1 **Assembly and phylogeographic analysis of novel *Taenia solium* mitochondrial**
2 **genomes reveal further differentiation between and within Asian and African-**
3 **American genotypes**
4
5 Gabriel Jiménez-Avalos[1]§, Alina Soto Obando[1]§, Maria Solis[1], Robert H Gilman[2], Vitaliano Cama[3],
6 Armando E Gonzalez[4], Hector H García[5,6], Patricia Sheen[1], David Requena[7*], and Mirko Zimic[1*], for the
7 Cysticercosis Working Group in Peru.

8 § These authors contributed equally.

9 [1]Laboratorio de Bioinformática, Biología Molecular y Desarrollos Tecnológicos. Laboratorios
10 de Investigación y Desarrollo.
11 Facultad de Ciencias y Filosofía.
12 Universidad Peruana Cayetano Heredia.
13 Av. Honorio Delgado 430.
14 San Martín de Porres, 15102.
15 Lima, Perú.
16
17 [2]Department of International Health.
18 Bloomberg School of Public Health.
19 Johns Hopkins University.
20 615 North Wolfe St., Room 5515.
21 Baltimore, 21205.
22 Maryland, USA.
23
24 [3]Division of Parasitic Diseases and Malaria.
25 Center for Global Health.
26 Centers for Disease Control and Prevention.
27 1600 Clifton Rd. MS D-65.
28 Atlanta, GA 30329.
29 The USA.
30
31 [4]Facultad de Medicina Veterinaria.
32 Universidad Nacional Mayor de San Marcos.
33 Av Circunvalación 28
34 San Borja, 15021
35 Lima, Peru
36
37 [5]Departamento de Microbiología.
38 Universidad Peruana Cayetano Heredia
39 Av. Honorio Delgado 430.
40 San Martín de Porres, 15102.
41 Lima, Perú
42
43 [6]Cysticercosis Unit, Instituto Nacional de Ciencias Neurológicas
44 Jr. Ancash 1271.
45 Cercado de Lima 15003.
46 Lima, Perú
47
48 [7]Laboratory of Cellular Biophysics.
49 The Rockefeller University.
50 1230 York Avenue.
51 New York, NY 10065.
52 USA.
53
54 * corresponding authors
55 Laboratorio de Bioinformática y Biología Molecular.
56 Facultad de Ciencias y Filosofía.
57 Universidad Peruana Cayetano Heredia.

58    Av. Honorio Delgado 430.
59    San Martín de Porres, 15102.
60    Lima, Perú.
61    Phone: (511) 3190000 ext. 2604.
62
63    Laboratory of Cellular Biophysics.
64    The Rockefeller University.
65    1230 York Avenue.
66    New York, NY 10065.
67    USA
68

# Abstract

70

71

**Background**

73

74    *Taenia solium* is a parasite that hampers human health, causing taeniasis and cysticercosis. The
75    genetic variability in its mitochondrial genome is related to the geographical origin of the specimen. Two
76    main genotypes have been identified: The Asian and the African-American. The geographic genetic
77    variability is expected to cause different clinical manifestations. Thus, characterizing differences
78    between and within genotypes is crucial for completing the epidemiology of *T. solium* diseases.

79

**Methods/Principal Findings**

81

82    Here, two Peruvian (one complete and one partial; 7,811X and 42X of coverage, respectively) and one
83    Mexican (complete, 3,395X) *T. solium* mitochondrial genomes were assembled using the Chinese
84    reference. Variant calling with respect to the reference was performed. Thirteen SNPs that involved a
85    change in the amino acid physicochemical nature were identified. Those were present in all the
86    assembled genomes and might be linked to differences in aerobic respiration efficiency between Latin
87    American (African-American) and Asian genotypes. Then, phylogeographic studies were conducted
88    using Cytochrome C oxidase subunit I and cytochrome B from these genomes and other isolates. The
89    analysis showed that Indonesian samples are the most ancient and related to the modern *T. solium*
90    ancestor of the Asian genotype. Finally, a consistent subdivision of the African-American genotype into
91    two subgroups was found. One subgroup relates to East African countries, while the other is West
92    Africa. The East African linage suggests a previously unnoticed influence of the Indian Ocean trade in
93    the genetic structure of Latin America *T. solium*.

94

**Conclusions/Significance**

96

97    Overall, this study reports novel mitochondrial genomes valuable for further studies. New Latin
98    American SNPs were identified and suggest metabolic differences between parasites of the Asian and
99    African-American genotypes. Moreover, the phylogeographic analysis revealed differences within each
100   genotype that shed light on *T. solium's* historical spread. Overall, the results represent an important
101   step in completing *T. solium* genetic epidemiology.

102

103

# Author Summary

105

106   *Taenia solium* is a human parasite that causes taeniasis and cysticercosis. Eradicated from developed
107   countries, they are still a public health problem in developing nations. *T. solium* differences in the
108   mitochondrial genetic material depend on its geographical origin. This is expected to cause different
109   clinical manifestations. Despite the importance of genetics to the epidemiology of *T. solium* diseases,
110   few efforts have been made to assemble and compare their genomes. We aimed to help fill this
111   knowledge gap by assembling three mitochondrial genomes from Latin America and comparing them
112   to the Chinese reference. Additionally, two genes from the Latin American genomes and from other

113    isolates were employed to assess *T. solium* genetic distribution. We found thirteen mutations with
114    respect to the Chinese genome present in all Latin American samples, which involved a change in the
115    amino acid physicochemical nature. Those might be causing metabolic differences between Asian and
116    Latin American parasites that could change their affinity to specific human tissues. Moreover, we
117    determined that Indonesian samples are the most ancient and related to the modern *T. solium* ancestor.
118    Finally, we identified a previously unnoticed influence of East African countries in *T. solium* phylogeny,
119    with which our assembled genomes are closely related.

120

121    **KEYWORDS:** Phylogenetics, Phylogeography, Haplotypes, Taeniasis, Cysticercosis,
122    Genetics, Genomics, Genetic Epidemiology, Mitochondrial genome.

123

124    **INTRODUCTION**

125

126    *Taenia solium* is a parasite responsible for two critical diseases in humans: taeniasis
127    and cysticercosis. The former refers to infection with the adult stage of the parasite.
128    The latter is the infection with its larvae and represents a major risk to human health.
129    Cysticercosis could progress to the central nervous system causing
130    neurocysticercosis, the leading cause of acquired adult epilepsy in developing
131    countries [1]. Humans are the only known definitive host, harboring the adult tapeworm
132    and releasing infectious eggs to the environment [2–4].

133

134    *T. solium* has spread globally, being endemic and highly prevalent in Asia, Africa, and
135    Latin America [5]. Interestingly, it has been shown that *T. solium* intraspecies variability
136    in the mitochondrial genome is strongly related to the geographical origin of the
137    specimens [6–9]. Two main genotypes have been identified, the Asian and the African-
138    American [6]. This geographic genetic variability is expected to result in clinical
139    heterogeneity in *T. solium* diseases between regions [10]. Therefore, an exhaustive
140    study of it is crucial for completing the epidemiology of taeniasis and cysticercosis
141    [6,11]. Despite this, there is a lack of characterization of differences between whole *T.*
142    *solium* mitochondrial genomes from different genotypes due to the few assembled
143    genomes available.

144

145    Mitochondrial genes, specially Cytochrome C oxidase subunit 1 (COX1) and
146    cytochrome B (CYTB), have proved to be useful markers to assess intraspecies
147    variability and phylogeography of *T. solium* [6–9,12–14]. Both genes have low
148    variability [6]; however, CYTB is suggested to be slightly more variable than COX1 in
149    *T. solium* [6,15]. The extremely low variability of COX1 limits its use in intraspecies
150    studies of this parasite, being CYTB more suitable for this purpose [15].
151    Comprehensive analysis of multiple *T. solium* mitochondrial genomes from different
152    origins will confirm if CYTB is the absolute best sequence for this kind of analysis.

153

154    The geographic distribution of *T. solium*'s genotypes was shaped by human migrations
155    and trade [6,7,9,12–14,16]. For example, the similarity and gene flow between Latin
156    American, African, and Philippine *T. solium* populations resulted from the European
157    maritime trade routes of the XV and XIX centuries. Furthermore, the sympatric
158    coexistence of the Asian and African-American genotypes in Madagascar is explained
159    by two independent human groups that migrated to this island and introduced both
160    lineages [9,12]. In that sense, geographic genetic differences (or similarities) between
161    *T. solium* populations depend on the connection degree of those places' human
162    groups. These differences can be used to assess the impact of human migration and
163    trade on the spread of this parasite [9,12,14], which is essential to prevent its

164  dissemination. However, there is also a lack of research on differentiation within each
165  genotype, although previous work has suggested, for instance, that two sublineages
166  could exist within the African-American genotype [8]. Including newly assembled
167  sequences and the ones reported worldwide in phylogeographic studies could help fill
168  this knowledge gap.
169
170  Hence, the present study assembled and annotated the *T. solium* mitochondrial
171  genomes of two Peruvian and one Mexican isolate. Those and the Chinese reference
172  mitochondrial genome were compared to provide a detailed characterization of the
173  differences between Asian and African-American genotypes' genomes. Finally, the
174  COX1 and CYTB sequences from these isolates and from others reported worldwide
175  were included in a phylogeographic reconstruction to analyze differentiation within the
176  Asian and African-American genotypes. COX1 and CYTB were used as they are the
177  genes with the most *T. solium* sequences available from diverse geographical origins.
178
179  **MATERIALS AND METHODS**
180
181  **Assembly and annotation of the mitochondrial genomes**
182
183  There are currently three Latin American *T. solium* mitochondrial reads available in
184  the NCBI Sequence Read Archive. Two from Peru [17], and one from Mexico (NCBI
185  BioProject: PRJNA170813). Additionally, an assembled reference mitochondrial
186  genome from a Chinese isolate has been published [18].
187
188  Whole-genome sequencing reads of the three Latin American samples were collected
189  (Accession codes: SRR644531, SRR650708, and SRR524725). Each sample was
190  independently mapped, using as reference the *T. solium* reference mitochondrial
191  genome of the Chinese isolate (GenBank ID: NC_004022). The mapped reads were
192  subjected to quality control in FastQC v.0.11.9
193  (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/) and trimmed using a
194  quality threshold of 0.001 in CLC Genomics Workbench v. 21.05.5
195  (https://digitalinsights.qiagen.com/). These cleaned reads were re-mapped against the
196  reference for each sample. The average coverage of this re-mapping was computed
197  using the pileup script of BBtools v. 38.91 (http://sourceforge.net/projects/bbmap/).
198  The consensus sequences were extracted, inserting the ambiguity symbol "N" to
199  handle conflicts in low coverage regions and inserting the most frequent base (voting)
200  when conflicts occurred in high coverage sections.
201
202  In addition, the cleaned reads were *de-novo* assembled to detect genetic
203  rearrangements. We also used CLC Genomics Workbench to perform quality control
204  on each contig (identifying chimeric sequences, misassemblies, and artifacts),
205  measure the depth coverage and %GC content per contig, and calculate the N50 of
206  the assembly.
207
208  The consensus sequences of each sample were manually curated to correct
209  misassemblies, which resulted in the final assembled genomes being used in further
210  steps. The annotation of the final assembled genomes was performed using CLC Main
211  Workbench v.21.05.5 using the annotation of the Chinese mitochondrial genome as a
212  reference.
213

**Variability and selective pressure**

Variant calling of the assembled Peruvian and Mexican mitochondrial genomes was performed by whole-genome multiple alignment in the software Mauve v. 2.4.0 [19], using the Chinese *T. solium* mitochondrial genome as reference and default parameters. Single Nucleotide Polymorphisms (SNPs) were manually curated to rule out possible sequencing, alignment, or variant calling errors; and classified into three categories: synonym, non-synonym, and mutations in non-coding regions. This was performed using the software DNAsp v. 6.12.03 [20]. SNPs were graphically represented on a scaled circular map using Circos v. 0.69 [21].

