression level of putative outer membrane proteins
525    (Giacani et al., 2015). Previous work has demonstrated that diversity in the *tp0470* and *arp*
526    repeat length, and repeat type usage in *arp*, is likely generated by recombination, although
527    modification of the number of repeats in the *tp0470* could also be possible via a polymerase
528    slippage mechanism. Our current study extends these findings to a large cohort of clinical
529    samples with near-complete genomes available, enabling examination of differences between
530    subclades and correlation with genome features.
531        From our results, it is clear that a very wide distribution of *tp0470* repeat lengths is
532    possible but with no sequence variation within the repeat. Although *T. pallidum* has an

533  extremely low rate of mutation and most genes are highly conserved, the absence of sequence
534  variation within the *tp0470* gene outside of repeat length variation suggests the protein may be
535  under purifying selection. The *arp* gene has multiple sequence variants generated by using
536  different repeat module types in a tandem arrangement, but is highly enriched for sequences
537  with fourteen repeats. For both genes, some differences between Nichols- and SS14-like clades
538  are observed: in *tp0470*, there is a slight increase in repeat length in the Nichols-like clade vs
539  SS14. In the *arp* gene, variants are generally limited to a single subclade, particularly in the
540  Nichols-like clade, which has far greater genetic diversity than the SS14-like clade, with an
541  average pairwise SNP distance of 42, as compared to an average pairwise SNP distance of 10 for
542  the SS14-like clade. It is unclear whether differences between *tp0470* lengths or repeat module
543  pattern in *arp* between subclades have functional consequences and are being selected for, or
544  whether the differences simply reflect random events during diversification. In the case of *arp*,
545  we did not find any SNPs throughout the genome that correlated with an unexpected repeat
546  sequence.
547      Until very recently (Romeis et al., 2021), no reverse genetics system for *T. pallidum*
548  existed, therefore, traditional bacteriological genetic tools, such as mutants and knockout
549  strains, to interrogate gene functions have not been available for the syphilis spirochete. Prior
550  to the development of an epithelial cell co-culture system in 2018 (Edmondson et al., 2018), *T.
551  pallidum* could only be passaged through rabbit testes, precluding forward genetics screens.
552  While proteome-wide bioinformatic structural predictions have helped to shed light on the
553  likely role of conserved structural domains (Houston et al., 2018), the structure and function *T.
554  pallidum* proteins containing novel motifs, such as the repeat sequences found in TP0470 and
555  ARP, remain unknown. To gain insight into their biological function and possible role in syphilis
556  pathogenesis, we employed several *in silico* tools to predict the topology and structure of full-
557  length ARP and TP0470 proteins.
558      The presence of a signal peptide on TP0470 is strong evidence that it is localized to the
559  periplasm, where it likely binds other as-yet-undermined protein(s) via its N-terminal
560  tetratricopeptide motif, which was predicted by both trRosetta and AlphaFold. Both algorithms
561  predicted an extended $\alpha$-helix with regions of alternating positive and negative surface
562  electrostatic potential, regardless of the length of the repeats. Although TP0470 does not have
563  any known interacting partners, it seems plausible that in addition to interactions formed by
564  tetratricopeptide motif, the unusually long, very polar $\alpha$-helix that comprises the repeats also
565  serves to mediate protein-protein interactions.
566      Although a signal peptide was not confidently predicted by one tool (SignalP 6.0) and
567  not predicted at all by a second (PSORT), the presence of a C-terminal SPOR domain, which
568  binds denuded peptidoglycan, strongly suggests ARP must be present in the periplasm.
569  However, it is less clear whether ARP is free in the periplasm or is acylated at cys-29 (weakly
570  predicted by SignalP 6.0), tethering it to the inner membrane. Both topologies are consistent
571  with other SPOR domain-containing proteins (Yahashiri et al., 2017); mass spectrometry, Edman
572  degradation, or other biochemical techniques will be necessary to resolve this question.
573      In addition to the unclear localization of ARP, the structure formed by the ARP repeat
574  modules remains murky. The most biologically plausible structure generated by the modeling
575  software is of the modular repeats forming a parallel beta sheet, with periodicity of 20 residues,
576  the same length as the repeats. The $\beta$-sheet and the loops that connect the strands form an

577 extremely negatively charged surface; it seems likely that whatever the ARP repeat domain
578 binds, it will be positively charged.
579     Figure 6 summarizes our current understanding of the structure and topology of TP0470
580 and ARP. Because TP0470 is predicted to be soluble, it may be able to traverse the
581 peptidoglycan layer through pores. In contrast, ARP may be associated with the inner
582 membrane, constraining its movement within the periplasm. However, these models do not
583 resolve why there is tremendous diversity of *tp0470* repeat lengths, and *arp* repeat lengths and
584 repeat module usage. To answer these questions and determine the biological function of the
585 repeat domains, extensive biochemical and biophysical studies of different variants will be
586 necessary.
587     One of the primary limitations of our genomic dataset is that of bias introduced by
588 unequal sampling. For example, although we and others have reported that 14-repeat arp
589 variants predominate, particularly variant 14.A, this may reflect the increased sampling in
590 geographical regions (Europe and the United States) where SS14-like clade strains predominate.
591 Furthermore, many of the strains included in our study were not collected with extensive
592 clinical histories or patient characteristics, limiting our ability to infer functional differences
593 from sequence variation. Finally, the use of PCR to interrogate genes containing repetitive
594 sequences is always challenging due to the generation of truncated artifact products, which can
595 amplify preferentially over the "real" product, as we saw for *tp0470* in some samples (SFig1B).
596 However, these products were readily apparent in both the gel images and histograms of
597 mapped reads; therefore the data remained interpretable, and were consistent in technical
598 replicates (Supporting information).
599     In addition to developing a novel method to examine repeat length and sequence in two
600 challenging genomic loci in *T. pallidum*, our study has demonstrated extensive *tp0470* and *arp*
601 repeat diversity among more than 200 clinical strains with whole genome sequence, by far the
602 largest study of these genes to date. Importantly, we found that more than 10% of strains
603 contained an *arp* variant that had a different length and sequence than the dominant variant in
604 the subclade, which builds on concerns about the utility of using the number of *arp* repeats as
605 part of strain typing tools for epidemiology (Mikalová et al., 2013). Finally, we have proposed a
606 possible mechanism by which each may interact with peptidoglycan and/or other periplasmic
607 factors and influence morphogenesis. Although additional genetic, biophysical, and biochemical
608 interaction studies will be necessary to characterize their function and elucidate their binding
609 partners, our study of *tp0470* and *arp* lays the foundation to directly link the genotype to
610 function of two novel genes that may influence *T. pallidum* pathogenesis.
611
612 **Data Availability**
613 Sequencing data is available under NCBI BioProjects PRJNA723099 and PRJNA815321.
614 Supporting Information contains BioSample and accession information for nanopore reads.
615
616 **Funding**
617 This work was supported, in whole or in part, by the Bill & Melinda Gates Foundation INV-
618 036560. Under the grant conditions of the Foundation, a Creative Commons Attribution 4.0
619 Generic License has already been assigned to the Author Accepted Manuscript version

