1 **Article**

2 **Whole-genome-based *Helicobacter pylori* geographic surveillance: a visualized and**

3 **expandable webtool**

4 Xiaosen Jiang[1,2,3] [†], Zheng Xu[1,4] [†], Tongda Zhang[1], Yuan Li [1], Wei Li[1,2,3] Hongdong Tan[1*]

5 [1] BGI-Shenzhen, Shenzhen, 518083, China

6 [2] BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083,

7 China

8 [3] College of Life Sciences, University of Chinese Academy of Sciences, Beijing, 100049,

9 China

10 [4] Shenzhen Key Laboratory of Unknown Pathogen Identification, BGI-Shenzhen, Shenzhen,

11 China

12

13 [†] These authors contributed equally: Xiaosen Jiang and Zheng Xu

14

15 * Corresponding Author

16 Hongdong Tan

17 Email: rtan@mgi-tech.com

18

19

20

21

22

23    **Abstract**

24    *Helicobacter pylori* exhibits specific geographic distributions that related to the clinical

25    outcomes. Despite the high infection rate of *H. pylori* throughout the world, the genetic

26    epidemiology surveillance of *H. pylori* still needs to be improved. Here, we used single

27    nucleotide polymorphisms (SNPs) profiling approach based on whole genome sequencing

28    (WGS) that facilitates genomic population analyses of *H. pylori* and encourages the

29    dissemination of microbial genotyping strategies worldwide. A total number of 1,211 public

30    *H. pylori* genomes were downloaded and used to construct the typing tool, named as HPTT

31    (*H. pylori* Typing Tool). Combined with the metadata, we developed two levels of genomic

32    typing, including a continent scale and a country scale that nested in the continent scale.

33    Results showed that Asia was the largest isolates source in our dataset, while isolates from

34    Europe and Oceania were comparatively more widespread. More specifically, Switzerland

35    and Australia are the main source of widespread isolates in their corresponding continents.

36    To integrate all the typing information and enable researchers to compare their own dataset

37    against the existing global database in an easy and rapid way, a user-friendly website

38    (https://db.cngb.org/HPTT/) was developed with both genomic typing tool and visualization

39    tool. To further confirm the validity of the website, ten newly assembled genomes were

40    downloaded and tested precisely located on the branch as we expected. In summary, *H. pylori*

41    typing tool (HPTT) is a novel genomic epidemiological tool that can achieve high resolution

42    analysis of genomic typing and visualizing simultaneously, providing insights into the

43    genetic population structure analysis, evolution analysis and epidemiological surveillance of

44    *H. pylori*.

45

## Introduction

*Helicobacter pylori* is one of the most sophisticated colonizers in the world that infects more than half of world's population ranged from infants to elders (Suerbaum and Michetti 2002). It is a Gram-negative bacterium that normally colonises at the gastric mucosa of human with about 10% infection result in diseases. The typical diseases were reported as gastritis, peptic ulcer, mucosa-associated lymphoid tissue (MALT) lymphoma and gastric cancer (Ernst and Gold 2000). Globally speaking, the risks of disease and the incidence and mortality of the gastric cancer were geographically different (Group 1993).

*H. pylori* displays a distinguished mutation rate among bacterial pathogens due to the lack of genes that initiates classical methyl-directed mismatch repair (MMR) (Alm, et al. 1999). The high mutation and recombination rate made *H. pylori* genomes with enormous plasticity, facilitating this pathogen perfectly adapted to its host (Kang and Blaser 2006; Didelot, et al. 2013). It has been reported that *H. pylori* in chronic infection could be taken place through vertical and familial transmission (Agnew and Koella 1997; Messenger, et al. 1999). In within-host evolution, the mutation rate could reach ~ 30 single nucleotide polymorphisms (SNPs) per genome per year (Kennemann, et al. 2011), comparing to *Escherichia coli* at ~ 1 SNP per genome per year (Reeves, et al. 2011). With the occurrence of large recombination events, a simple and efficient way to define the geographical pattern and epidemiological surveillance of *H. pylori* is crucially needed (Yamaoka 2009; Jolley, et al. 2018).

Among all genetic typing methods recorded in previous studies (Salama, et al. 2007; Yamaoka 2009), seven-gene multi-locus sequence typing (MLST) for *H. pylori* is a current popular tool due to its simple and rapid typing strategy. The 7-gene MLST covers genes including *atpA*, *efp*, *mutY*, *ppa*, *trpC*, *urel*, *yphC* that categorize *H. pylori* into different sequence types (STs) (Achtman, et al. 1999). This 7-gene MLST typing method enables regional specific recognition based on defined STs, through which geographical pattern may be linked with the different risk of clinical disease. For example, non-African and African lineage could be in association with different risk of gastric disease (Campbell, et al. 2001). However, the resolution of seven-gene MLST was still low, which limited us to trace the epidemiological origins of *H. pylori* strains. For users, submitting the microbial genomes is essential to get the allele number before getting the typing results. The seven-gene genotypes of *H. pylori* is diverse due to the high variability of *H. pylori* genomes, which hinders the recognition of patterns directly from the 4-digit code in 7-gene MLST. In addition, there is no

78    information of geographical patterns or visualization tool for seven-gene MLST, thus such
79    related geographic patterns were hard to find when a new ST was found.

80    Here, we describe a *H. pylori* genomic typing tool, HPTT (*H. pylori* Typing Tool) using the
81    SNP profiling based on whole-genome sequencing data. In addition to the genomic typing,
82    HPTT also provides a phylogenetic and geographic visualization tool based on the Nextstrain
83    framework (Hadfield, et al. 2018). This tool allows users to upload *H. Pylori* WGS data for
84    genomic typing and uncover the possible transmission events of *H. Pylori.* It is believed that
85    this tool can not only improve the genome typing resolutions, but also predict the origin of
86    the epidemic *H. pylori* isolates, enabling the global surveillance of *H. pylori*.

87

88    **Methods**

89    ***H. pylori* genomes downloaded and filtered in this study**

90    A total number of 1,654 assembled *H. pylori* genomes were downloaded from NCBI RefSeq
91    database (genomes available as of 4th May 2020) using ncbi-genome-download tool (version
92    0.2.12). The corresponding metadata of assembled genomes was searched by function using
93    Entrez Direct (version 10.9) (Kans 2020). By metadata filtering, 1,211 genomes were
94    selected with sample collection location available (**Table 1**). All genomes were scanned by
95    mlst (version 2.11) with the library of MLST updated on 31 December 2020 (Jolley and
96    Maiden 2010).