The variability of the 12 protein-coding genes was evaluated as the level of sequence difference (D) [22] for each pairwise combination of the genomes from China, Puno, Huancayo, and Mexico. Each of these genes was re-aligned using the web implementation of the EMBOSS Needle algorithm [23]. D was computed as D = 1 - (M/L), where M is the number of invariant sites and L is the difference between the alignment length and ambiguous bases.

The Ka/Ks ratio of each protein-coding gene against the Chinese reference sequence was computed to evaluate the level of selection pressure and evolutionary adaptation of the *T. solium* mitochondrial genomes. For this purpose, the multiple genome alignment was split into 12 protein-coding gene alignments and submitted to the DNSsp software to calculate the ratio.

**Phylogenetic analysis**

COX1 and CYTB mitochondrial genes from the genomes assembled and from different complete gene sequences reported worldwide were employed to perform a phylogenetic reconstruction. This comprises a total of 45 *T. solium* COX1 and 31 CYTB sequences available in GenBank, including, as the outgroup, sequences from *T. saginata* (COX1*:* AB066495.1 and NC_009938.1; CYTB*:* AB066581.1 and NC_00938. 1), *T. asiatica* (COX1*:* AB066494.1 and NC_004826.2; CYTB*:* AB066580.1 and NC_004826.2) and *Echinococcus multilocularis* (COX1 and CYTB*:* NP_000928.2).

Multiple global alignments were performed independently for COX1 and CYTB in the software MAFFT v. 7 [24,25], using the progressive G-INS-1 method. The informative coding regions of the alignments were extracted using the Gblocks server v. 0.91 [26,27] with default options.

Phylogenetic analysis was conducted for COX1 and CYTB*,* separately, using the Maximum Likelihood (ML) and Bayesian Inference (BI) algorithms. For ML, RaxML v. 8.2.12 [28] was used with the GTRCAT model, and 1000 bootstrap replicates to estimate the robustness of the branches. BI was conducted in BEAST2 v. 2.6.1 implemented on the CRIPRES online server platform [29].

The evolutionary model was estimated with jModelTest2 v. 2.1.6 for COX1 (GTR+G with four gamma categories) and CYTB (GTR+I, I = 0.5720) [30], using the Akaike criterion correction. A basic coalescent model for demographic history and a relaxed molecular clock model (uncorrelated lognormal) during the Markov Chain Monte Carlo

264    (MCMC) process were employed. The substitution rates were set to 0.0225 and
265    0.0195 substitutions per site per million years for COX1 and CYTB, respectively,
266    according to the genetic distance computed for *T. saginata* and *T. asiatica,* which
267    diverged 1.245 [0.78, 1.71] million years ago (Myra) [9,16]. The analysis was run for
268    50 million generations, sampling every 5000 states and using a burn-in of 10% to
269    obtain an effective sample size (ESS) greater than 200. Lastly, a Maximum Clade
270    Credibility tree (MCC) in TreeAnotator v. 2.6.0 [31] was generated.
271
272    **Haplotype network**
273
274    Haplotypes were identified in DnaSP using as input the multiple alignments of COX1
275    and CYTB previously generated and considering the total number of mutations as
276    nucleotide substitutions. To prevent adding ambiguity to the network, ambiguous
277    nucleotides were not considered gaps. The haplotype networks were calculated by
278    Median Joining using the software Networks v.10 (fluxus-engineering.com) [32]. The
279    genetic differentiation (φst) between the African-American subclades 1 and 2 was
280    calculated using a haplotype distance matrix in Arlequin v. 3.5.2.2 [33].
281
282    **RESULTS**
283
284    **Genome assembly and annotation**
285
286    Reads from the Peruvian and Mexican isolates mapped against the Chinese reference
287    were extracted, trimmed, filtered by quality, and re-mapped, resulting in 1,317,941
288    reads from isolates from Puno (7,811X coverage), 5,561 from Huancayo (42X), and
289    674,666 from Mexico (3,395X). Genomes from Puno and Mexico were complete, while
290    the one from Huancayo was partial. The three consensus genome sequences were of
291    similar length (13700-13709 nucleotides). Complete mapping statistics are present in
292    Supplementary Table 1).
293
294    Additionally, the trimmed and filtered reads were *de-novo* assembled. The same
295    mitochondrial gene arrangement as the Chinese reference assembly was confirmed
296    (Figure 1, Supplementary Table 2). Interestingly, the size of the protein-coding genes
297    in Latin American samples was identical to the Chinese reference (Supplementary
298    Table 2). An exception occurred for CYTB in the isolate from Huancayo, which has a
299    missing codon corresponding to positions 872-874 of the Chinese reference.
300
301    **Variability and selective pressure analysis**
302
303    Variant calling, computation of the level of sequence difference (D), and selective
304    pressure analysis were employed to provide a detailed comparison of the assembled
305    genomes and the Chinese reference. SNP distribution detected in Latin American
306    mitochondrial genomes was almost identical (Figure 1). Of the 34 non-synonymous
307    SNPs, 31 were present in the three Latin American samples (Supplementary Table 3).
308    They were distributed in all protein-coding genes except for NADH-ubiquinone
309    oxidoreductase chain 1 (ND1) and NADH-ubiquinone oxidoreductase chain (ND3).
310    Notably, 13 of the 31 SNPs changed the amino acid physicochemical nature. Those
311    were located within COX1, cytochrome C oxidase subunit 2 (COX2), NADH-
312    ubiquinone oxidoreductase chain 4 (ND4), and NADH-ubiquinone oxidoreductase

313     chain 5 (ND5) (Figure 1, Supplementary Table 3). These change-in-nature SNPs
314     represented 83.33% of the total mutations detected within ND5.

315

316     The similarity in the SNPs distribution along Latin American samples is also supported
317     by the small contribution of Huancayo, Puno, and Mexico pairwise comparisons in the
318     accumulative D (Figure 2). This was less than 0.01 for COX1, COX2, CYTB, NADH-
319     ubiquinone oxidoreductase chain 2 (ND2), and ND5; and 0 for the other protein-coding
320     genes (Figure 2). Five genes showed relatively high D values: ATP synthase subunit
321     6 (ATP6), COX2, CYTB, ND4, and NADH-ubiquinone oxidoreductase chain 6 (ND6).
322     CYTB was the only one with nonzero values for all the pairwise comparisons.

323

324     All the protein-coding genes presented a Ka/Ks ratio of less than 1, indicating a
325     purifying selection. In particular, NADH-ubiquinone oxidoreductase chain 1 (ND1) and
326     NADH-ubiquinone oxidoreductase chain (ND3) seem to be subject to absolute
327     purification (Ka/Ks = 0) in all the samples evaluated. The three Latin American
328     samples had the same Ka/Ks values for ATP6, cytochrome C oxidase subunit 3
329     (COX3), ND4, NADH-ubiquinone oxidoreductase chain 4L (ND4L), and ND6. In
330     contrast, COX2 and ND2 had higher Ka/Ks in the samples from Puno and Huancayo,
331     respectively, while ND5 had a higher Ka/Ks in the Mexican sample (Figure 3).

332

333     **Phylogenetic analysis**

334

335     To identify further differentiation within the Asian and African-American genotypes, a
336     phylogenetic reconstruction using COX1 and CYTB was performed, including the
337     isolates of this study and others reported worldwide. The phylogenetic analysis
338     distinguished two major clades for both markers: Asian and African-American. While
339     both genes supported the Asian clade (COX1: PP=0.95; CYTB: BS=96%, PP=1.0),
340     the African-American clade was only supported by the COX1 marker (BS=98%,
341     PP=0.99). See Figure 4.

342

343     For COX1, the Asian group was further subdivided into four subclades with high
344     support (Figure 4a). The first comprised countries from the East (China and Japan),
345     South (India and Nepal), and Southeast Asia (Thailand), along with the island of
346     Madagascar (Asian subclade 1; PP: 0.96). The second included identical sequences
347     from Nepal and China (Asian subclade 2; BS: 78). The third comprised sequences
348     exclusively from China (Asian subclade 3; PP: 1.0). A fourth group contained two
349     sequences from Indonesia (Asian subclade 4; BS: 78, PP: 1.0). In the CYTB-based
350     tree (Figure 4b), groups similar to the Asian subclades 2 (BS: 94, PP: 1.0) and 4 (BS:
351     99, PP: 1.0) were also present with high support. In addition, a group formed by just
352     Indian samples was present in the CYTB phylogeny (BS: 89), which might be
353     analogous to Asian subclade 1, as Indian sequences are only present in this subclade.

354

355     In the African-American genotype, two subclades appeared within the COX1
356     phylogeny. The first (African-American subclade 1; BS: 71) consisted of samples from
357     Tanzania and Mexico (Yucatan, Mexico State 1 and 2). The second (African-American
358     subclade 2; BS: 84 and PP: 1.0) comprised samples from Mexico (Mexico State 3 and
359     the Mexican sequence assembled in this study), Cameroon (West and North), Peru
360     (Puno and Huancayo), and Brazil. CYTB presented a similar topology for the African-
361     American genotype to that obtained with COX1. The group between all the countries

362  of the African-American subclade 1 is supported but not including Tanzania (BS = 96,
363  PP = 0.99).