### References

629    Arora, N., Schuenemann, V. J., Jäger, G., Peltzer, A., Seitz, A., Herbig, A., et al. (2017). Origin of
630        modern syphilis and emergence of a pandemic Treponema pallidum cluster. *Nat
631        Microbiol* 2, 16245. doi: 10.1038/nmicrobiol.2016.245.

632    Beale, M. A., Marks, M., Cole, M. J., Lee, M.-K., Pitt, R., Ruis, C., et al. (2021). Global phylogeny
633        of Treponema pallidum lineages reveals recent expansion and spread of contemporary
634        syphilis. *Nat Microbiol* 6, 1549–1560. doi: 10.1038/s41564-021-01000-z.

635    Beale, M. A., Marks, M., Sahi, S. K., Tantalo, L. C., Nori, A. V., French, P., et al. (2019). Genomic
636        epidemiology of syphilis reveals independent emergence of macrolide resistance across
637        multiple circulating lineages. *Nat Commun* 10, 3255. doi: 10.1038/s41467-019-11216-7.

638    Brinkman, M. B., McKevitt, M., McLoughlin, M., Perez, C., Howell, J., Weinstock, G. M., et al.
639        (2006). Reactivity of antibodies from syphilis patients to a protein array representing the
640        Treponema pallidum proteome. *J Clin Microbiol* 44, 888–891. doi:
641        10.1128/JCM.44.3.888-891.2006.

642    Bushnell, B. BBMap short read aligner, and other bioinformatic tools.

643    Cejková, D., Zobaníková, M., Chen, L., Pospíšilová, P., Strouhal, M., Qin, X., et al. (2012). Whole
644        genome sequences of three Treponema pallidum ssp. pertenue strains: yaws and
645        syphilis treponemes differ in less than 0.2% of the genome sequence. *PLoS Negl Trop Dis*
646        6, e1471. doi: 10.1371/journal.pntd.0001471.

647    Centurion-Lara, A., LaFond, R. E., Hevner, K., Godornes, C., Molini, B. J., Van Voorhis, W. C., et
648        al. (2004). Gene conversion: a mechanism for generation of heterogeneity in the tprK
649        gene of Treponema pallidum during infection. *Mol Microbiol* 52, 1579–1596. doi:
650        10.1111/j.1365-2958.2004.04086.x.

651    Cerveny, L., Straskova, A., Dankova, V., Hartlova, A., Ceckova, M., Staud, F., et al. (2013).
652        Tetratricopeptide repeat motifs in the world of bacterial pathogens: role in virulence
653        mechanisms. *Infect Immun* 81, 629–635. doi: 10.1128/IAI.01035-12.

654    Chen, W., Šmajs, D., Hu, Y., Ke, W., Pospíšilová, P., Hawley, K. L., et al. (2021). Analysis of
655        Treponema pallidum Strains From China Using Improved Methods for Whole-Genome

656  Sequencing From Primary Syphilis Chancres. *J Infect Dis* 223, 848–853. doi:
657  10.1093/infdis/jiaa449.

658  Currin, A., Swainston, N., Dunstan, M. S., Jervis, A. J., Mulherin, P., Robinson, C. J., et al. (2019).
659  Highly multiplexed, fast and accurate nanopore sequencing for verification of synthetic
660  DNA constructs and sequence libraries. *Synthetic Biology* 4, ysz025. doi:
661  10.1093/synbio/ysz025.

662  De Lay, B. D., Cameron, T. A., De Lay, N. R., Norris, S. J., and Edmondson, D. G. (2021).
663  Comparison of transcriptional profiles of Treponema pallidum during experimental
664  infection of rabbits and in vitro culture: Highly similar, yet different. *PLoS Pathog* 17,
665  e1009949. doi: 10.1371/journal.ppat.1009949.

666  Du, Z., Su, H., Wang, W., Ye, L., Wei, H., Peng, Z., et al. (2021). The trRosetta server for fast and
667  accurate protein structure prediction. *Nat Protoc* 16, 5634–5651. doi: 10.1038/s41596-
668  021-00628-9.

669  Edmondson, D. G., Hu, B., and Norris, S. J. (2018). Long-Term In Vitro Culture of the Syphilis
670  Spirochete Treponema pallidum subsp. pallidum. *mBio* 9. doi: 10.1128/mBio.01153-18.

671  Flasarová, M., Pospíšilová, P., Mikalová, L., Vališová, Z., Dastychová, E., Strnadel, R., et al.
672  (2012). Sequencing-based molecular typing of treponema pallidum strains in the Czech
673  Republic: all identified genotypes are related to the sequence of the SS14 strain. *Acta*
674  *Derm Venereol* 92, 669–674. doi: 10.2340/00015555-1335.

675  Giacani, L., Brandt, S. L., Ke, W., Reid, T. B., Molini, B. J., Iverson-Cabral, S., et al. (2015).
676  Transcription of TP0126, Treponema pallidum putative OmpW homolog, is regulated by
677  the length of a homopolymeric guanosine repeat. *Infect Immun* 83, 2275–2289. doi:
678  10.1128/IAI.00360-15.

679  Giacani, L., Molini, B. J., Kim, E. Y., Godornes, B. C., Leader, B. T., Tantalo, L. C., et al. (2010).
680  Antigenic variation in Treponema pallidum: TprK sequence diversity accumulates in
681  response to immune pressure during experimental syphilis. *J Immunol* 184, 3822–3829.
682  doi: 10.4049/jimmunol.0902788.