97    **SNP analysis**

98    The 1,211 assembled genomes were mapped to the reference genome *H. pylori* 26695
99    (GenBank: AE000511.1) (Tomb, et al. 1997) using MUMmer (version 3.23) (Kurtz, et al.
100   2004). SNPs were filtered with a minimum mapping quality cutoff at 0.90 across 1,211
101   assembled *H. pylori* genomes. 6,129 SNPs were found, and a SNP profile of *H. Pylori* is
102   established for the corresponding isolates.

103   **Phylogenetic analysis**

104   The maximum likelihood (ML) phylogenetic tree was constructed by iqtree (version 2.0.3)
105   (Nguyen, et al. 2015) based on 6,129 SNPs alignments of all 1,211 isolates. The reference

106   genome *H. pylori* 26695 was used as outgroup. The tree was generalized by the Gamma

107   distribution to model site-specific rate variation (the GTR model). Bootstrap pseudo-analyses

108   of the alignment were set at >= 1000. All ML trees were visualized and annotated using

109   Figtree (version 1.4.4). The minimum spanning tree was constructed by the GrapeTree

110   (v1.5.0) (Zhou, et al. 2018).

**Geographic typing system**

112   Based on the phylogenetic tree, two levels of geographic group were defined, including the

113   first level defined at the continent scale and the second level defined as country specific

114   scale. In the first level of genotyping, lineages that carrying more than seven isolates and >

115   75% isolates sourced from one major continent were defined as a continent specific group or

116   clade. A mixed continent group was defined when there was no major continent identified

117   with isolates at >75%. In the second level, lineages carrying more than one isolate and > 75%

118   isolates sourced from one major country were defined as a country specific group or

119   subclade. In addition, a mixed group was also defined at level two when there were more

120   than two isolates and not a major country identified with isolates at > 75%. The association of

121   the genomic lineage of *H. Pylori* with the geographic origin of isolates fully sequenced

122   provide a map to allow us tracing both the origin and evolution path of a detected or

123   sequenced *H. Pylori* genome.

**Establishment of *H. pylori* database**

125   The HPTT website was established based on two modules: 1) The genomic-geographical

126   typing tool of *H. pylori* isolates and 2) a visualization tool of both the genomic and

127   geographic typing results. The online typing tool was written in PHP, Javascript, css, and

128   html. The online visualization service was performed based on the CodeIgniter framework

129   (https://www.codeigniter.com/), tree visualization was analysed by the augur

130   (https://github.com/nextstrain/augur) bioinformatics tool and the auspice

131   (https://github.com/nextstrain/auspice) visualization tool imbedded in the Nextstrain

132   (Hadfield, et al. 2018) open source project. The *H. pylori* database was stored in a Mysql

133   database.

134

135

136 **Results**

137 **Definition of two levels of geographic genotypes for *H. pylori***

138 A total of 1,211 assembled genomes with available geographic information from NCBI

139 RefSeq database were downloaded and analyzed for establishing the *H. pylori* genotyping

140 database (**Supplementary Table 1**). All assembly genomes were mapped to the reference

141 genome *H. pylori* 26695. Based on the maximum likelihood tree, while 6,129 SNPs were

142 defined for further genomic typing. In terms of geographic information, 1,112 isolates were

143 grouped at two levels, including 37 continent-level groups (**Figure 1A**) and/or 236 country-

144 level groups (**Figure 1C**). The median pairwise distances between isolates were found as

145 follows: 319 SNPs within continent clades and 1,493 SNPs within country subclades. We

146 labelled these continent clades and country subclades using a structured hierarchical

147 nomenclature system similar to that used for *M. tuberculosis* (Coll, et al. 2014). For instance,

148 region 1 clade (G1) is subdivided into country subclades G1.C1 and G1.C2.

149 **A continent level genomic typing for *H. pylori***

150 A total number of 37 continent level of groups (n=1,112) were defined, including 25

151 continent specific groups and the 12 mixed continent groups (**Figure 1A&B**). Isolates across

152 the tree did not fall into the continent group but can be defined as country group were named

153 as G0 (n=74). Isolates across the tree neither fall into the country group nor continent group

154 was defined as non-grouped (n=25).

155 There were five continent specific groups contained more than 75% Asian isolates,

156 supporting Asia to be the continent with the largest isolate source (n=319, 26.34%) (**Figure

157 2A**). North America was found to be the second largest group of isolate pool which consisted

158 six continent specific groups (n=132, 10.90%). Although less isolates were found sourced

159 from Europe (n=109, 9.00%), these isolates were distributed in nine continent specific

160 groups. Two groups (G16 & G29) of isolates were found as Oceania specific group (n=39,

161 3.22%) and three groups (G1 & G26 & G35) were found as South America specific groups

162 (n=109, 9.00%). In addition, the 12 mixed groups of isolates contained 226 isolates

163 (18.66%). Among all G level groups, G2 was the largest continent specific group (n=223)

164 that mainly contained isolates from Asia (193/223, 82.83%), while G35 was the second

165 largest continent specific group (n=109) that mainly contained isolates from South America

166 (99/109, 87.61%). Apart from all the continent groups above, there was no Africa specific

167 group found, but only with isolates collected from Africa defined in G28 (n=2), G37 (n=7)

168 and G29 (n=1) (**Figure 2**).

169 Although the continent specific groups did not 100% stick to one continent in our typing

170 system, the transmission events were still possible to trace. While most of the Asian isolates

171 fell into the Asia groups, a small proportion of the Asian isolates belonged to the mixed

172 groups. Similarly, most of the isolates sourced from North America and South America fell

173 in their own region groups, while a minority of the isolates were in the mixed groups.

174 Interestingly, isolates from Oceania and Europe could be found across all 12 mixed continent

175 groups, reflecting that *H. pylori* isolates from these two continents were relatively wide

176 spread of across the globe.