364

365  **Divergence time**

366

367  Divergence times and their 95% highest posterior density intervals were calculated to
368  situate the differentiation events within the time scale. Divergence of Asian and
369  African-American genotypes occurred 0.458 [0.0405, 1.0625] Myra for COX1 or 0.634
370  [0.0667, 1.4574] Myra for CYTB.

371

372  Asian subclade 4 (Indonesia) diverged from the rest of Asia around 0.2450 [0.025,
373  0.5683] Myra for COX1 (Figure 4a) or 0.2496 [0.0284, 0.5804] Myra for CYTB (Figure
374  4b). Asian subclade 1 differentiated from Asian subclade 2 and 3 0.1565 [0.0161,
375  0.3526] Myra or 0.1613 [0.0200, 0.3826] Myra, for COX1 and CYTB, respectively.

376

377  The earliest divergence within Asian subclade 1 occurred at 0.0991 [0.0087, 0.2259]
378  Myra according to COX1 or 0.0909 [0.0051, 0.229] Myra according to CYTB. For
379  COX1, this divergence formed the common ancestor of a Chinese and Thailandese
380  sample. Within Asian subclade 2, differentiation of Chinese samples from Nepalese
381  (COX1) or Nepalese and Vietnamese (CYTB) occurred 0.0759 [0.0058, 0.1796] Myra
382  for COX1 or 0.0664 [0.0035, 0.1738] Myra for CYTB. Asian subclade 3 was only
383  present in the COX1 phylogeny. Within it, one Chinese sample diverged 0.0786 Myra.

384

385  According to COX1, the African-American subclade 1 diverged from subclade 2 about
386  0.1425 [0.0105, 0.3504] Myra. The clade that included both African-American
387  subclades was not supported in the CYTB phylogeny, so their divergence time is not
388  specified. A list of the divergence times with their 95% highest posterior density
389  intervals are present in Supplementary Table 4.

390

391  **Haplotype network**

392

393  To confirm if the samples of the subclades identified formed differentiated
394  subpopulations, haplotype networks using COX1 and CYTB were constructed. From
395  the COX1 multiple alignment, 43 positions of high variability were identified, supporting
396  25 haplotypes. These were diagrammed according to their genetic distances in a
397  haplotype network (Figure 5). Sequences comprising each haplotype are listed in
398  Supplementary Table 5. In contrast, the alignment of 31 CYTB sequences collapsed
399  just into 11 haplotypes, which were generated from 23 polymorphic sites. The
400  haplotype diversity (Hd) was 0.93 for COX1 and 0.88 for CYTB, respectively.

401

402  Both haplotype networks suggested that Asian subclades 1, 2, and 3 were closely
403  related. They formed a unique subpopulation with a dispersion center with a high
404  Indian component. In the COX1 network (Figure 5a), haplotype 6 (H6) was the
405  dispersion center. It was composed of sequences from Japan, India, Nepal, and
406  Madagascar, all separated by the ocean (except Nepal and India). Around it, unique
407  haplotypes were found distributed in India (H11, H16, H8) and Madagascar (H5, H7,
408  H9, H10). One branch connected to an unknown haplotype, which diverged into a
409  haplotype from Thailand (H1) and China (H17). The unknown haplotype was also
410  connected to Chinese haplotype H2, from which other Chinese and Nepali haplotypes
411  diverged. For CYTB (Figure 5b), the dispersion center (H6) was exclusively composed

412  of Indian samples. H6 was connected to Chinese, Nepali, and Vietnamese samples
413  from Asian subclade 3 through an unknown haplotype. Notably, the samples of Asian
414  subclade 3 showed a strong interconnection between them.
415
416  Indonesian samples remained somewhat isolated from the rest of Asian countries in
417  both networks, forming another subpopulation. For COX1, the isolate from Papua
418  (former Irian Jaya, H18) seemed to be more basal than the one of Bali (H19). For the
419  CYTB network, all Indonesian samples were clustered together (H11).
420
421  The COX1 network distinguished two separated groups related to the African-
422  American subclades 1 and 2, which differentiated in 3 punctual mutations. The CYTB
423  network also distinguished between African-American subclades. However, it
424  separated a haplotype conformed by a unique Tanzanian sample from the rest of the
425  countries of African-Subclade 1. To determine if the countries that formed African-
426  American subclade 1 were genetically different from those of the African-American
427  subclade 2, a computation of the $\phi$st value between these two groups was made.
428  Values were 0.83 for COX1 and 0.62 ($p < 0.05$) for CYTB.
429
430  **DISCUSSION**
431
432  The mitochondrial genomes from Puno and Mexico had 7,811X and 3,395X of
433  coverage, respectively. They had no ambiguous nucleotides and were the same size
434  as the reference. The sample from Huancayo had a lower coverage (42X). Although
435  this resulted in a partial genome, it was still adequate for the rest of the analysis. As
436  expected, the genome size and the %GC are similar between these three isolates,
437  supporting the assembly method. No gene rearrangements were detected (Figure 1,
438  Supplementary Table 2), and the size of each protein-coding gene is the same, except
439  for the gene CYTB in Huancayo, which has one codon less (Supplementary Table 2).
440  Nakao *et al.* [18] previously reported the presence of an abbreviated stop codon U (or
441  T in DNA) at the ND1 gene. Noteworthy, all the mitochondrial genomes assembled in
442  the present study present this stop coding, confirming this observation.
443
444  Variant calling showed that differences with respect to the Chinese reference are
445  almost identical in all the Latin American genomes and concentrated in protein-coding
446  genes, as revealed by the SNPs distribution (Figure 1). For instance, of the 34 non-
447  synonymous SNPs, 31 are present in all the assembled genomes (Supplementary
448  Tables 3).
449
450  Surprisingly, almost half of the Latin American SNPs involve a change in the amino
451  acid physicochemical nature (Figure 1, Supplementary Table 3). Those are located in
452  four mitochondrial protein-coding genes and especially in ND5. It is well-established
453  that change-in-nature mutations affect the structure and function of proteins [33].
454  Hence, the change-in-nature SNPs in the Latin American samples could have altered
455  the folding of their mitochondrial proteins compared to the proteins of Asian *T. solium*.
456  This may be linked to a reduced aerobic respiration efficiency that may affect their
457  fitness in oxygen-rich environments. The situation above-described could lead to a
458  negative tropism towards the oxygen-rich subcutaneous tissue, which has an $O_2$
459  partial pressure of 40-80 mmHg [34] higher than the ~23 mmHg [35] and <11 mmHg
460  [36,37] of the brain and intestinal lumen, respectively. Accordingly, it has been
461  suggested that subcutaneous cysticercosis is uncommon in Latin America but not in

462  Asia [38–40]. The evidence reported here calls for studies to confirm this clinical
463  difference between regions and relate it to the change-in-nature SNPs found.
464
465  The present results suggest that the CYTB sequence is the most variable of the whole
466  mitochondrial genome. It has the highest density of SNPs, the highest accumulative
467  D, and different D values for all pairwise comparisons (Figure 2). Besides, it has a
468  medium length of 1068 bp (Supplementary Table 2). Taking these features together
469  show that CYTB is a more suitable molecular marker for intraspecific variability
470  analysis that will differentiate isolates of the same or different genotypes. In
471  agreement, it has been stated that, within *taeniidae*, CYTB is a better marker for
472  reconstructing phylogenies among closely related groups (such as intraspecific
473  variations) because of its higher evolutionary rate [15]. Indeed, other studies have
474  reported results that support that CYTB has higher variability than, for example, COX1.
475  For instance, 28 SNPs (1.7% variability rate) in the COX1 gene were found in contrast
476  to the 31 SNPs (2.9% variability rate) present in CYTB [6]. Despite this, more complete
477  sequences for COX1 are available compared to CYTB. Additionally, COX1 sequences
478  have a greater variety of geographical origins. Extra CYTB sequences would allow
479  better and more informative phylogeographic reconstructions of *T. solium* lineages.
480
481  ATP6, COX2, and ND6 also have a relatively high D value; however, the three of them
482  have a short length. D values could overestimate the variability for small genes as the
483  percentage of identity is inversely correlated to the alignment length. Thus, high D
484  values for small sequences as these three should be taken cautiously and do not
485  necessarily imply high variability.
486
487  The region corresponding to the small and large ribosomal RNA (rRNA) and the
488  cysteine transfer RNA (tRNA-Cys) showed the lowest quantity of SNPs. This and the
489  fact that these SNPs are present in all the Latin American genomes suggest the low
490  variability of the region. Moreover, this region contains an internal sequence that
491  remains identical among the Chinese and Latin American genomes. The internal
492  sequence could be the target of conserved primers to specifically amplify the *T. solium*
493  mitochondrial DNA, as some portions differ from *T. saginata* and *T. asiatica*
494  (Supplementary Figure 1).
495
496  ND1 and ND3 are subjected to absolute purification (Figure 3), which is corroborated
497  by the fact that only synonymous SNPs were detected. In that sense, mutations in
498  these genes seem deleterious and therefore negatively selected. ND1 is a crucial
499  subunit of the mitochondrial respiratory complex I because it allocates the entrance of
500  the quinone reaction chamber and the first half part of the first proton translocation
501  channel, which receives input from the cytoplasm [41,42]. Moreover, ND3 provides the
502  flexibility needed for a concerted rearrangement that generates the driving force for
503  proton pumping [43]. Hence, mutations in these genes could affect the quinone
504  reductase activity and collapse the proton translocation system on the inner
505  mitochondrial membrane. The importance of both subunits is in agreement with the
506  fact that these genes are "cold spots" for amino acid mutations.
507
508  The phylogenetic analysis of both genes showed two main lineages: the Asian and the
509  African-American (Figure 4), which has been reported by other studies [6,8,13].
510  Divergence between those lineages occurred 0.458 [0.0405, 1.0625] (COX1, Figure

511 4a) to 0.634 [0.0667, 1.4574] (CYTB, Figure 4b) Myra during the Pleistocene, in
512 agreement with previous works [6,9,13].