683  Gray, R. R., Mulligan, C. J., Molini, B. J., Sun, E. S., Giacani, L., Godornes, C., et al. (2006).
684  Molecular evolution of the tprC, D, I, K, G, and J genes in the pathogenic genus
685  Treponema. *Mol Biol Evol* 23, 2220–2233. doi: 10.1093/molbev/msl092.

686  Grillová, L., Oppelt, J., Mikalová, L., Nováková, M., Giacani, L., Niesnerová, A., et al. (2019).
687  Directly Sequenced Genomes of Contemporary Strains of Syphilis Reveal
688  Recombination-Driven Diversity in Genes Encoding Predicted Surface-Exposed Antigens.
689  *Front Microbiol* 10, 1691. doi: 10.3389/fmicb.2019.01691.

690  Grimes, M., Sahi, S. K., Godornes, B. C., Tantalo, L. C., Roberts, N., Bostick, D., et al. (2012). Two
691  mutations associated with macrolide resistance in Treponema pallidum: increasing

692    prevalence and correlation with molecular strain type in Seattle, Washington. *Sex*
693        *Transm Dis* 39, 954–958. doi: 10.1097/OLQ.0b013e31826ae7a8.

694    Harper, K. N., Liu, H., Ocampo, P. S., Steiner, B. M., Martin, A., Levert, K., et al. (2008). The
695        sequence of the acidic repeat protein (arp) gene differentiates venereal from
696        nonvenereal Treponema pallidum subspecies, and the gene has evolved under strong
697        positive selection in the subspecies that causes syphilis. *FEMS Immunol Med Microbiol*
698        53, 322–332. doi: 10.1111/j.1574-695X.2008.00427.x.

699    Hook III, E. W., Behets, F., Van Damme, K., Ravelomanana, N., Leone, P., Sena, A. C., et al.
700        (2010). A Phase III Equivalence Trial of Azithromycin versus Benzathine Penicillin for
701        Treatment of Early Syphilis. *J INFECT DIS* 201, 1729–1735. doi: 10.1086/652239.

702    Hopkins, S., Lyons, F., Coleman, C., Courtney, G., Bergin, C., and Mulcahy, F. (2004). Resurgence
703        in Infectious Syphilis in Ireland: An Epidemiological Study. *Sexually Transmitted Diseases*
704        31, 317–321. doi: 10.1097/01.OLQ.0000123653.84940.59.

705    Houston, S., Lithgow, K. V., Osbak, K. K., Kenyon, C. R., and Cameron, C. E. (2018). Functional
706        insights from proteome-wide structural modeling of Treponema pallidum subspecies
707        pallidum, the causative agent of syphilis. *BMC Struct Biol* 18, 7. doi: 10.1186/s12900-
708        018-0086-3.

709    Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., et al. (2021). Highly
710        accurate protein structure prediction with AlphaFold. *Nature* 596, 583–589. doi:
711        10.1038/s41586-021-03819-2.

712    Kumar, S., Caimano, M. J., Anand, A., Dey, A., Hawley, K. L., LeDoyt, M. E., et al. (2018).
713        Sequence Variation of Rare Outer Membrane Protein β-Barrel Domains in Clinical Strains
714        Provides Insights into the Evolution of Treponema pallidum subsp. pallidum, the Syphilis
715        Spirochete. *mBio* 9. doi: 10.1128/mBio.01006-18.

716    Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
717        transform. *Bioinformatics* 25, 1754–1760. doi: 10.1093/bioinformatics/btp324.

718    Lieberman, N. A. P., Lin, M. J., Xie, H., Shrestha, L., Nguyen, T., Huang, M.-L., et al. (2021).
719        Treponema pallidum genome sequencing from six continents reveals variability in
720        vaccine candidate genes and dominance of Nichols clade strains in Madagascar. *PLoS*
721        *Negl Trop Dis* 15, e0010063. doi: 10.1371/journal.pntd.0010063.

722    Liu, D., He, S.-M., Zhu, X.-Z., Liu, L.-L., Lin, L.-R., Niu, J.-J., et al. (2020). Molecular
723        Characterization Based on MLST and ECDC Typing Schemes and Antibiotic Resistance
724        Analyses of Treponema pallidum subsp. pallidum in Xiamen, China. *Front Cell Infect*
725        *Microbiol* 10, 618747. doi: 10.3389/fcimb.2020.618747.

726    Liu, H., Rodes, B., George, R., and Steiner, B. (2007). Molecular characterization and analysis of
727            a gene encoding the acidic repeat protein (Arp) of Treponema pallidum. *J Med Microbiol*
728            56, 715–721. doi: 10.1099/jmm.0.46943-0.

729    Lu, S., Wang, J., Chitsaz, F., Derbyshire, M. K., Geer, R. C., Gonzales, N. R., et al. (2020).
730            CDD/SPARCLE: the conserved domain database in 2020. *Nucleic Acids Res* 48, D265–
731            D268. doi: 10.1093/nar/gkz991.

732    Lukehart, S. A., Godornes, C., Molini, B. J., Sonnett, P., Hopkins, S., Mulcahy, F., et al. (2004).
733            Macrolide resistance in Treponema pallidum in the United States and Ireland. *N Engl J*
734            *Med* 351, 154–158. doi: 10.1056/NEJMoa040216.

735    Marra, C. M., Sahi, S. K., Tantalo, L. C., Godornes, C., Reid, T., Behets, F., et al. (2010). Enhanced
736            molecular typing of treponema pallidum: geographical distribution of strain types and
737            association with neurosyphilis. *J Infect Dis* 202, 1380–1388. doi: 10.1086/656533.

738    McKevitt, M., Brinkman, M. B., McLoughlin, M., Perez, C., Howell, J. K., Weinstock, G. M., et al.
739            (2005). Genome scale identification of Treponema pallidum antigens. *Infect Immun* 73,
740            4445–4450. doi: 10.1128/IAI.73.7.4445-4450.2005.

741    Mikalová, L., Pospíšilová, P., Woznicová, V., Kuklová, I., Zákoucká, H., and Smajs, D. (2013).
742            Comparison of CDC and sequence-based molecular typing of syphilis treponemes: tpr
743            and arp loci are variable in multiple samples from the same patient. *BMC Microbiol* 13,
744            178. doi: 10.1186/1471-2180-13-178.