177 **The nested country level genomic typing for *H. pylori***

178 A total number of 859 isolates were grouped into 216 geographic patterns at country level,

179 which was predominant in 29 countries across six continents (**Figure 3**). Among these 29

180 countries encompassing 216 groups, 20 countries found in 168 groups were defined as

181 country-specific groups, while the rest 9 countries were scattered in the rest 48 country-level

182 mixed groups.

183 G35.C07 was the largest country specific group that contained 49 isolates from Colombia,

184 followed by the G35.C05 (n=35) dominated in Colombia as well. These isolates from

185 Colombia were mainly collected from the NCBI Bioproject PRJNA352848, which study

186 contained the population structure of *H. pylori* in regional evolution in South America

187 (Muñoz-Ramírez, et al. 2017). The isolates from group G35.C07 and G35.C05 were mainly

188 found from Colombia, Mexico and Spain (**Figure 3**). This result provided the evidence that

189 the *H. pylori* isolates were possibly transmitted from Spain and spread locally in South

190 America and North America. In comparison, Australia and Switzerland were the largest

191 countries of isolate source which isolates scattered across more than half of the country

192 specific groups.

193 When comparing the percentage of isolates from different countries, those isolates from

194 France, Germany, Malaysia, Nicaragua, Sweden, and UK were found scattered in more than

195 one continent group, while isolates from Cambodia, Colombia, India, Peru, Spain and US

196 were focused in one continent group when they were also found in other continent groups.

197  More importantly, Australia and Switzerland were two countries that mostly found with

198  scattered isolates in different regional specific groups.

199  Three clusters were observed in the percentage of different isolate sources at continent scale

200  (G32 to G25 with red branches in Figure 3), consisting of groups from Europe and mixed

201  continents. Specifically, those isolates from mixed groups were mainly sourced from

202  European and Oceania countries, making this cluster as Europe-Oceania dominated. The

203  second cluster was the mixed by Asian, Oceanian, European and mixed groups (G4 to G2

204  with green branches in Figure 3) but dominated by isolates from Australia and Asian

205  countries. Therefore, cluster two was specified as Asian-Pacific cluster. The third cluster was

206  formed by North American groups (G31 to G37 with purple branches in Figure 3), while

207  South American branches was next to the North America cluster.

**Comparing with seven-gene MLST**

209  Seven-gene MLST was implied to get sequence types (STs) of all 1,211 isolates.

210  Unfortunately, due to the high mutation rate of the *H. pylori* strains, most of the seven-gene

211  allele were only found with high similarity instead of an accurate type, as the result, a large

212  number of isolates (n=876, 72.3%) were untyped in our dataset (**Supplementary Table 1 &**

213  **Figure 1**). However, despite the most of undefined isolates, the typed isolates with exact ST

214  number would still be grouped closely by minimum spanning tree.

**A user-friendly typing website**

216  In order to support our *H. pylori* geographic typing tool, a user-friendly typing website was

217  established and available at https://db.cngb.org/HPTT/. Our HTTP approach is compatible

218  with any whole-genome sequencing (WGS) data with metadata (**Figure 4**). For the

219  sequencing data from pure-cultured isolates, the assembled genomes can be directly

220  submitted to our website. However, it is worth noting that sequences or assembled genomes

221  needed to be extracted from metagenome samples before submission (Parks, et al. 2017;

222  Olekhnovich, et al. 2019). Except for the sequenced genome data, the available assembled

223  contigs from NCBI Sequence Read Archive (SRA) or assembly database (RefSeq), or other

224  genome databases (e.g., European Nucleotide Achieve) can also be directly uploaded to our

225  website. By using MUMmer alignment and blast process, the uploaded genome can be

226  located to the closest matching genomes, further facilitating the possible transmission route

227  analysis across the globe. In addition, our database can be also linked to the NCBI genome

228    database, helping the user easily locate the metadata information from the available database

229    (**Supplementary materials**).

230    Except for the typing tool, the Nextstrain framework was also embedded in our website. By

231    clicking the uploaded genome number, information can be linked to the phylogenetic tree

232    with corresponding continent and country. Possible evolution relationships and interactive

233    located functions have made our typing tools easy to be applied and understood.

234    Ten genomes that newly uploaded to NCBI were downloaded and tested for the accuracy of

235    the typing method and the efficiency of our website (**Supplementary Table 2**). Since our

236    typing tool was established based on the MPS (Massive Parallel Sequencing) data, the first

237    genome (GCF_002206465.1) sequenced by Pacbio was failed to be assigned groups. The rest

238    nine genomes were typed successfully.

239

240    **Discussion**

241    The epidemiological patterns of *H. pylori* isolates have been reported with specific

242    geographic characteristics. In this study, the new typing webtool HPTT not only illustrated

243    the population structure of *H. pylori* but also made the genomic typing easy to approach. In

244    the continent level of typing, 1,112 isolates were grouped into 37 continent specific patterns.

245    Except for 12 continent mixed groups, the rest could be defined as continent specific groups

246    across the five continents. Isolates from Europe and Oceania were universally found in most

247    of the continent-level groups (Europe 33/37, 89.19% and Oceania 26/37, 70.27%),

248    illustrating that isolates from these two continents were widely spread across the world.

249    In the country level of typing, 1,045 isolates were grouped into 216 country level of groups.

250    Most of the isolates were defined as country specific groups (168/216, 77.77%), while the

251    rest of the isolates were grouped as country mixed groups (48/216, 22.22%). Australian and

252    Swiss isolates were found to be widespread around the world, while isolates from Columbia

253    was more regional specific. It has been reported that *H. pylori* in South America was

254    originally transmitted from Spain (Muñoz-Ramírez, et al. 2017), this data perfectly aligned

255    with our results in G35.C05 and G35.C07, giving the support of the accuracy of our genomic

256    typing method.

257  In this study, except for the novel typing tool, a user-friendly website was also established.

258  By using this typing tool, users can achieve fast and precise genomic typing, easily locating

259  the possible origins and transmission events across the world. When located in the actual

260  geographic group, it is easily for users to check the details of the corresponding composition

261  of the branches in our database. The genome with the highest identity can be easily linked to

262  the NCBI database as well as the visualization tool where the dynamic evolution of *H. pylori*

263  was shown. At the same time, seven-gene MLST results were displayed for each genome in

264  database.