514 One study has reported archaeological evidence that places modern humans in
515 Daoxian, China, 0.12 to 0.08 Myra [44]. However, it has not been confirmed if this was
516 a successful (persistent) settlement in south Asia. Interestingly, the present results
517 suggest that Chinese *T. solium* started to diverge around a similar period, 0.0786 to
518 0.0664 Myra according to Asian subclades 2 (both phylogenies, Figure 4a,b) and 3
519 (COX1 phylogeny, Figure 4a). The divergence within a geographical region requires
520 this region to be persistently settled by humans infected with *T. solium*. Hence, the
521 present data suggest that humans successfully populated south China around this
522 period and that the archaeological evidence might have arisen from these early
523 settlements.

525 The phylogenetic analyses suggested that the Asian genotype is further subdivided
526 into four (COX1) or three (CYTB) subclades. However, both haplotype networks
527 (Figure 5) indicated that Asian subclades 1, 2, and 3 are closely related, forming a
528 unique subpopulation with India as the center of dispersal [7]. That subpopulation is
529 differentiated from another formed by Indonesian samples. The different degrees of
530 relationship between Asian *T. solium* samples suggest heterogeneous gene flow.

532 In the COX1 haplotype network (Figure 5a), samples from Japan, India, Nepal, and
533 Madagascar are grouped in the same haplotype (H6). Considering that *T. solium* is
534 not endemic in Japan, its relation with H6 samples is likely the result of a recent
535 reintroduction, a not-so-rare event in the last years [45]. The Madagascan isolates'
536 close association with Nepali and Indian samples suggests that the parasite was
537 introduced into Madagascar from the Indian subcontinent [12]. A particular case
538 occurred with Nepali samples. One group of samples was included in H6 in close
539 association with Indian isolates; however, the other was included in H4 in close
540 association with Chinese sequences. These two genetic subpopulations suggest the
541 existence of two gene flows towards Nepal, one from the north (from China) and
542 another from the south (from India). They remain separated, possibly due to the
543 geographical barrier that the Himalayas constitute.

545 Asian subclade 4, which is formed by Indonesian samples, was the first to diverge
546 ~0.25 Myra according to COX1 and CYTB phylogenies. A similar divergence time was
547 found in previous research [9]. This correlated with a basal position concerning the
548 other Asian subclades in both phylogenies, which has also been reported [6,9].
549 Interestingly, both haplotypes networks grouped Indonesian samples in haplotypes
550 that remain distant from other Asian sequences. Furthermore, both haplotype
551 networks directly linked them to an unknown haplotype, which bridged the African to
552 the Asian lineages. Given that the unknown haplotype is the only bridge between the
553 two genotypes, it is proposed that it corresponds (or is related) to the *T. solium* modern
554 ancestor that independently dispersed in Asia and Africa. These results suggest that
555 Indonesian samples are more ancient and related to the modern *T. solium* ancestor
556 than other Asian samples. This could be explained by the introduction of the parasite
557 by early human migrations followed by isolation of the population (lack of gene flow).
558 Studying more *T. solium* samples from Indonesia could help identify plesiomorphies
559 and, by comparing them with more recent isolates, obtain an insight into how *T. solium*
560 evolves.

561

562     Indeed, *T. solium* subpopulations in Indonesia are isolated, as shown by the fact that
563 this parasite is restricted to Bali and Papua (former Irian Jaya). The COX1 haplotype
564 network suggested that the sample from Papua is more ancient than the one of Bali
565 because it is closer to other Asian isolates. Nonetheless, this proposal seems
566 inconsistent with the epidemiological evidence that suggests that the introduction of
567 *T. solium* to Papua was made 50 years ago from Bali [14,46]. Interestingly, even
568 though CYTB is more variable than COX1, the CYTB network shows no resolution to
569 distinguish between Indonesian haplotypes, while the COX1 network does. This
570 inconsistency might suggest that the haplotype differentiation seen in the COX1
571 haplotype network was an artifact attributed to a random selection when *T. solium* was
572 introduced into Papua, as it has been hypothesized in other work [14].

573

574     Remarkably, phylogenies and haplotype networks constructed in this work suggest
575 that the African-American lineage is further subdivided into two groups (African-
576 American subclades 1 and 2). The genetic differentiation between the two is confirmed
577 by the fact that $\phi$st values were high and significant ($p < 0.05$). Interestingly, the
578 subdivision was observed in a previous study [8]. Nonetheless, they did not include
579 African samples, making it impossible to perform inferences about the origin of both
580 subclades. The haplotype networks further confirmed the subdivision, allowing one to
581 visualize two different clusters.

582

583     Regarding African-American subclade 2, both haplotype networks showed a close
584 relation between samples from Brazil, Peru, Mexico, and Cameroon, considering the
585 majority were included in the same haplotype. Of note, this group only had sequences
586 from West Africa (Cameroon). The geographic composition and degree of association
587 of the samples of this subclade are coherent with the conversion between the Trans-
588 Atlantic slave and trade routes [7]. These trade routes mainly connected West Africa,
589 Europe, and the Americas.

590

591     As for the African-American subclade 1, the COX1 network revealed a close
592 association between samples from Mexico, Tanzania, Madagascar, and Ecuador. The
593 low differentiation and high genetic flow between one Mexican and one Tanzanian
594 isolate were reported previously [7]. Of note, this group exclusively included
595 sequences from East Africa (Tanzania and Madagascar). East African countries were
596 not direct participants in the Trans-Atlantic trade routes. Hence, their link with Latin
597 American samples might have been caused by another source of gene flow, such as
598 the one generated by the Indian Ocean slave trade. This trade connected East Africa,
599 Indian Ocean countries, the Middle East, and later the Americas.

600

601     All in all, the subdivision of the African-American genotype reveals that two different
602 sublineages exist in Latin America: one that derives from West Africa and the other
603 originated in East Africa. Two separate gene flows created from the Trans-Atlantic and
604 Indian Ocean trade routes may have caused the observed distribution. Notably,
605 Mexican samples were present in both lineages, which agrees with the proposal that
606 at least two genetic subpopulations coexist in Mexico [13]. Although the Peruvian
607 isolates were only included in the East African linage, subpopulations from the West
608 African linage are not discarded to exist, given that Mexico and Peru were trade
609 centers during European colonization.

610

611 In conclusion, thirteen SNPs with respect to the Chinese reference that involved a
612 change in the amino acid physicochemical nature were identified. They might be
613 related to differences in aerobic respiration efficiency between Asian and Latin
614 American *T. solium*. Further differences within the Asian genotype were also reported.
615 For instance, its differentiation and basal position combined with its divergence time
616 suggest that the subclade formed by Indonesian samples is the most ancient and
617 closely related to the modern *T. solium* ancestor than other Asian sequences.
618 Strikingly, all phylogeographic analyses showed that the African-American genotype
619 is subdivided into two subgroups. One has a strong relation with East African countries
620 while the other with West Africa, which might reflect the influence of the Trans-Atlantic
621 and the Indian Ocean trade routes. Of note, the isolates whose genomes were
622 assembled were part of the West African sublineage. In summary, the present study
623 shows that a detailed comparison of the mitochondrial variability of *T. solium* within
624 and between Asian and African-American genotypes still reveals interesting features
625 that could be used to combat *T. solium* diseases.

626
627
628 ## DATA AVAILABILITY
629
630 The assembled and annotated mitochondrial genomes from Puno and Huancayo were
631 uploaded to the GeneBank with accession numbers
632 IN_PROCESS_OF_SUBMISSION and KT591612, respectively. Regarding the
633 assembled genome from Mexico, nucleotide sequence data reported are available in
634 the Third Party Annotation Section of the DDBJ/ENA/GenBank databases under the
635 accession number TPA: BK061219. Other raw data will be available upon request.

636
637 ## ACKNOWLEDGEMENTS
638

643
644 ## FINANCIAL SUPPORT
645

648
649 ## CONFLICT OF INTEREST
650
651 The authors declare no conflict of interest.

652
653 ## REFERENCES
654
655 1.  Singh G, Burneo JG, Sander JW. From seizures to epilepsy and its substrates:
656     Neurocysticercosis. Epilepsia. 2013 May;54(5):783–92.
657 2.  Yoshino K. Studies on the post-embryonal development of Taenia solium. Part I. On the hatching
658     of the eggs of Taenia solium. J Med Assoc Formosa. 1933 Oct 18;32(10 (343)):1392–409.
659 3.  Yoshino K. Studies on the Postembryonal Development of Taenia solium. Part II. On the
660     Youngest Form of Cysticercus cellulosae and on the Migratory Course of the Oncosphaera of
661     Taenia solium within the Intermediate Host. J Med Assoc Formosa. 1933;32(11 (344)):1569–