745    Molepo, J., Pillay, A., Weber, B., Morse, S. A., and Hoosen, A. A. (2007). Molecular typing of
746            Treponema pallidum strains from patients with neurosyphilis in Pretoria, South Africa.
747            *Sex Transm Infect* 83, 189–192. doi: 10.1136/sti.2006.023895.

748    Noda, A. A., Méndez, M., Rodríguez, I., and Šmajs, D. (2022). Genetic Recombination in
749            Treponema pallidum: Implications for Diagnosis, Epidemiology, and Vaccine
750            Development. *Sex Transm Dis* 49, e7–e10. doi: 10.1097/OLQ.0000000000001497.

751    Osbak, K. K., Houston, S., Lithgow, K. V., Meehan, C. J., Strouhal, M., Šmajs, D., et al. (2016).
752            Characterizing the Syphilis-Causing Treponema pallidum ssp. pallidum Proteome Using
753            Complementary Mass Spectrometry. *PLoS Negl Trop Dis* 10, e0004988. doi:
754            10.1371/journal.pntd.0004988.

755    Parr, J. B. (2021). IDEELResearch/tpallidum_genomics: Tpallidum_genomics. doi:
756            10.5281/ZENODO.5773174.

757    Pětrošová, H., Zobaníková, M., Čejková, D., Mikalová, L., Pospíšilová, P., Strouhal, M., et al.
758            (2012). Whole genome sequence of Treponema pallidum ssp. pallidum, strain Mexico A,
759            suggests recombination between yaws and syphilis strains. *PLoS Negl Trop Dis* 6, e1832.
760            doi: 10.1371/journal.pntd.0001832.

761  Pillay, A., Liu, H., Chen, C. Y., Holloway, B., Sturm, A. W., Steiner, B., et al. (1998). Molecular
762      subtyping of Treponema pallidum subspecies pallidum. *Sex Transm Dis* 25, 408–414. doi:
763      10.1097/00007435-199809000-00004.

764  Pillay, A., Liu, H., Ebrahim, S., Chen, C. Y., Lai, W., Fehler, G., et al. (2002). Molecular typing of
765      Treponema pallidum in South Africa: cross-sectional studies. *J Clin Microbiol* 40, 256–
766      258. doi: 10.1128/JCM.40.1.256-258.2002.

767  Pinto, M., Borges, V., Antelo, M., Pinheiro, M., Nunes, A., Azevedo, J., et al. (2016). Genome-
768      scale analysis of the non-cultivable Treponema pallidum reveals extensive within-
769      patient genetic variation. *Nat Microbiol* 2, 16190. doi: 10.1038/nmicrobiol.2016.190.

770  Pope, V., Fox, K., Liu, H., Marfin, A. A., Leone, P., Seña, A. C., et al. (2005). Molecular subtyping
771      of Treponema pallidum from North and South Carolina. *J Clin Microbiol* 43, 3743–3746.
772      doi: 10.1128/JCM.43.8.3743-3746.2005.

773  Reid, T. B., Molini, B. J., Fernandez, M. C., and Lukehart, S. A. (2014). Antigenic variation of TprK
774      facilitates development of secondary syphilis. *Infect Immun* 82, 4959–4967. doi:
775      10.1128/IAI.02236-14.

776  Romeis, E., Tantalo, L., Lieberman, N., Phung, Q., Greninger, A., and Giacani, L. (2021). Genetic
777      engineering of Treponema pallidum subsp. pallidum, the Syphilis Spirochete. *PLoS
778      Pathog* 17, e1009612. doi: 10.1371/journal.ppat.1009612.

779  Smajs, D., McKevitt, M., Howell, J. K., Norris, S. J., Cai, W.-W., Palzkill, T., et al. (2005).
780      Transcriptome of Treponema pallidum: gene expression profile during experimental
781      rabbit infection. *J Bacteriol* 187, 1866–1874. doi: 10.1128/JB.187.5.1866-1874.2005.

782  Šmajs, D., Strouhal, M., and Knauf, S. (2018). Genetics of human and animal uncultivable
783      treponemal pathogens. *Infect Genet Evol* 61, 92–107. doi:
784      10.1016/j.meegid.2018.03.015.

785  Staudová, B., Strouhal, M., Zobaníková, M., Cejková, D., Fulton, L. L., Chen, L., et al. (2014).
786      Whole genome sequence of the Treponema pallidum subsp. endemicum strain Bosnia
787      A: the genome is related to yaws treponemes but contains few loci similar to syphilis
788      treponemes. *PLoS Negl Trop Dis* 8, e3261. doi: 10.1371/journal.pntd.0003261.

789  Sutton, M. Y., Liu, H., Steiner, B., Pillay, A., Mickey, T., Finelli, L., et al. (2001). Molecular
790      subtyping of Treponema pallidum in an Arizona County with increasing syphilis
791      morbidity: use of specimens from ulcers and blood. *J Infect Dis* 183, 1601–1606. doi:
792      10.1086/320698.

793  Taouk, M. L., Taiaroa, G., Pasricha, S., Herman, S., Chow, E. P. F., Azzatto, F., et al. (2022).
794      Characterisation of Treponema pallidum lineages within the contemporary syphilis
795      outbreak in Australia: a genomic epidemiological analysis. *Lancet Microbe* 3, e417–e426.
796      doi: 10.1016/S2666-5247(22)00035-0.

797  Teufel, F., Almagro Armenteros, J. J., Johansen, A. R., Gíslason, M. H., Pihl, S. I., Tsirigos, K. D., et
798       al. (2022). SignalP 6.0 predicts all five types of signal peptides using protein language
799       models. *Nat Biotechnol*. doi: 10.1038/s41587-021-01156-3.

800  Thurlow, C. M., Joseph, S. J., Ganova-Raeva, L., Katz, S. S., Pereira, L., Chen, C., et al. (2022).
801       Selective Whole-Genome Amplification as a Tool to Enrich Specimens with Low
802       Treponema pallidum Genomic DNA Copies for Whole-Genome Sequencing. *mSphere*,
803       e0000922. doi: 10.1128/msphere.00009-22.