265  The most interesting part of HPTT tool and methodology allow us to perform genome typing

266  with assembled genomes from the metagenomics samples, as illustrated in Figure 4. Due to

267  rapid mutation of *H. pylori*, it is most likely that the sample from one's gut are heterogeneity

268  in nature. The whole genome sequencing by combining sequencing libraries labelled with

269  different barcodes on a meta sample, and a cultured pure isolate could yield enough data from

270  one single run to perform the epidemiological surveillance of *H. pylori* on a global level to

271  find the origins in evolution profile. An open-source assay protocol will be developed and

272  shared in the future to combine with this HTTP tool to enable the epidemiological

273  surveillance of *H. pylori*.

274  Although our typing tool filled the gap of genetic epidemiological surveillance of *H. pylori*,

275  some of the functions still need to be improved. For example, cytotoxin-associated gene A

276  (*cagA*) and *vacA* were the two crucial genes that reported to be correlated with geographic

277  patterns of *H. pylori* (Yamaoka 2009; Breurec, et al. 2011). The *cagA* gene is one of the most

278  important virulence genes in *H. pylori*, located at the end of cag pathogenicity island (cag

279  PAI) that encodes 120–145 kDa CagA protein (Šterbenc, et al. 2019). Another virulence

280  factor was vacuolating cytotoxin encoded by the gene *vacA* (Šterbenc, et al. 2019). The

281  variation of these two genes were widely reported by the *H. pylori* groups that can reflect the

282  genomic different for different geographic patterns. However, such rapid typing method on a

283  website for these two genes are still lacking, which could be considered in the further HPTT

284  version 2.

285  *H. pylori* is normally treated by the antibiotics without antimicrobial susceptibility testing

286  (Pohl, et al. 2019). Antibiotics-resistant *H. pylori* has been reported related to several

287  mutations within the genes *pbp1A*, *23S rRNA*, *gyrA*, *rdxA*, *frxA*, and *rpoB* (Domanovich-

288  Asor, et al. 2021). In version 2, these antibiotics-resistant genes will be included in our

289    second version despite an antibiotic-resistant specific tool was available now (Yusibova, et

290    al. 2020). As more or more strains or isolates are deposited into our database with the

291    geographic information, the HPTT tool will be evolute into a more powerful tool to associate

292    the genomic typing information with its origin and phenotypes.

293    In summary, this work illustrates the efforts in global epidemiological study of *H. pylori*

294    isolates. Two functions were designed for the web typing tool, one for genomic typing and

295    the other for phylogenetic and geographic visualization. The accuracy of our genomic typing

296    system was proved by ten unused genomes as well as another published study (Muñoz-

297    Ramírez, et al. 2017). Together with the visualization tool, the genomic population structure

298    of *H. pylori* with geographic documents were described. Future studies based on this

299    approach will be expanded by the crucial virulence gene and antibiotic related genes. This

300    tool would be beneficial for the surveillance of *H. pylori* for public health and the monitoring

301    of its epidemic development.

302    **Acknowledgements**

306    **Data Availability**

307    All assembled *H. pylori* genomes were downloaded from NCBI assembly database.

308

309    **References**

310    Achtman M, Azuma T, Berg DE, Ito Y, Morelli G, Pan ZJ, Suerbaum S, Thompson SA, Van

311    Der Ende A, Van Doorn LJ. 1999. Recombination and clonal groupings within Helicobacter

312    pylori from different geographical regions. Molecular microbiology 32:459-470.

313    Agnew P, Koella JC. 1997. Virulence, parasite mode of transmission, and host fluctuating

314    asymmetry. Proceedings of the Royal Society of London. Series B: Biological Sciences

315    264:9-15.

316     Alm RA, Ling L-SL, Moir DT, King BL, Brown ED, Doig PC, Smith DR, Noonan B, Guild

317     BC, Dejonge BL. 1999. Genomic-sequence comparison of two unrelated isolates of the

318     human gastric pathogen Helicobacter pylori. Nature 397:176-180.

319     Breurec S, Guillard B, Hem S, Papadakos KS, Brisse S, Huerre M, Monchy D, Oung C,

320     Sgouras DN, Tan TS. 2011. Expansion of European vacA and cagA alleles to East-Asian

321     Helicobacter pylori strains in Cambodia. Infection, Genetics and Evolution 11:1899-1905.

322     Campbell DI, Warren BF, Thomas JE, Figura N, Telford JL, Sullivan PB. 2001. The African

323     enigma: low prevalence of gastric atrophy, high prevalence of chronic inflammation in West

324     African adults and children. Helicobacter 6:263-267.

325     Coll F, McNerney R, Guerra-Assunção JA, Glynn JR, Perdigao J, Viveiros M, Portugal I,

326     Pain A, Martin N, Clark TG. 2014. A robust SNP barcode for typing Mycobacterium

327     tuberculosis complex strains. Nature communications 5:1-5.

328     Didelot X, Nell S, Yang I, Woltemate S, Van der Merwe S, Suerbaum S. 2013. Genomic

329     evolution and transmission of Helicobacter pylori in two South African families. Proceedings

330     of the National Academy of Sciences 110:13880-13885.

331     Domanovich-Asor T, Craddock HA, Motro Y, Khalfin B, Peretz A, Moran-Gilad J. 2021.

332     Unraveling antimicrobial resistance in Helicobacter pylori: Global resistome meets global

333     phylogeny. Helicobacter:e12782.

334     Ernst PB, Gold BD. 2000. The disease spectrum of Helicobacter pylori: the

335     immunopathogenesis of gastroduodenal ulcer and gastric cancer. Annual Reviews in

336     Microbiology 54:615-640.

337     Group ES. 1993. Epidemiology of, and risk factors for, Helicobacter pylori infection among

338     3194 asymptomatic subjects in 17 populations. The EUROGAST Study Group. Gut 34:1672-

339     1676.

340     Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, Sagulenko P, Bedford T,

341     Neher RA. 2018. Nextstrain: real-time tracking of pathogen evolution. bioinformatics

342     34:4121-4123.

343     Jolley KA, Bray JE, Maiden MC. 2018. Open-access bacterial population genomics: BIGSdb

344     software, the PubMLST. org website and their applications. Wellcome open research 3.