662   86.
663   4.   Yoshino K. Studies on the postembryonal development of Taenia solium. Part III. On the
664        development of Cysticercus cellulosae within the definite intermediate host. J Med Assoc
665        Formosa. 1933;32(12 (345)):166–9.
666   5.   WHO. WHO estimates of the global burden of foodborne diseases: foodborne disease burden
667        epidemiology reference group 2007–2015. WHO Executive Summary. 2015 [cited 2022 Feb 18].
668   6.   Nakao M, Okamoto M, Sako Y, Yamasaki H, Nakaya K, Ito A. A phylogenetic hypothesis for the
669        distribution of two genotypes of the pig tapeworm Taenia solium worldwide. Parasitology. 2002
670        Jun 1;124(Pt 6):657–62.
671   7.   Martinez-Hernandez F, Jimenez-Gonzalez DE, Chenillo P, Alonso-Fernandez C, Maravilla P,
672        Flisser A. Geographical widespread of two lineages of Taenia solium due to human migrations:
673        Can population genetic analysis strengthen this hypothesis? Infect Genet Evol. 2009
674        Dec;9(6):1108–14.
675   8.   Solano D, Navarro JC, León-Reyes A, Benítez-Ortiz W, Rodríguez-Hidalgo R. Molecular
676        analyses reveal two geographic and genetic lineages for tapeworms, Taenia solium and Taenia
677        saginata, from Ecuador using mitochondrial DNA. Exp Parasitol. 2016 Dec;171:49–56.
678   9.   Michelet L, Carod J-F, Rakontondrazaka M, Ma L, Gay F, Dauga C. The pig tapeworm Taenia
679        solium, the cause of cysticercosis: Biogeographic (temporal and spacial) origins in Madagascar.
680        Mol Phylogenet Evol. 2010 May;55(2):744–50.
681   10.  Ito A, Budke CM. Genetic Diversity of Taenia solium and its Relation to Clinical Presentation of
682        Cysticercosis. Yale J Biol Med. 2021;94(2):343–9.
683   11.  Campbell G, Garcia HH, Nakao M, Ito A, Craig PS. Genetic variation in Taenia solium. Parasitol
684        Int. 2006 Jan;55(SUPPL.):S121–6.
685   12.  Yanagida T, Carod J-F, Sako Y, Nakao M, Hoberg EP, Ito A. Genetics of the Pig Tapeworm in
686        Madagascar Reveal a History of Human Dispersal and Colonization. Yao Y-G, editor. PLoS
687        One. 2014 Oct 15;9(10):e109002.
688   13.  Michelet L, Dauga C. Molecular evidence of host influences on the evolution and spread of
689        human tapeworms. Biol Rev. 2012 Aug;87(3):731–41.
690   14.  Yanagida T, Swastika K, Dharmawan NS, Sako Y, Wandra T, Ito A, et al. Origin of the pork
691        tapeworm Taenia solium in Bali and Papua, Indonesia. Parasitol Int. 2021 Aug;83(November
692        2020):102285.
693   15.  Okamoto M, Nakao M, Sako Y, Ito A. Molecular variation of Taenia solium in the world.
694        Southeast Asian J Trop Med Public Health. 2001 [cited 2022 Feb 13];32 Suppl 2(2):90–3.
695   16.  Hoberg EP, Alkire NL, Queiroz AD, Jones A. Out of Africa: origins of the Taenia tapeworms in
696        humans. Proc R Soc London Ser B Biol Sci. 2001 Apr 22;268(1469):781–7.
697   17.  Pajuelo MJ, Eguiluz M, Dahlstrom E, Requena D, Guzmán F, Ramirez M, et al. Identification
698        and Characterization of Microsatellite Markers Derived from the Whole Genome Analysis of
699        Taenia solium. Brehm K, editor. PLoS Negl Trop Dis. 2015 Dec 23;9(12):e0004316.
700   18.  Nakao M, Sako Y, Ito A. The mitochondrial genome of the tapeworm Taenia solium: a finding of
701        the abbreviated stop codon U. J Parasitol. 2003 Jun;89(3):633–5.
702   19.  Darling ACE, Mau B, Blattner FR, Perna NT. Mauve: Multiple Alignment of Conserved Genomic
703        Sequence With Rearrangements. Genome Res. 2004 Jul;14(7):1394–403.
704   20.  Rozas J, Ferrer-Mata A, Sánchez-DelBarrio JC, Guirao-Rico S, Librado P, Ramos-Onsins SE,
705        et al. DnaSP 6: DNA Sequence Polymorphism Analysis of Large Data Sets. Mol Biol Evol. 2017
706        Dec 1;34(12):3299–302.
707   21.  Krzywinski M, Schein J, Birol İ, Connors J, Gascoyne R, Horsman D, et al. Circos: An information
708        aesthetic for comparative genomics. Genome Res. 2009 Sep;19(9):1639–45.
709   22.  Lavikainen A, Haukisalmi V, Lehtinen MJ, Henttonen H, Oksanen A, Meri S. A phylogeny of
710        members of the family Taeniidae based on the mitochondrial cox1 and nad1 gene data.
711        Parasitology. 2008 Oct 21;135(12):1457–67.
712   23.  Madeira F, Park YM, Lee J, Buso N, Gur T, Madhusoodanan N, et al. The EMBL-EBI search
713        and sequence analysis tools APIs in 2019. Nucleic Acids Res. 2019 Jul 2;47(W1):W636–41.
714   24.  Katoh K, Rozewicki J, Yamada KD. MAFFT online service: multiple sequence alignment,
715        interactive sequence choice and visualization. Brief Bioinform. 2019 Jul 19;20(4):1160–6.
716   25.  Kuraku S, Zmasek CM, Nishimura O, Katoh K. aLeaves facilitates on-demand exploration of
717        metazoan gene family trees on MAFFT sequence alignment server with enhanced interactivity.
718        Nucleic Acids Res. 2013 Jul 1;41(W1):W22–8.
719   26.  Castresana J. Selection of Conserved Blocks from Multiple Alignments for Their Use in
720        Phylogenetic Analysis. Mol Biol Evol. 2000 Apr 1;17(4):540–52.
721   27.  Talavera G, Castresana J. Improvement of Phylogenies after Removing Divergent and

Ambiguously Aligned Blocks from Protein Sequence Alignments. Kjer K, Page R, Sullivan J, editors. Syst Biol. 2007 Aug 1;56(4):564–77.

28. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. Bioinformatics. 2014 May 1;30(9):1312–3.

29. Miller MA, Pfeiffer W, Schwartz T. Creating the CIPRES Science Gateway for inference of large phylogenetic trees. In: 2010 Gateway Computing Environments Workshop (GCE). IEEE; 2010. p. 1–8.

30. Darriba D, Taboada GL, Doallo R, Posada D. jModelTest 2: more models, new heuristics and parallel computing. Nat Methods. 2012 Aug 30;9(8):772–772.

31. Drummond AJ, Rambaut A. BEAST: Bayesian evolutionary analysis by sampling trees. BMC Evol Biol. 2007 Nov 8;7(1):214.

32. Bandelt HJ, Forster P, Rohl A. Median-joining networks for inferring intraspecific phylogenies. Mol Biol Evol. 1999 Jan 1;16(1):37–48.

33. Volkenstein M V. Coding of Polar and Non-polar Amino-acids. Nature. 1965 Jul;207(4994):294–5.

34. Völker H-U, Röper G, Sterk J, Willy C. Long-term invasive measurement of subcutaneous oxygen partial pressure above the sacrum on lying healthy volunteers. Wound Repair Regen. 2006 Sep;14(5):542–7.

35. Pennings FA, Schuurman PR, van den Munckhof P, Bouma GJ. Brain Tissue Oxygen Pressure Monitoring in Awake Patients during Functional Neurosurgery: The Assessment of Normal Values. J Neurotrauma. 2008 Oct;25(10):1173–7.

36. He G, Shankar RA, Chzhan M, Samouilov A, Kuppusamy P, Zweier JL. Noninvasive measurement of anatomic structure and intraluminal oxygenation in the gastrointestinal tract of living mice with spatial and spectral EPR imaging. Proc Natl Acad Sci. 1999 Apr 13;96(8):4586–91.

37. Fisher EM, Khan M, Salisbury R, Kuppusamy P. Noninvasive Monitoring of Small Intestinal Oxygen in a Rat Model of Chronic Mesenteric Ischemia. Cell Biochem Biophys. 2013 Nov 1;67(2):451–9. Available from: http://link.springer.com/10.1007/s12013-013-9611-y

38. Feng Y, Ouyang S, Zhou X, Yang S. Clinicoelectroencephalographic studies of cerebral cysticercosis 158 cases. Chin Med J (Engl). 1979 Nov;92(11):770–86.

39. Simanjuntak GM, Margono SS, Okamoto M, Ito A. Taeniasis/cysticercosis in Indonesia as an emerging disease. Parasitol Today. 1997 Sep;13(9):321–3.

40. Cruz I, Cruz ME, Teran W, Schantz PM, Tsang V, Barry M. Human subcutaneous Taenia solium cysticercosis in an Andean population with neurocysticercosis. Am J Trop Med Hyg. 1994 Oct 1 [cited 2022 Feb 13];51(4):405–7.

41. Baradaran R, Berrisford JM, Minhas GS, Sazanov LA. Crystal structure of the entire respiratory complex I. Nature. 2013 Feb 17;494(7438):443–8.

42. Parey K, Haapanen O, Sharma V, Köfeler H, Züllig T, Prinz S, et al. High-resolution cryo-EM structures of respiratory complex I: Mechanism, assembly, and disease. Sci Adv. 2019 Dec 6;5(12):1–11.

43. Cabrera-Orefice A, Yoga EG, Wirth C, Siegmund K, Zwicker K, Guerrero-Castillo S, et al. Locking loop movement in the ubiquinone pocket of complex I disengages the proton pumps. Nat Commun. 2018 Dec 29;9(1):4500.

44. Liu W, Martinón-Torres M, Cai Y, Xing S, Tong H, Pei S, et al. The earliest unequivocally modern humans in southern China. Nature. 2015 Oct 14;526(7575):696–9.

45. Yanagida T, Sako Y, Nakao M, Nakaya K, Ito A. Taeniasis and cysticercosis due to Taenia solium in Japan. Parasit Vectors. 2012 Dec 17;5(1):18.

46. Sutisna P, Kapti IN, Wandra T, Dharmawan NS, Swastika K, Raka Sudewi AA, et al. Towards a cysticercosis-free tropical resort island: A historical overview of taeniasis/cysticercosis in Bali. Acta Trop. 2019 Feb;190(October 2018):273–83.

## FIGURE LEGENDS

**Figure 1. Nucleotide and amino acid differences in the 5 *T. solium* genomes.** The thick outer circle is the Chinese *T. solium* reference mitochondrial genome (NCBI Reference Sequence: NC_004022.1), where inner boxes represent the coding sequences (CDS). The color code represents the CDS type: purple for protein-coding genes, green for tRNA's, orange for ribosomal RNAs, and gray for non-coding regions. The inner circles represent the 3 *T. solium* mitogenomes assembled in this study (black: Puno, light blue: Huancayo, and pink: Mexico). The blue bars indicate synonymous nucleotide substitutions on each inner circle, while red bars highlight those causing amino acid substitutions. Flanking arrows in specific red bars indicate substitutions that involve a change in the amino acid nature. The circular segments shaded in transparent blue, and red indicate low and high variability regions, respectively. Darker blue within the blue-shaded circular segment indicates an identical region conserved in the four mitochondrial genomes. SNPs in the low variability region suggest that the region could differentiate Asian isolates from African-American isolates (as SNPs are differences with respect to the Chinese genome). However, the identical SNP distribution in all Latin American genomes suggests that it could not differentiate isolates of the same genotype, such as the African-American.

**Figure 2. D values per protein-coding gene of each possible pairwise combination of 4 mitochondrial genomes (Chinese reference, Huancayo, Puno, and Mexico).** D was used to indicate the level of sequence difference between the four mitochondrial genomes used in this study. D values were computed by performing all possible pairwise alignments and applying the D = 1 - (M/L) formula. L is the difference between the alignment length and the number of ambiguous codons. M is the number of invariant sites in the alignment. Different D values for each combination were stacked and presented per protein-coding gene in a bar plot. Thus, the height of a particular bar of a gene corresponds to the sum of D values for the different pairwise combinations; in other words, the bar height is an accumulative D value.