804  Van Damme, K., Behets, F., Ravelomanana, N., Godornes, C., Khan, M., Randrianasolo, B., et al.
805       (2009). Evaluation of Azithromycin Resistance in Treponema pallidum Specimens From
806       Madagascar. *Sexually Transmitted Diseases* 36, 775–776. doi:
807       10.1097/OLQ.0b013e3181bd11dd.

808  Varadi, M., Anyango, S., Deshpande, M., Nair, S., Natassia, C., Yordanova, G., et al. (2022).
809       AlphaFold Protein Structure Database: massively expanding the structural coverage of
810       protein-sequence space with high-accuracy models. *Nucleic Acids Res* 50, D439–D444.
811       doi: 10.1093/nar/gkab1061.

812  Wang, L.-G., Lam, T. T.-Y., Xu, S., Dai, Z., Zhou, L., Feng, T., et al. (2020). Treeio: An R Package for
813       Phylogenetic Tree Input and Output with Richly Annotated and Associated Data.
814       *Molecular Biology and Evolution* 37, 599–603. doi: 10.1093/molbev/msz240.

815  Wick, R. R., Judd, L. M., Gorrie, C. L., and Holt, K. E. (2017). Completing bacterial genome
816       assemblies with multiplex MinION sequencing. *Microbial Genomics* 3. doi:
817       10.1099/mgen.0.000132.

818  Wick, R. R., Judd, L. M., and Holt, K. E. (2019). Performance of neural network basecalling tools
819       for Oxford Nanopore sequencing. *Genome Biol* 20, 129. doi: 10.1186/s13059-019-1727-
820       y.

821  Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. 2nd ed. 2016. Cham: Springer
822       International Publishing : Imprint: Springer doi: 10.1007/978-3-319-24277-4.

823  Yahashiri, A., Jorgenson, M. A., and Weiss, D. S. (2017). The SPOR Domain, a Widely Conserved
824       Peptidoglycan Binding Domain That Targets Proteins to the Site of Cell Division. *J
825       Bacteriol* 199, e00118-17. doi: 10.1128/JB.00118-17.

826  Yang, J., Anishchenko, I., Park, H., Peng, Z., Ovchinnikov, S., and Baker, D. (2020). Improved
827       protein structure prediction using predicted interresidue orientations. *Proc Natl Acad Sci
828       U S A* 117, 1496–1503. doi: 10.1073/pnas.1914677117.

829  Yu, G., Smith, D. K., Zhu, H., Guan, Y., and Lam, T. T. (2017). ggtree: an R package for
830       visualization and annotation of phylogenetic trees with their covariates and other
831       associated data. *Methods Ecol Evol* 8, 28–36. doi: 10.1111/2041-210X.12628.

832  Yu, N. Y., Wagner, J. R., Laird, M. R., Melli, G., Rey, S., Lo, R., et al. (2010). PSORTb 3.0: improved
833       protein subcellular localization prediction with refined localization subcategories and
834       predictive capabilities for all prokaryotes. *Bioinformatics* 26, 1608–1615. doi:
835       10.1093/bioinformatics/btq249.

836  Zhou, L., Feng, T., Xu, S., Gao, F., Lam, T. T., Wang, Q., et al. (2022). ggmsa: a visual exploration
837       tool for multiple sequence alignment and associated data. *Briefings in Bioinformatics*,
838       bbac222. doi: 10.1093/bib/bbac222.

839  Zobaníková, M., Strouhal, M., Mikalová, L., Cejková, D., Ambrožová, L., Pospíšilová, P., et al.
840       (2013). Whole genome sequence of the Treponema Fribourg-Blanc: unspecified simian
841       isolate is highly similar to the yaws subspecies. *PLoS Negl Trop Dis* 7, e2172. doi:
842       10.1371/journal.pntd.0002172.

843
844  **Figure Legends**
845
846  **Figure 1: Variation in *tp0470* repeat length**. A) Recombination masked whole genome
847  phylogeny (left) with the number of *tp0470* repeats for each strain (right). Sequence variant
848  number is included as text to the right of each length bar. All data are also included in a tabular
849  representation in Supporting Info. Number of *tp0470* repeats by subclade (B) or country (C).
850
851  **Figure 2: Variation in *arp* repeat length**. A) Recombination masked whole genome phylogeny
852  (left) with the number of *arp* repeats for each strain (right). Sequence variant number is
853  included as text to the right of each length bar. All data are also included in a tabular
854  representation in Supplementary Table 2. Number of *arp* repeats by subclade (B) or country (C).
855  D) Multiple sequence alignment of the nine variants with 14 *arp* repeats. Variant positions are
856  highlighted and bases colored red, blue, yellow, or green for A, C, G, or T, respectively. The
857  number of strains with each variant sequence is included in the bar graph to the right of the
858  multiple sequence alignment.
859
860  **Figure 3: *arp* repeat type usage**. A) Nucleotide sequence of the arp repeat module types.
861  Variable positions are highlighted. B) Amino acid sequence of the arp repeat module types.
862  Variable positions are highlighted. C) Recombination masked whole genome phylogeny (left)
863  with the repeat type usage per strain (right). The 60 bp arp repeats are colored by type.
864
865  **Figure 4: Structure predictions of TP0470**. A) trRosetta and B) AlphaFold predictions of
866  structure of 15 repeat TP0470 variant. N terminal tetratricopeptide repeat domain is shown in
867  green, repeats are in purple, and C-terminal region is in gold. C) APBS electrostatic surface
868  potential (top) and stick model of sidechains (bottom) for portion of repeat helix.
869
870  **Figure 5: Structure predictions ARP14A**. A) trRosetta and B) AlphaFold predictions of structure
871  of ARP14A. C-terminal SPOR domain is shown in magenta, repeats in cyan. C) APBS electrostatic

872 surface potential for trRosetta ARP structure from A. Red denotes negative charge (acidic) and
873 blue denotes positive charge (basic).

874

875 **Figure 6: Model showing ARP and TP0470 cellular location and putative interactions.** Both
876 ARP and TP0470 are localized to the periplasm. The ARP N terminus may be acylated at cysteine
877 29. OM: Outer Membrane. IM: Inner Membrane. PG: Peptidoglycan.