345     Jolley KA, Maiden MC. 2010. BIGSdb: scalable analysis of bacterial genome variation at the

346     population level. BMC bioinformatics 11:1-11.

347     Kang J, Blaser MJ. 2006. Bacterial populations as perfect gases: genomic integrity and

348     diversification tensions in Helicobacter pylori. Nature Reviews Microbiology 4:826-836.

349   Kans J. 2020. Entrez direct: E-utilities on the UNIX command line. In. Entrez Programming

350   Utilities Help [Internet]: National Center for Biotechnology Information (US).

351   Kennemann L, Didelot X, Aebischer T, Kuhn S, Drescher B, Droege M, Reinhardt R, Correa

352   P, Meyer TF, Josenhans C. 2011. Helicobacter pylori genome evolution during human

353   infection. Proceedings of the National Academy of Sciences 108:5033-5038.

354   Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL. 2004.

355   Versatile and open software for comparing large genomes. Genome biology 5:R12.

356   Messenger SL, Molineux IJ, Bull J. 1999. Virulence evolution in a virus obeys a trade off.

357   Proceedings of the Royal Society of London. Series B: Biological Sciences 266:397-404.

358   Muñoz-Ramírez ZY, Mendez-Tenorio A, Kato I, Bravo MM, Rizzato C, Thorell K, Torres R,

359   Aviles-Jimenez F, Camorlinga M, Canzian F. 2017. Whole genome sequence and

360   phylogenetic analysis show Helicobacter pylori strains from Latin America have followed a

361   unique evolution pathway. Frontiers in cellular and infection microbiology 7:50.

362   Nguyen L-T, Schmidt HA, Von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective

363   stochastic algorithm for estimating maximum-likelihood phylogenies. Molecular biology and

364   evolution 32:268-274.

365   Olekhnovich EI, Manolov AI, Samoilov AE, Prianichnikov NA, Malakhova MV, Tyakht

366   AV, Pavlenko AV, Babenko VV, Larin AK, Kovarsky BA. 2019. Shifts in the human gut

367   microbiota structure caused by quadruple Helicobacter pylori eradication therapy. Frontiers

368   in microbiology 10:1902.

369   Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, Hugenholtz P,

370   Tyson GW. 2017. Recovery of nearly 8,000 metagenome-assembled genomes substantially

371   expands the tree of life. Nature microbiology 2:1533-1542.

372   Pohl D, Keller PM, Bordier V, Wagner K. 2019. Review of current diagnostic methods and

373   advances in Helicobacter pylori diagnostics in the era of next generation sequencing. World

374   journal of gastroenterology 25:4629.

375   Reeves PR, Liu B, Zhou Z, Li D, Guo D, Ren Y, Clabots C, Lan R, Johnson JR, Wang L.

376   2011. Rates of mutation and host transmission for an Escherichia coli clone over 3 years.

377   PloS one 6:e26907.

378   Salama NR, Gonzalez-Valencia G, Deatherage B, Aviles-Jimenez F, Atherton JC, Graham

379   DY, Torres J. 2007. Genetic analysis of Helicobacter pylori strain populations colonizing the

380   stomach at different times postinfection. Journal of bacteriology 189:3834-3845.

381   Šterbenc A, Jarc E, Poljak M, Homan M. 2019. Helicobacter pylori virulence genes. World

382   journal of gastroenterology 25:4870.

383    Suerbaum S, Michetti P. 2002. Helicobacter pylori infection. New England journal of

384    medicine 347:1175-1186.

385    Tomb J-F, White O, Kerlavage AR, Clayton RA, Sutton GG, Fleischmann RD, Ketchum KA,

386    Klenk HP, Gill S, Dougherty BA. 1997. The complete genome sequence of the gastric

387    pathogen Helicobacter pylori. Nature 388:539-547.

388    Yamaoka Y. 2009. Helicobacter pylori typing as a tool for tracking human migration.

389    Clinical Microbiology and Infection 15:829-834.

390    Yusibova M, Hasman H, Clausen PTLC, Imkamp F, Wagner K, Andersen LP. 2020. CRHP

391    Finder, a webtool for the detection of clarithromycin resistance in Helicobacter pylori from

392    whole-genome sequencing data. Helicobacter 25:e12752.

393    Zhou Z, Alikhan N-F, Sergeant MJ, Luhmann N, Vaz C, Francisco AP, Carriço JA, Achtman

394    M. 2018. GrapeTree: visualization of core genomic relationships among 100,000 bacterial

395    pathogens. Genome research 28:1395-1404.

396

397

398     **Table 1 Summary of 1,211 *H. pylor*i genomes**

| Continent | Country (region) of origin | Number of isolates |
|---|---|---|
| Asia | | 312 (25.76%) |
| | Cambodia | 53 |
| | China | 74 |
| | China (Taiwan) | 8 |
| | India | 47 |
| | Indonesia | 1 |
| | Japan | 31 |
| | Kuwait | 2 |
| | Malaysia | 79 |
| | North Korea | 1 |
| | Singapore | 14 |
| | South Korea | 1 |
| | Vietnam | 1 |
| Africa | | 10 (0.82%) |
| | Morocco | 6 |
| | Nigeria | 1 |
| | South Africa | 3 |
| Europe | | 294 (24.28%) |
| | Belarus | 2 |
| | Belgium | 6 |
| | France | 37 |
| | Germany | 31 |
| | Ireland | 1 |
| | Poland | 2 |
| | Portugal | 1 |
| | Russia | 3 |
| | Spain | 54 |
| | Sweden | 19 |
| | Switzerland | 130 |
| | United Kingdom | 8 |
| Oceania | | 178 (14.70%) |
| | Australia | 177 |
| | Papua New Guinea | 1 |
| North America | | 233 (19.24%) |
| | Canada | 2 |
| | El Salvador | 1 |
| | Mexico | 118 |
| | Nicaragua | 24 |
| | United States of America | 88 |
| South America | | 184 (15.19%) |

| Angola | 1 |
| Colombia | 172 |
| Peru | 11 |

399

400

401     **Figures and Figure legends**



402

403     **Figure 1. Two Clades of geographic typing based on the WGS.** The HPTT enrolled 1,211

404     *H. pylori* genomes downloaded from NCBI. The clade nodes in each figure are corresponding

405     to A) G groups for continent level of typing, B) the continent that isolate collected from, C) C

406     groups for country level of typing, D) the country that isolates collected from. Numbers in

407     parenthesis refer to the number of isolates in each genogroups.