**Figure 3. Ka/Ks ratios per protein-coding gene.** Ka/Ks ratios against the Chinese reference mitochondrial genomes were calculated using the whole-genome multiple alignment in the DNAsp software for each protein-coding gene of the mitochondrial genomes. Ka/Ks values of the same protein-coding gene were stacked together (bars indicate accumulative Ka/Ks).

**Figure 4. Phylogenetic trees constructed from COX1 and CYTB complete sequences.** (a) Bayesian Inference (BI) tree constructed with 45 *T. solium* COX1 complete sequences. (b) BI tree constructed with 31 CYTB complete sequences. Posterior probabilities (PP) above 0.7 are shown until the level of Asian and African-American subclades. Maximum Likelihood (ML) trees were also constructed using the same sequences. If a clade has specified a Bootstrap (BS) value, it means it was present in the BI and ML tree, with a BS greater than 70% for the latter. BS values above 70% are shown only until the level of Asian and African-American subclades. For both trees, the distances of the branches are based on the timeline at a scale of one million years (Myr). Only the nodes supported in the ML and BI with BS > 70% or PP > 0.7 have their divergence times (DT) shown in Myr. Again, DT is only shown until subclades. 95% highest posterior densities are specified in Supplementary Table 4. (*) BS: 78, DT: 0.0759. (+) BS: 86, PP: 1.00, DT: 0.0495.

**Figure 5. Haplotype network of COX1 and CYTB.** (a) COX1 network. (b) CYTB network. The geographic origins of the samples included in each haplotype are color-coded. Colored squares indicate the most important clades and subclades.

## SUPPLEMENTARY MATERIAL LEGENDS

**Supplementary Figure 1. Selected portions of the multiple alignment of the identical internal region.** The sequences corresponding to the identical internal region conserved in the *T. solium* mitochondrial genome from China (NC_004022.1), Puno, Huancayo, and Mexico (see Figure 1) were aligned with the mitochondrial genomes of *T. asiatica* (NC_004826.2), *T. saginata* (NC_009938.1) and *E. multilocularis* (NC_000928.2) in the online version of Clustal Omega [23]. The alignment was then cropped to include only matching sections using CLC Genomics Workbench v. 21.05.5 (https://digitalinsights.qiagen.com/). Selected portions of the alignment that showed significant differences between *T. solium* and the other organisms are presented.

840 **Supplementary Table 1. Mapping statistics of the genome assembly.** The number of reads mapped