878

879

880 **Figure S1: *tp0470* PCR band and nanopore pipeline call concordance of select samples**. A) PCR
881 bands following barcoding. Bands should have a size of 24 bp times the number of repeats plus
882 273, including both flanking regions (225 bp) and two barcodes (48 bp total). The red box
883 highlights non-specific amplification of low molecular weight fragments, while the red arrow
884 shows the correct band size. B) Histogram of the distribution of reads aligned to each length
885 variant in the mapping reference file.

886

887 **Figure S2: *arp* PCR band and nanopore pipeline call concordance of select samples**. A) PCR
888 bands following barcoding. Bands should have a size of 60 bp times the number of repeats plus
889 874, including both flanking regions (826 bp) and two barcodes (48 bp total). B) Histogram of
890 the distribution of reads aligned to each length variant in the mapping reference file.

891

892 **Figure S3: *arp* repeat length is not correlated with *tp0470* repeats.** No correlation was seen
893 between number of *arp* and *tp0470* repeats (Pearson coefficient = 0.008).

894

895 **Figure S4: Terminal branch lengths are longer for strains with non-dominant arp sequences.**
896 Average branch lengths from tip to ancestral node were determined for strains with dominant
897 and non-dominant arp variants. **, $p$=0.0046, Welch's t-test.

898

899 **Figure S5: Possible single recombination events to generate arp variants from the dominant
900 arp variant in each subclade.** The dominant sequence from each subclade is shown in the top
901 position for each pair, and the non-dominant (possibly nascent) variant below. Dotted lines
902 show possible junctions**.** Only a single possibility per variant is shown.

903

904 **Figure S6: Longer *tp0470* repeats are associated with secondary syphilis.** Out of 79 samples
905 with stage information, samples from secondary syphilis had on average approximately 5 more
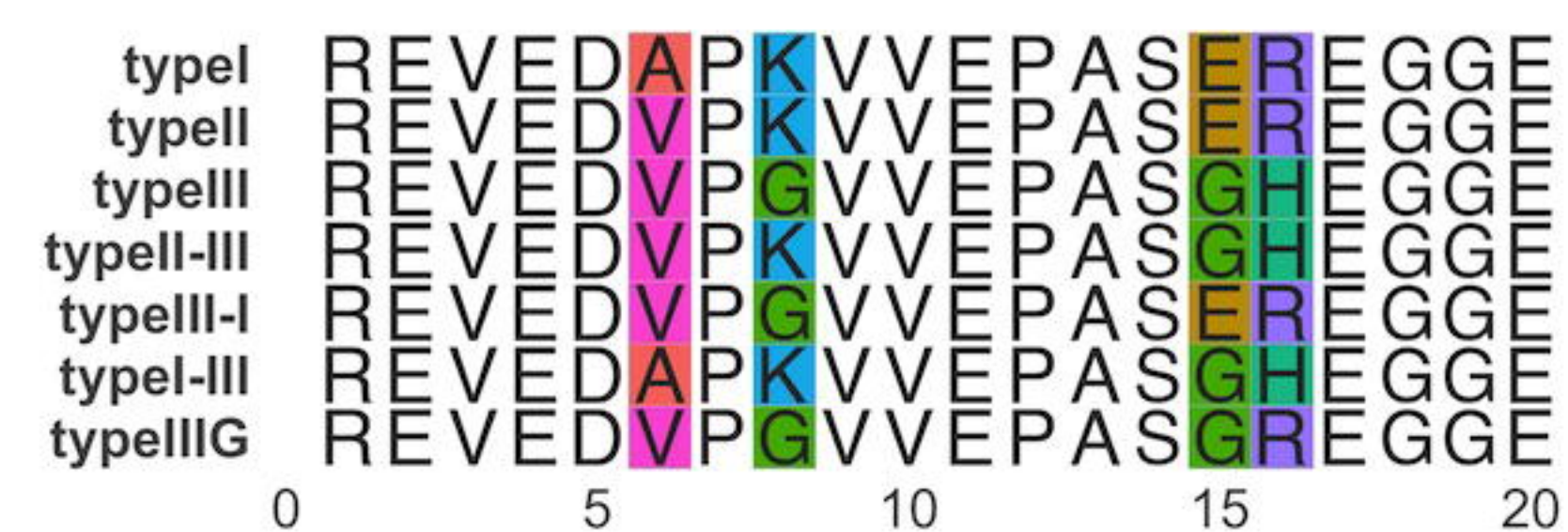906 repeats than primary (*, $p$=0.03624, Welch's t-test).

907

908 **Figure S7: TP0470 predictions**: A) trRosetta and B) AlphaFold predictions of structure of 51
909 repeat TP0470 variant. Structures are colored from blue to red N term to C term. C) Distribution
910 of *tp0470* variants in phylogeny.