408

409

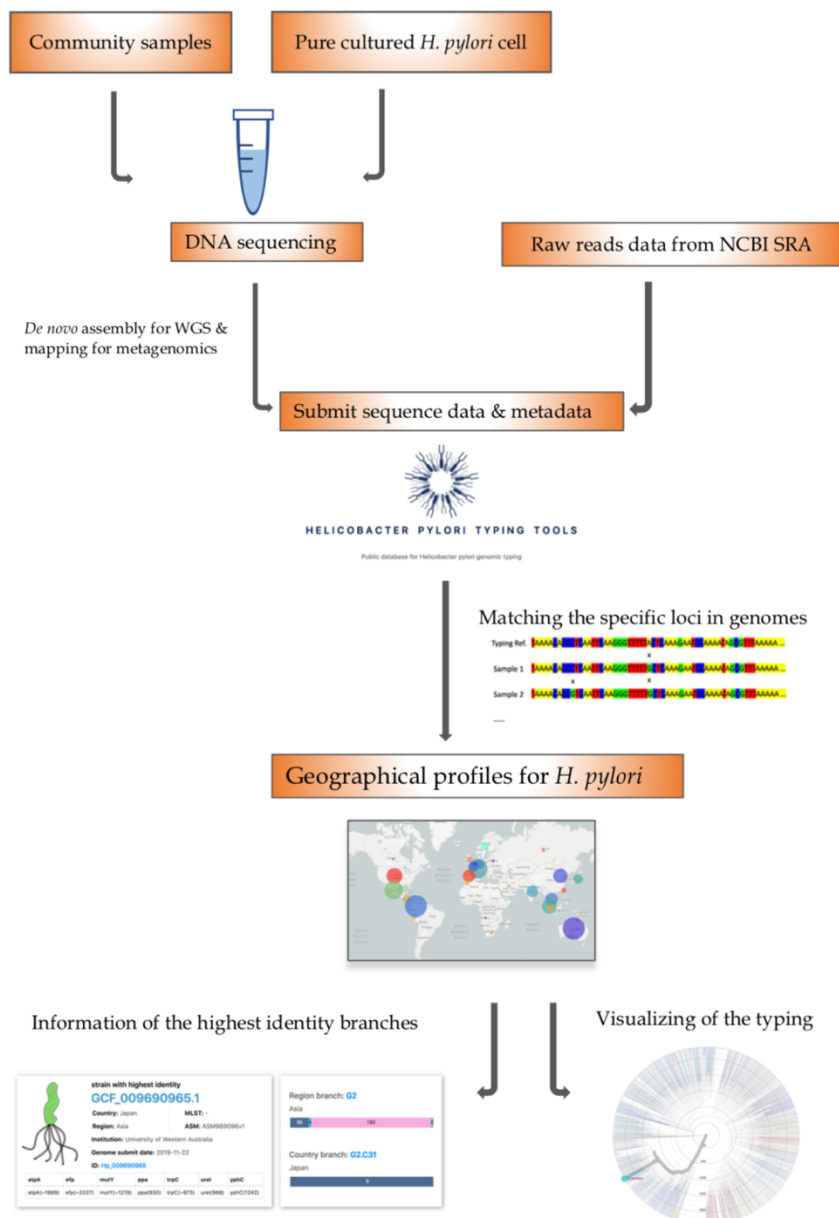| Africa | Asia | Europe | North America | Oceania | South America | |
|---|---|---|---|---|---|---|
| 0 |  | 2 | 0 | 4 | 0 | G3_Asia |
| 0 | 75 | 0 | 0 | 25 | 0 | G8_Asia |
| 0 | 79 | 5 | 0 | 16 | 0 | G11_Asia |
| 0 | 83 | 3 | 2 | 13 | 0 | G2_Asia |
| 0 | 82 | 0 | 0 | 18 | 0 | G4_Asia |
| 0 | 0 | 82 | 18 | 0 | 0 | G27_Europe |
| 0 | 0 |  | 0 | 0 | 0 | G17_Europe |
| 0 | 0 |  | 0 | 0 | 0 | G18_Europe |
| 0 | 0 |  | 0 | 7 | 0 | G23_Europe |
| 0 | 0 |  | 0 | 11 | 0 | G32_Europe |
| 0 | 0 |  | 0 | 10 | 0 | G12_Europe |
| 0 | 0 |  | 0 | 10 | 0 | G14_Europe |
| 0 | 0 | 81 | 0 | 19 | 0 | G10_Europe |
| 0 | 0 | 76 | 0 | 24 | 0 | G9_Europe |
| 0 | 0 | 70 | 0 | 30 | 0 | G19_Mix |
| 0 | 9 | 63 | 6 | 22 | 0 | G25_Mix |
| 0 | 10 | 50 | 10 | 30 | 0 | G13_Mix |
| 0 | 7 | 56 | 15 | 22 | 0 | G24_Mix |
| 0 | 14 | 50 | 0 | 36 | 0 | G15_Mix |
| 0 | 0 | 56 | 0 | 44 | 0 | G21_Mix |
| 0 | 33 | 22 | 0 | 44 | 0 | G6_Mix |
| 0 | 36 | 45 | 0 | 18 | 0 | G7_Mix |
| 20 | 0 | 20 | 10 | 40 | 10 | G28_Mix |
| 0 | 0 | 34 | 9 | 57 | 0 | G20_Mix |
| 0 | 0 | 25 | 0 | 75 | 0 | G16_Oceania |
| 4 | 0 | 21 | 0 | 75 | 0 | G29_Oceania |
| 0 | 0 | 8 | 17 | 0 | 75 | G1_South_America |
| 0 | 0 | 11 | 7 | 0 | 82 | G26_South_America |
| 0 | 0 | 4 | 9 | 0 | 88 | G35_South_America |
| 0 | 0 | 0 | 64 | 0 | 36 | G34_Mix |
| 5 | 0 | 8 | 52 | 13 | 23 | G37_Mix |
| 0 | 24 | 0 | 76 | 0 | 0 | G5_North_America |
| 0 | 0 | 13 | 87 | 0 | 0 | G31_North_America |
| 0 | 0 | 10 | 83 | 7 | 0 | G22_North_America |
| 0 | 0 | 6 | 88 | 6 | 0 | G30_North_America |
| 0 | 0 | 6 | 88 | 0 | 6 | G33_North_America |
| 0 | 0 | 5 |  | 0 | 0 | G36_North_America |