841 to the Chinese mitochondrial genome (reference) and their mean quality scores before trimming (pre)
842 are given. The number of reads mapped after trimming (post), the coverage, the genome length in base
843 pairs (bp), the %GC, and the N50 of the assembly in nucleotides (nt) are also shown.
844
845 **Supplementary Table 2. Gene arrangement of *T. solium* mitochondrial genomes from Peru,**
846 **Mexico, and China (reference).** The size of each genome in base pairs (bp) is given in parenthesis
847 next to the mitochondrial genome name. Position intervals per gene are shown. The size of each gene
848 (in bp) is specified in parentheses next to each position interval. Start and stop codons of each protein-
849 coding gene per mitochondrial genome are also specified.
850
851 **Supplementary Table 3. Non-synonymous mutations in the Latin American mitochondrial**
852 **genomes with respect to the Chinese reference.** Amino acid mutations in all the assembled genomes
853 of this work are given per protein-coding genes. In the present work, standard amino acids were
854 classified into five groups, considering their major species at pH of 7: polar and uncharged (S, T, N, and
855 Q), polar and positively charged (R, H, and K), Polar and negatively charged (D, E), nonpolar (A, V, I,
856 L, M, F, Y, and W) and special cases (C, G, and P). A mutation that changes the amino acid nature
857 involves that the mutant amino acid belongs to a different group than the original. Mutations highlighted
858 in red and blue represent changes that do and do not involve a change in the amino acid nature,
859 respectively.
860
861 **Supplementary Table 4. Divergence events and their dating.** Divergence events in COX1 and CYTB
862 phylogenies are described. The gene in whose phylogeny the event occurred is specified in bold at the
863 start of the divergence event description. If nothing is specified, the same event occurred in both
864 phylogenies. The divergence times with 95% highest posterior density intervals (brackets) in millions of
865 years ago (Myra) are also reported for each event.
866
867 **Supplementary Table 5. Sequence composition of the haplotypes formed in the COX1 and CYTB**
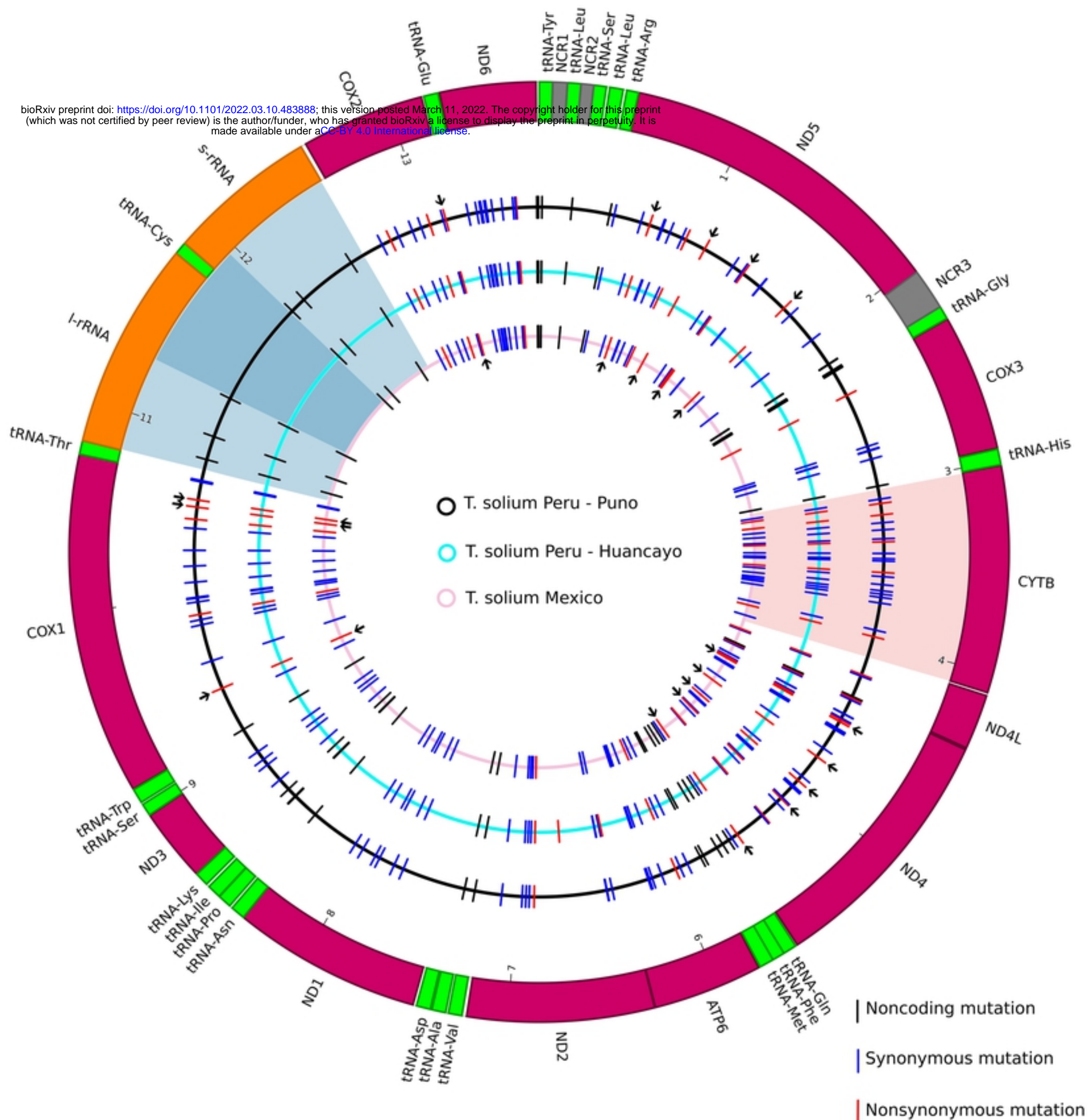868 **networks.** Sequences (with their accession numbers) included in each haplotype are specified.
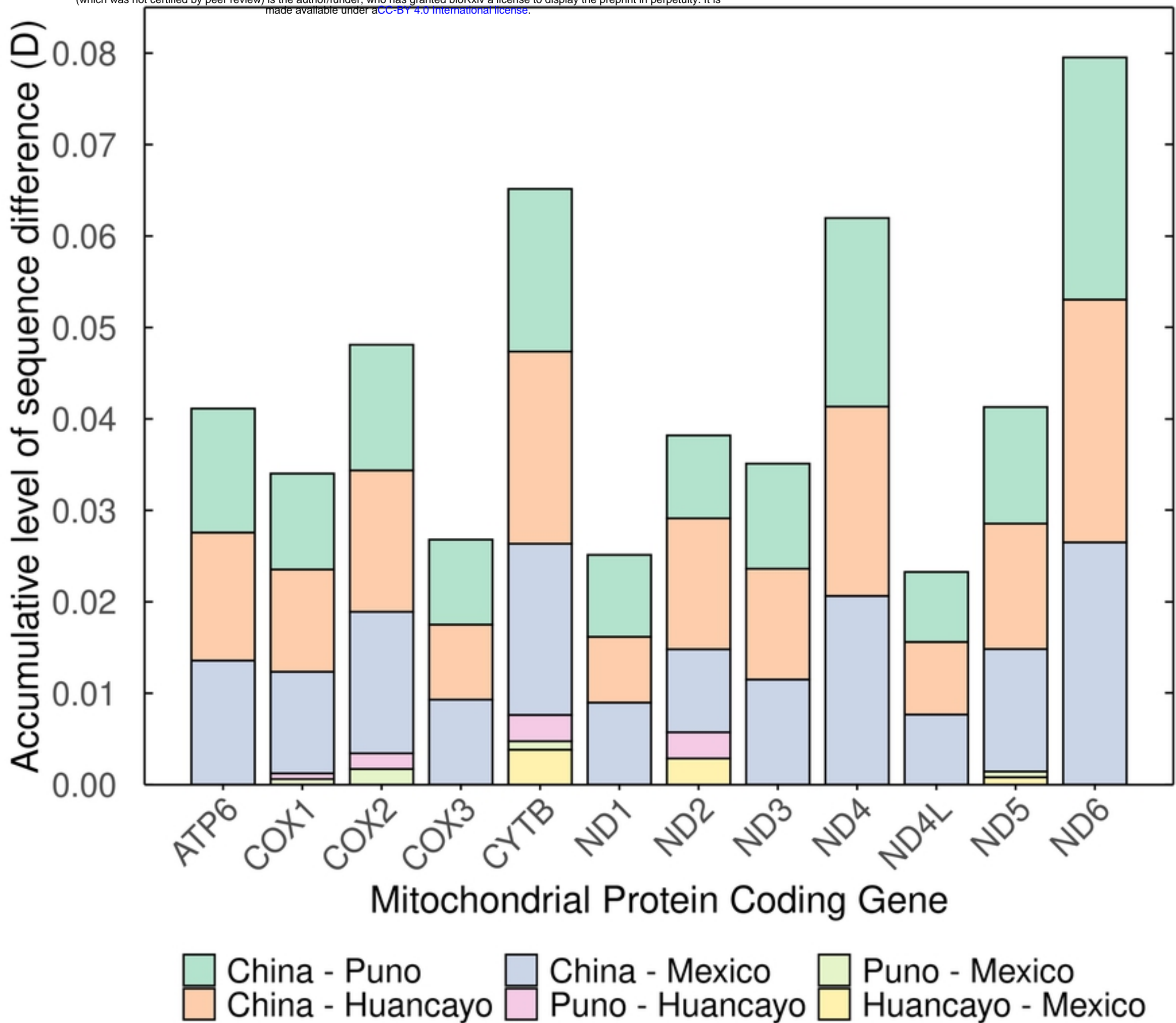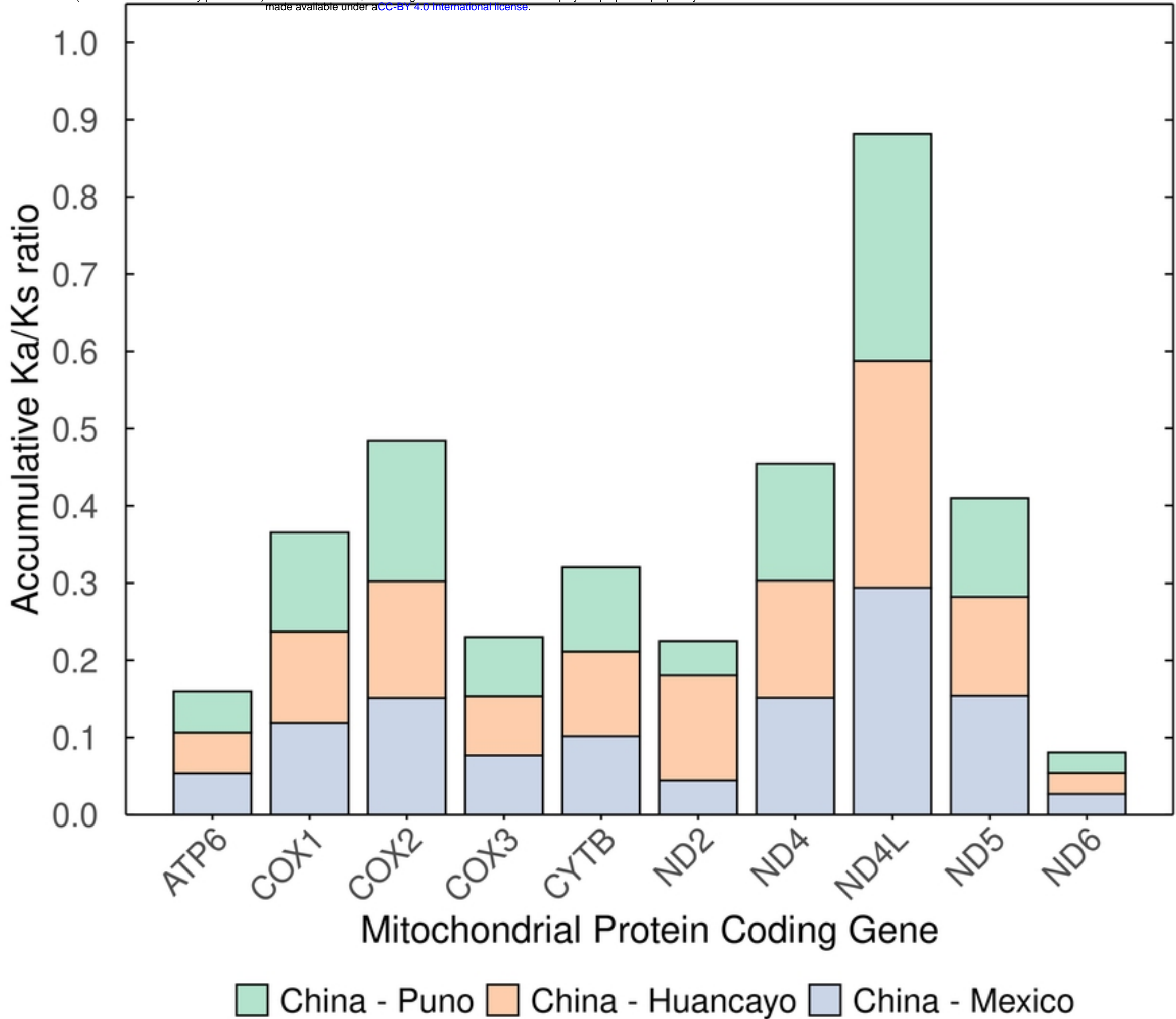
Figure 1