911

912 **Figure S8: trRosetta predictions for additional ARP variants.** Structures are colored from blue
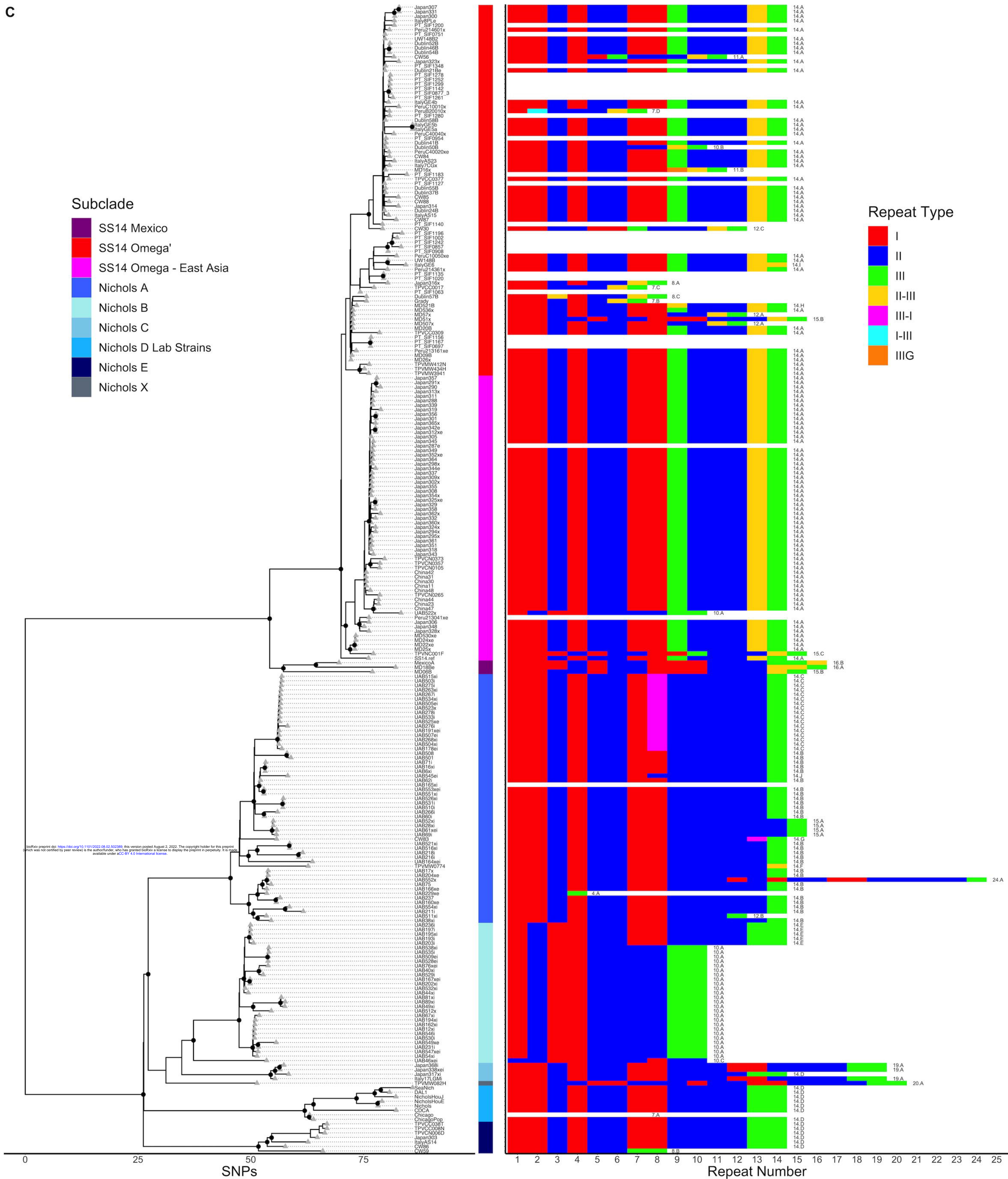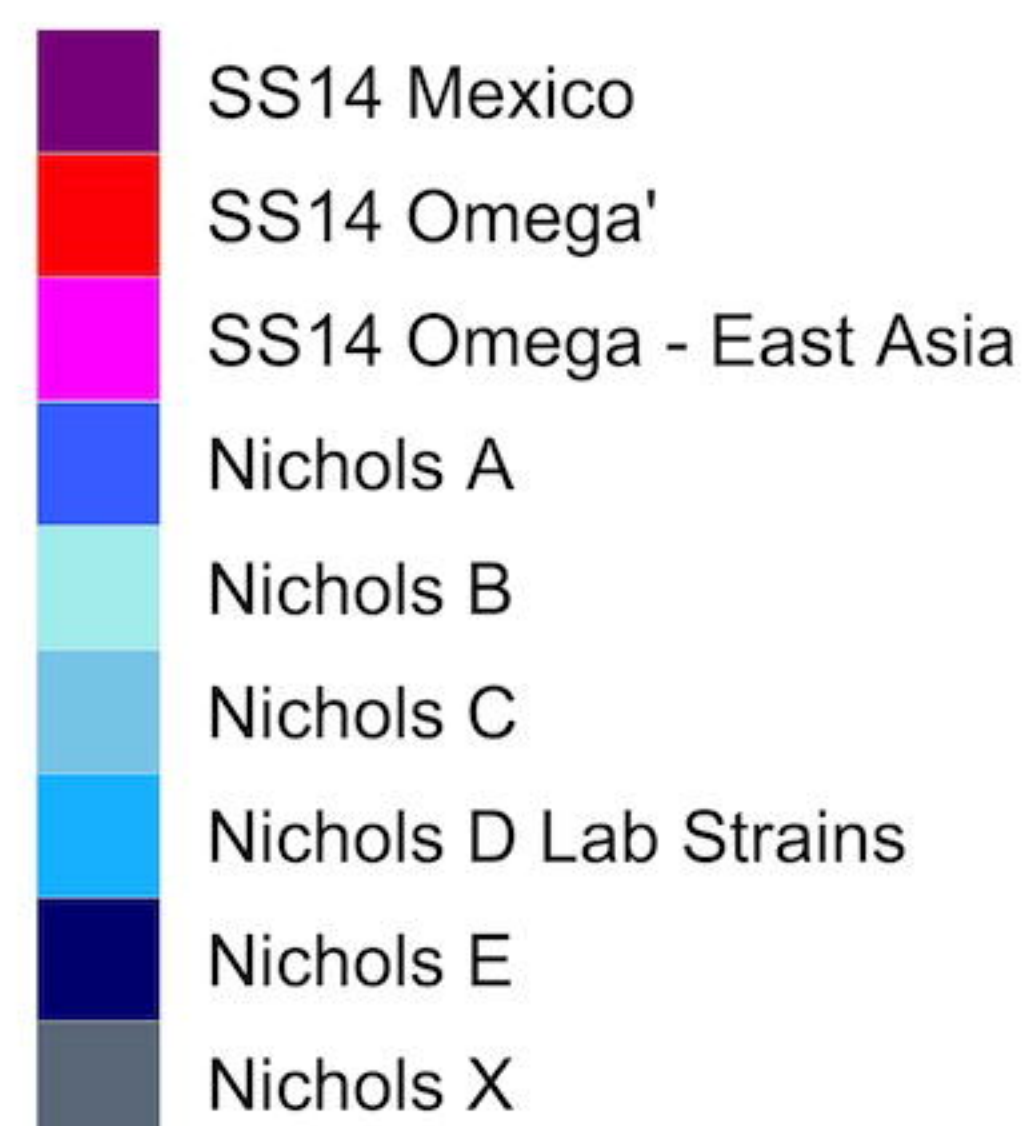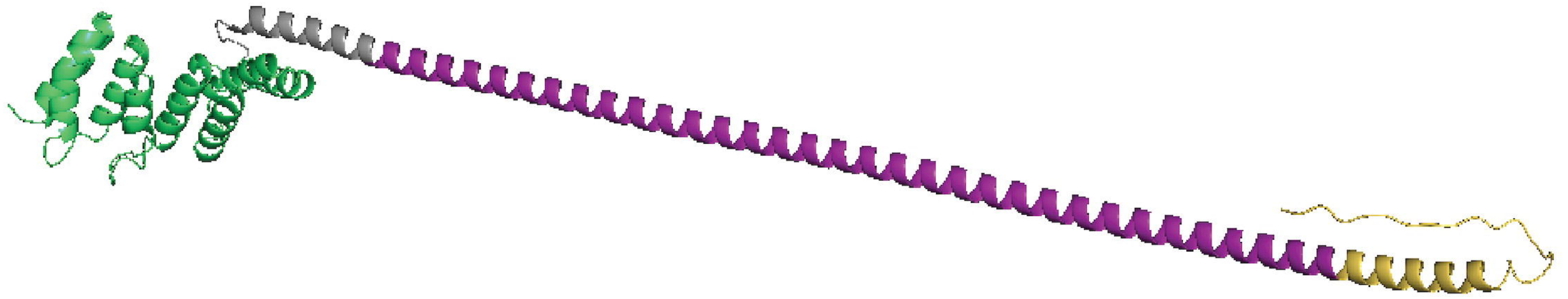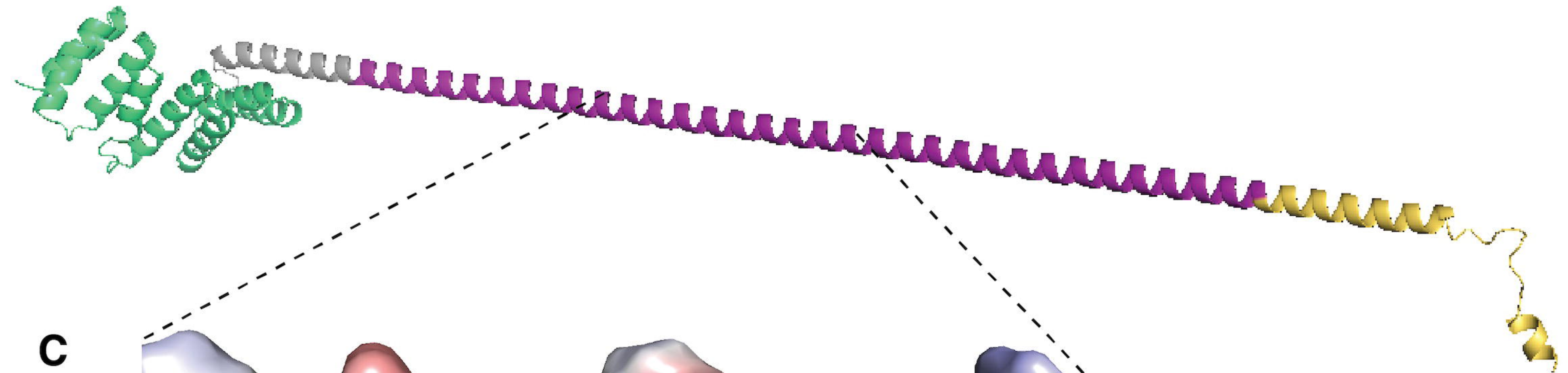913 to red N term to C term.