410 **Figure 2. Geographical clustering of *H. pylori* continent clades.** The number in each cube

411 represents the percentage of unique isolates sourced from each of the continents. A total

412 number of 37 continent level of groups were defined. The deeper the colour, the higher the

413 percentage of the isolates in that continent level of clade groups. Also, a phylogenetic tree is

414 shown in the left side of the table. The background information of isolates is provided in

415 Supplementary Table 1.

416

**Figure 3. Geographical clustering of *H. pylori* country subclades.** The number in each cube represents the percentage of unique isolates sourced from each of the country in that continent groups. A total number of 216 country level of groups were defined. The deeper the colour, the higher the percentage of the isolates sourced from that country in continent level of groups. The background information of isolates is provided in Supplementary Table 1.

422

**Figure 4. The HPTT workflow.** The SNP based genotyping approach can be used with the Whole Genome Sequencing (WGS) data, which can be acquired in following ways: DNA can be extracted from a pure cultured bacterial cell with WGS data or a community sample with metagenomic sequencing data. After being sequenced by an appropriate platform, the assembled genomes can be directly submitted to our database. In addition, the public assembled data also can be directly submitted to our database. The downstream analyses of the aligned sequence data can be linked to the phylogenetic and geographic page.

**A** — Continent group (Ggroup)

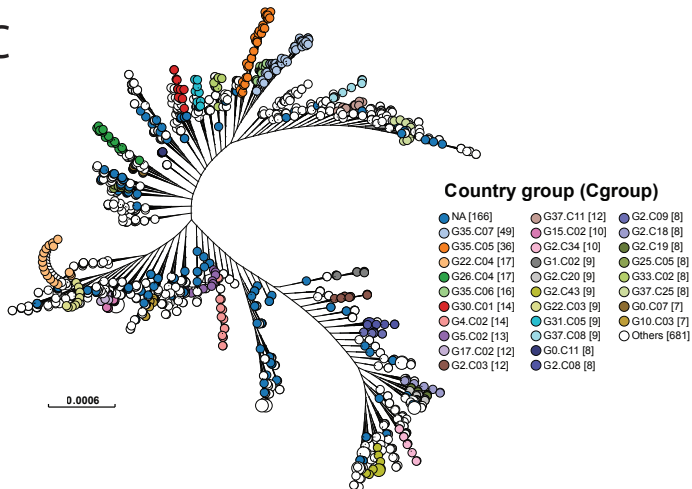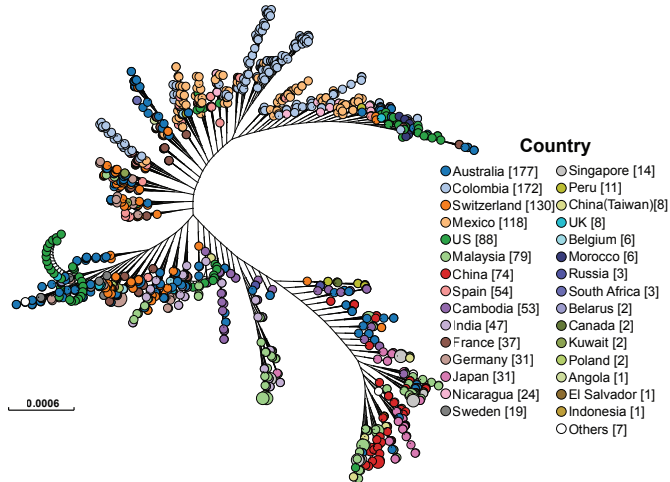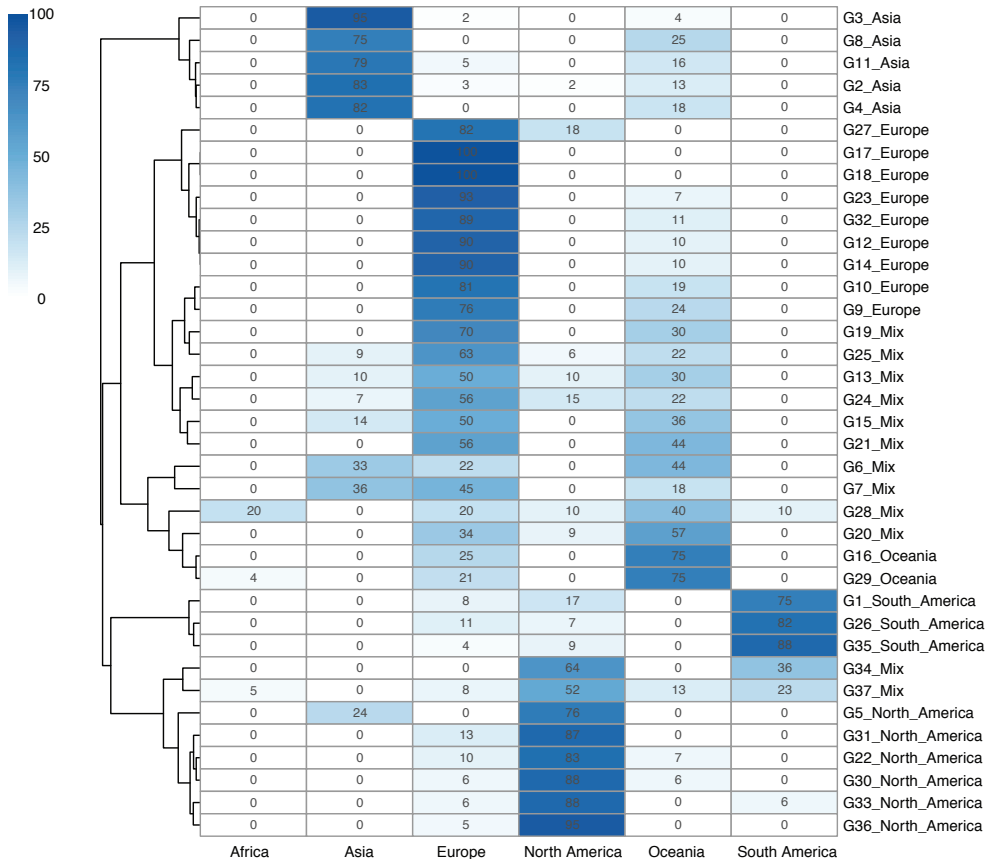| | | |
|---|---|---|
| G2 [233] | G11 [19] | G34 [11] |
| G37 [141] | G36 [19] | G7 [11] |
| G35 [113] | G17 [17] | G12 [10] |
| NA [99] | G4 [17] | G13 [10] |
| G3 [57] | G5 [17] | G14 [10] |
| G25 [54] | G10 [16] | G19 [10] |
| G20 [35] | G16 [16] | G28 [10] |
| G22 [30] | G30 [16] | G21 [9] |
| G9 [29] | G33 [16] | G32 [9] |
| G26 [28] | G15 [14] | G6 [9] |
| G24 [27] | G23 [14] | G8 [8] |
| G29 [24] | G1 [12] | G18 [7] |
| G31 [23] | G27 [11] | |