Figure 2

Figure 3

**Legend:**
— Asian — Asian subclade 1 — Asian subclade 2 — Asian subclade 3 — Asian subclade 4
— African-American — African-American subclade 1 — African-American subclade 2
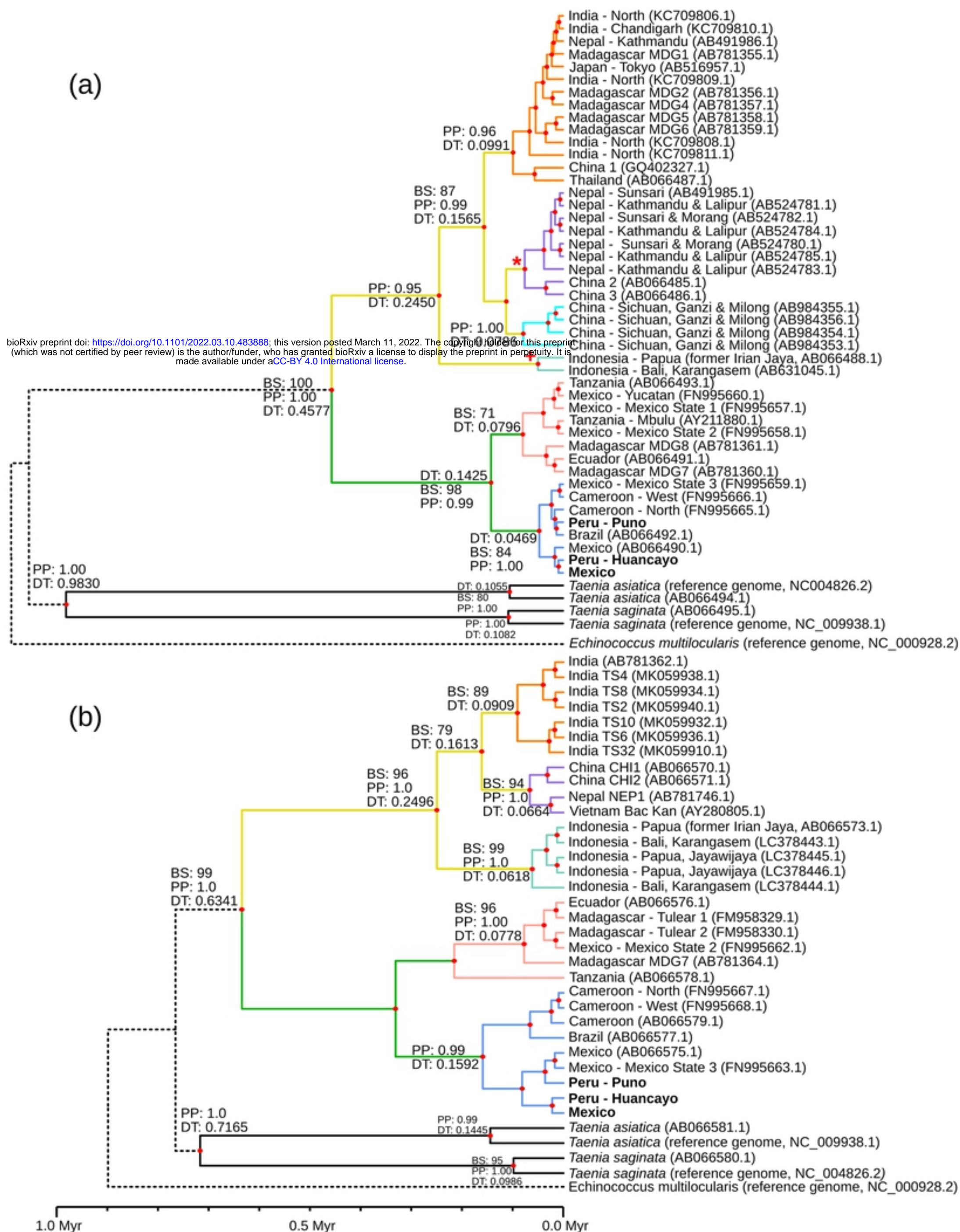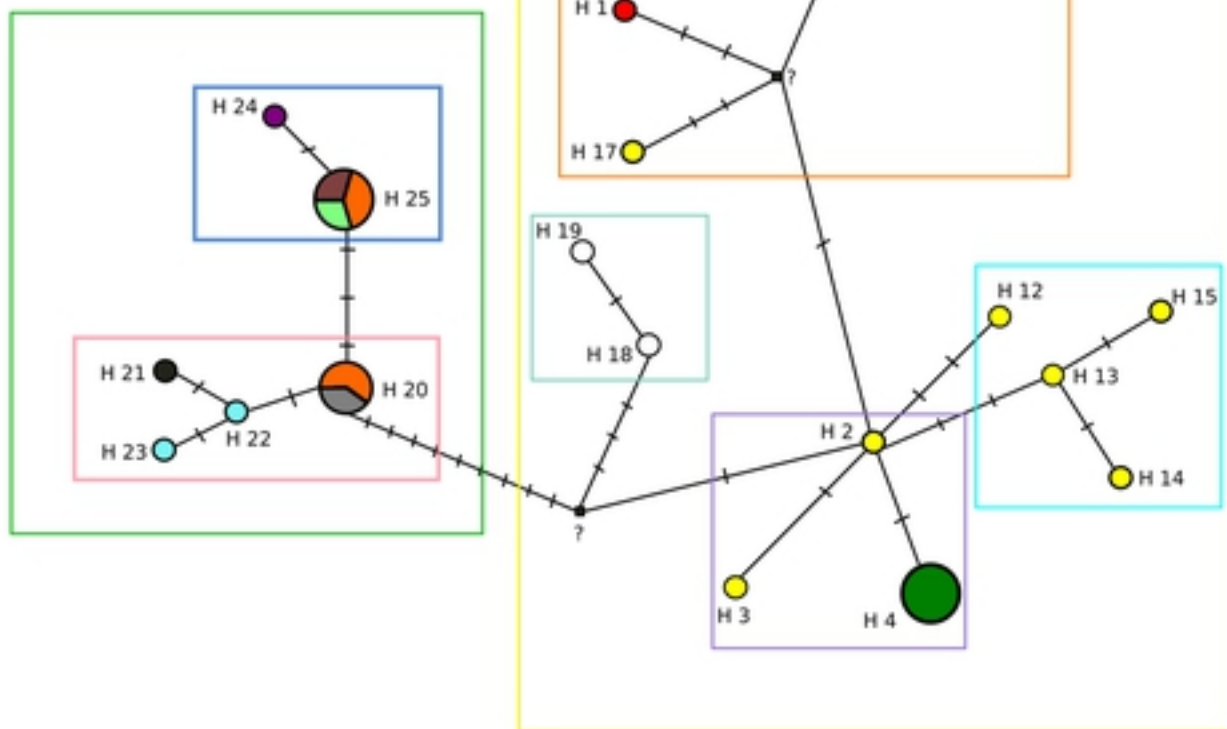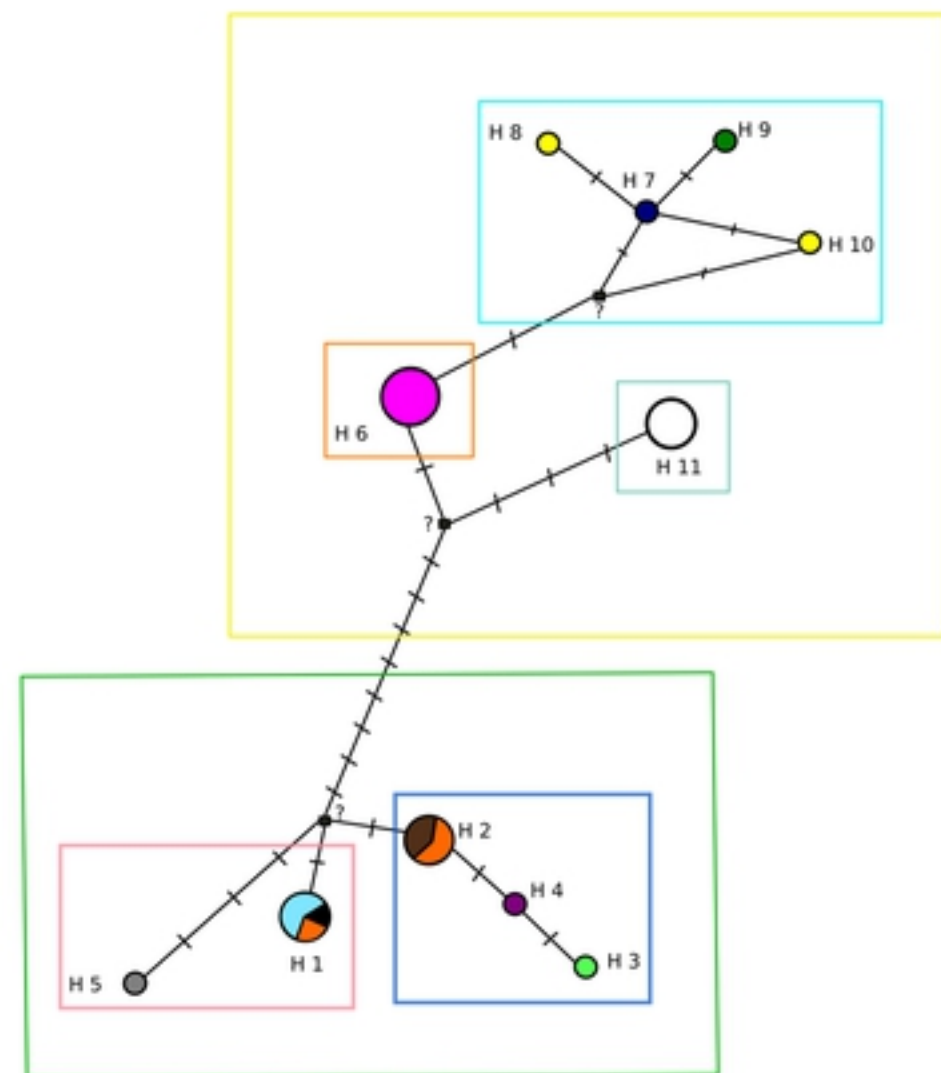
**(a)**

India - North (KC709806.1)
India - Chandigarh (KC709810.1)
Nepal - Kathmandu (AB491986.1)
Madagascar MDG1 (AB781355.1)
Japan - Tokyo (AB516957.1)
India - North (KC709809.1)
Madagascar MDG2 (AB781356.1)
Madagascar MDG4 (AB781357.1)
Madagascar MDG5 (AB781358.1)
Madagascar MDG6 (AB781359.1)
India - North (KC709808.1)
India - North (KC709811.1)
China 1 (GQ402327.1)
Thailand (AB066487.1)
Nepal - Sunsari (AB491985.1)
Nepal - Kathmandu & Lalipur (AB524781.1)
Nepal - Sunsari & Morang (AB524782.1)
Nepal - Kathmandu & Lalipur (AB524784.1)
Nepal - Sunsari & Morang (AB524780.1)
Nepal - Kathmandu & Lalipur (AB524785.1)
Nepal - Kathmandu & Lalipur (AB524783.1)
China 2 (AB066485.1)
China 3 (AB066486.1)
China - Sichuan, Ganzi & Milong (AB984355.1)
China - Sichuan, Ganzi & Milong (AB984356.1)
China - Sichuan, Ganzi & Milong (AB984354.1)
China - Sichuan, Ganzi & Milong (AB984353.1)
Indonesia - Papua (former Irian Jaya, AB066488.1)
Indonesia - Bali, Karangasem (AB631045.1)
Tanzania (AB066493.1)
Mexico - Yucatan (FN995660.1)
Mexico - Mexico State 1 (FN995657.1)
Tanzania - Mbulu (AY211880.1)
Mexico - Mexico State 2 (FN995658.1)
Madagascar MDG8 (AB781361.1)
Ecuador (AB066491.1)
Madagascar MDG7 (AB781360.1)
Mexico - Mexico State 3 (FN995659.1)
Cameroon - West (FN995666.1)
Cameroon - North (FN995665.1)
**Peru - Puno**
Brazil (AB066492.1)
Mexico (AB066490.1)
**Peru - Huancayo**
**Mexico**
*Taenia asiatica* (reference genome, NC004826.2)
*Taenia asiatica* (AB066494.1)
*Taenia saginata* (AB066495.1)
*Taenia saginata* (reference genome, NC_009938.1)
*Echinococcus multilocularis* (reference genome, NC_000928.2)

PP: 0.96 / DT: 0.0991
BS: 87 / PP: 0.99 / DT: 0.1565
PP: 0.95 / DT: 0.2450
PP: 1.00
BS: 100 / PP: 1.00 / DT: 0.4577
BS: 71 / DT: 0.0796
DT: 0.1425 / BS: 98 / PP: 0.99
DT: 0.0469 / BS: 84 / PP: 1.00
PP: 1.00 / DT: 0.9830
DT: 0.1055 / BS: 80 / PP: 1.00 / PP: 1.00 / DT: 0.1082

**(b)**

India (AB781362.1)
India TS4 (MK059938.1)
India TS8 (MK059934.1)
India TS2 (MK059940.1)
India TS10 (MK059932.1)
India TS6 (MK059936.1)
India TS32 (MK059910.1)
China CHI1 (AB066570.1)
China CHI2 (AB066571.1)
Nepal NEP1 (AB781746.1)
Vietnam Bac Kan (AY280805.1)
Indonesia - Papua (former Irian Jaya, AB066573.1)
Indonesia - Bali, Karangasem (LC378443.1)
Indonesia - Papua, Jayawijaya (LC378445.1)
Indonesia - Papua, Jayawijaya (LC378446.1)
Indonesia - Bali, Karangasem (LC378444.1)
Ecuador (AB066576.1)
Madagascar - Tulear 1 (FM958329.1)
Madagascar - Tulear 2 (FM958330.1)
Mexico - Mexico State 2 (FN995662.1)
Madagascar MDG7 (AB781364.1)
Tanzania (AB066578.1)
Cameroon - North (FN995667.1)
Cameroon - West (FN995668.1)
Cameroon (AB066579.1)
Brazil (AB066577.1)
Mexico (AB066575.1)
Mexico - Mexico State 3 (FN995663.1)
**Peru - Puno**
**Peru - Huancayo**
**Mexico**
*Taenia asiatica* (AB066581.1)
*Taenia asiatica* (reference genome, NC_009938.1)
*Taenia saginata* (AB066580.1)
*Taenia saginata* (reference genome, NC_004826.2)
*Echinococcus multilocularis* (reference genome, NC_000928.2)

BS: 89 / DT: 0.0909
BS: 79 / DT: 0.1613
BS: 96 / PP: 1.0 / DT: 0.2496
BS: 94 / PP: 1.0 / DT: 0.0664
BS: 99 / PP: 1.0 / DT: 0.0618
BS: 99 / PP: 1.0 / DT: 0.6341
BS: 96 / PP: 1.00 / DT: 0.0778
PP: 0.99 / DT: 0.1592
PP: 1.0 / DT: 0.7165
PP: 0.99 / DT: 0.1445
BS: 95 / PP: 1.00 / DT: 0.0986

1.0 Myr          0.5 Myr          0.0 Myr

**Figure 4**

Figure 5