914

Surface Potential (KT/e)

Rotate 180 →
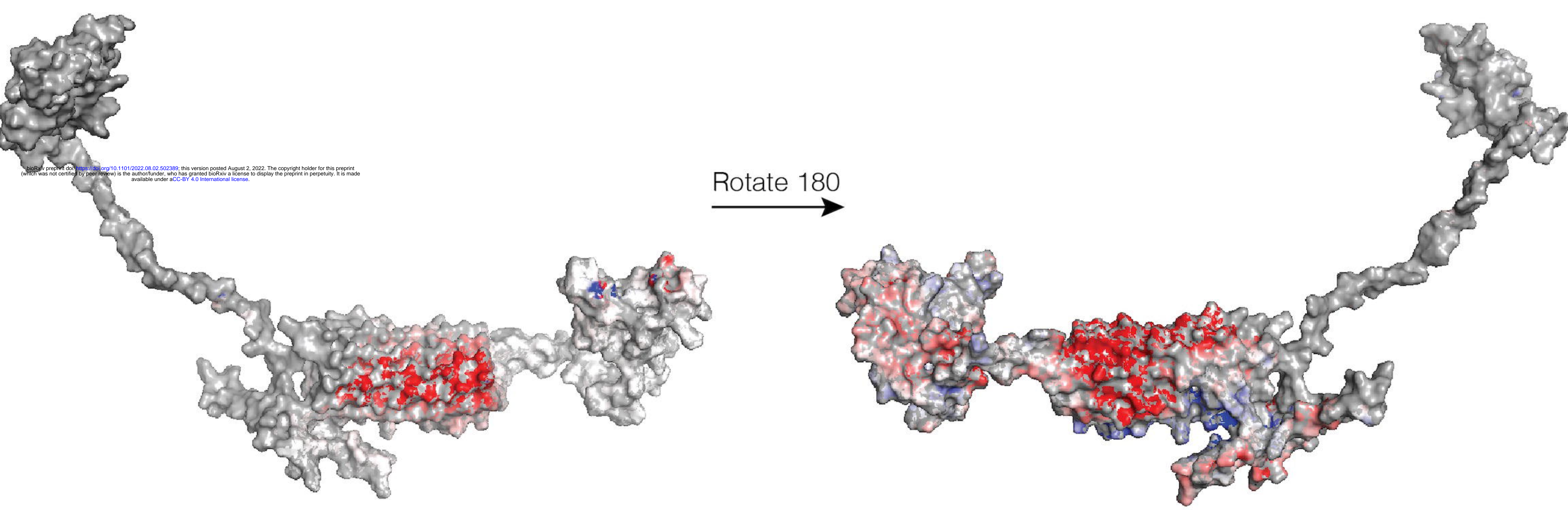
Surface Potential (KT/e)

−5.000          5.000

OM

denuded PG

crosslinked PG

TPR domain

SPOR domain

ARP

TPR domain

TP0470

Cys

IM