0.0006

**B** — Continent

- Asia [312]
- Europe [294]
- North America [233]
- South America [184]
- Oceania [178]
- Africa [10]

0.0006

**C** — Country group (Cgroup)

| | | |
|---|---|---|
| NA [166] | G37.C11 [12] | G2.C09 [8] |
| G35.C07 [49] | G15.C02 [10] | G2.C18 [8] |
| G35.C05 [36] | G2.C34 [10] | G2.C19 [8] |
| G22.C04 [17] | G1.C02 [9] | G25.C05 [8] |
| G26.C04 [17] | G2.C20 [9] | G33.C02 [8] |
| G35.C06 [16] | G2.C43 [9] | G37.C25 [8] |
| G30.C01 [14] | G22.C03 [9] | G0.C07 [7] |
| G4.C02 [14] | G31.C05 [9] | G10.C03 [7] |
| G5.C02 [13] | G37.C08 [9] | Others [681] |
| G17.C02 [12] | G0.C11 [8] | |
| G2.C03 [12] | G2.C08 [8] | |

0.0006

**D** — Country

| | |
|---|---|
| Australia [177] | Singapore [14] |
| Colombia [172] | Peru [11] |
| Switzerland [130] | China(Taiwan)[8] |
| Mexico [118] | UK [8] |
| US [88] | Belgium [6] |
| Malaysia [79] | Morocco [6] |
| China [74] | Russia [3] |
| Spain [54] | South Africa [3] |
| Cambodia [53] | Belarus [2] |
| India [47] | Canada [2] |
| France [37] | Kuwait [2] |
| Germany [31] | Poland [2] |
| Japan [31] | Angola [1] |
| Nicaragua [24] | El Salvador [1] |
| Sweden [19] | Indonesia [1] |
| | Others [7] |

0.0006

| | Africa | Asia | Europe | North America | Oceania | South America | |
|---|---|---|---|---|---|---|---|
| | 0 | 35 | 2 | 0 | 4 | 0 | G3_Asia |
| | 0 | 75 | 0 | 0 | 25 | 0 | G8_Asia |
| | 0 | 79 | 5 | 0 | 16 | 0 | G11_Asia |
| | 0 | 83 | 3 | 2 | 13 | 0 | G2_Asia |
| | 0 | 82 | 0 | 0 | 18 | 0 | G4_Asia |
| | 0 | 0 | 82 | 18 | 0 | 0 | G27_Europe |
| | 0 | 0 | 100 | 0 | 0 | 0 | G17_Europe |
| | 0 | 0 | 100 | 0 | 0 | 0 | G18_Europe |
| | 0 | 0 | 93 | 0 | 7 | 0 | G23_Europe |
| | 0 | 0 | 89 | 0 | 11 | 0 | G32_Europe |
| | 0 | 0 | 90 | 0 | 10 | 0 | G12_Europe |
| | 0 | 0 | 90 | 0 | 10 | 0 | G14_Europe |
| | 0 | 0 | 81 | 0 | 19 | 0 | G10_Europe |
| | 0 | 0 | 76 | 0 | 24 | 0 | G9_Europe |
| | 0 | 0 | 70 | 0 | 30 | 0 | G19_Mix |
| | 0 | 9 | 63 | 6 | 22 | 0 | G25_Mix |
| | 0 | 10 | 50 | 10 | 30 | 0 | G13_Mix |
| | 0 | 7 | 56 | 15 | 22 | 0 | G24_Mix |
| | 0 | 14 | 50 | 0 | 36 | 0 | G15_Mix |
| | 0 | 0 | 56 | 0 | 44 | 0 | G21_Mix |
| | 0 | 33 | 22 | 0 | 44 | 0 | G6_Mix |
| | 0 | 36 | 45 | 0 | 18 | 0 | G7_Mix |
| | 20 | 0 | 20 | 10 | 40 | 10 | G28_Mix |
| | 0 | 0 | 34 | 9 | 57 | 0 | G20_Mix |
| | 0 | 0 | 25 | 0 | 75 | 0 | G16_Oceania |
| | 4 | 0 | 21 | 0 | 75 | 0 | G29_Oceania |
| | 0 | 0 | 8 | 17 | 0 | 75 | G1_South_America |
| | 0 | 0 | 11 | 7 | 0 | 82 | G26_South_America |
| | 0 | 0 | 4 | 9 | 0 | 88 | G35_South_America |
| | 0 | 0 | 0 | 64 | 0 | 36 | G34_Mix |
| | 5 | 0 | 8 | 52 | 13 | 23 | G37_Mix |
| | 0 | 24 | 0 | 76 | 0 | 0 | G5_North_America |
| | 0 | 0 | 13 | 87 | 0 | 0 | G31_North_America |
| | 0 | 0 | 10 | 83 | 7 | 0 | G22_North_America |
| | 0 | 0 | 6 | 88 | 6 | 0 | G30_North_America |
| | 0 | 0 | 6 | 88 | 0 | 6 | G33_North_America |
| | 0 | 0 | 5 | 95 | 0 | 0 | G36_North_America |