



The *Helicobacter pylori* Genome Project: insights into *H. pylori* population structure from analysis of a worldwide collection of complete genomes

Received: 5 September 2023

Accepted: 13 November 2023

Published online: 11 December 2023

Check for updates

Kaisa Thorell^{1,204}✉, Zilia Y. Muñoz-Ramírez^{2,204}, Difei Wang^{3,4}, Santiago Sandoval-Motta^{5,6,7}, Rajiv Boscolo Agostini⁸, Silvia Ghirotto⁸, Roberto C. Torres⁹, HpGP Research Network*, Daniel Falush⁹, M. Constanza Camargo^{4,205} & Charles S. Rabkin^{4,205}

Helicobacter pylori, a dominant member of the gastric microbiota, shares co-evolutionary history with humans. This has led to the development of genetically distinct *H. pylori* subpopulations associated with the geographic origin of the host and with differential gastric disease risk. Here, we provide insights into *H. pylori* population structure as a part of the *Helicobacter pylori* Genome Project (HpGP), a multi-disciplinary initiative aimed at elucidating *H. pylori* pathogenesis and identifying new therapeutic targets. We collected 1011 well-characterized clinical strains from 50 countries and generated high-quality genome sequences. We analysed core genome diversity and population structure of the HpGP dataset and 255 worldwide reference genomes to outline the ancestral contribution to Eurasian, African, and American populations. We found evidence of substantial contribution of population hpNorthAsia and subpopulation hspUral in Northern European *H. pylori*. The genomes of *H. pylori* isolated from northern and southern Indigenous Americans differed in that bacteria isolated in northern Indigenous communities were more similar to North Asian *H. pylori* while the southern had higher relatedness to hpEastAsia. Notably, we also found a highly clonal yet geographically dispersed North American subpopulation, which is negative for the *cag* pathogenicity island, and present in 7% of sequenced US genomes. We expect the HpGP dataset and the corresponding strains to become a major asset for *H. pylori* genomics.

Helicobacter pylori has co-existed with humans for more than 100,000 years. It is the primary etiologic agent associated with gastric diseases such as ulcers and gastric cancer. Still, while over half the world population is colonized with *H. pylori*, less than 2% will end up with gastric cancer^{1,2}. The intimate symbiotic relationship with humans,

together with predominantly vertical transmission, has led *H. pylori* to evolve into multiple distinct geographic populations^{3–5}. The phylogeographic structure of *H. pylori* is classified into major populations (“hp”) and subpopulations (“hsp”) that correlate with ancient human migrations^{3,4,6}. However, most worldwide efforts in this regard have

A full list of affiliations appears at the end of the paper. *A list of authors and their affiliations appears at the end of the paper. ✉e-mail: kaisa.thorell@gu.se

been based on the analysis of only a handful of genes rather than whole genomes^{3,4}. The risk of developing disease from *H. pylori* infection varies greatly by geography⁷ and genomic studies of both humans and *H. pylori* are required to identify the factors that modify this risk.

The *Helicobacter pylori* Genome Project (*HpGP*) is an international and multidisciplinary initiative to sequence and map *H. pylori* population structure by collecting strains worldwide. Here we analyze 1011 *H. pylori* genomes, sequenced with PacBio Single Molecule, Real-Time long-read technology, which made it possible to acquire complete assemblies. By relating the *HpGP* dataset to a reference set of known population assignment, we were able to quantify, with great resolution, the different inferred ancestral sources of *H. pylori* subpopulations and the recent and ongoing admixture among subpopulations.

Results

HpGP is a dataset of high quality and worldwide representation

The *HpGP* has assembled clinical strains from 50 countries, including 12 countries from which no *H. pylori* genome sequences have previously been published (Table 1). Out of the 1011 genomes, all but seven were completely circularized (Supplementary Data 1).

To investigate the population structure of the *HpGP* dataset, we performed fineSTRUCTURE (FS), chromosome painting, and network analyses of shared core genome features as described⁸, and discriminant analysis of principal components (DAPC)⁹. To anchor the dataset, we used 255 *H. pylori* reference genomes with known Hp/hsp population assignments, representing 17 global subpopulations (Supplementary Data 2). In total, the core genome (set of homologous genes present in >95% of genomes) of the *HpGP* dataset, the *HpGP*-26695 reference genome, and 255 worldwide references, consisted of 1227 genes.

The fineSTRUCTURE global analysis revealed four main *H. pylori* population clusters: (i) Southwest Europe, including Latin America and Northeast Africa, (ii) Northern and Central Europe, Middle East, and Central Asia, (iii) Western and Southern Africa, including Africa2 and North, South and Central America, and (iv) North, Central and East Asia, and Indigenous populations in America. In total, these formed 17 main subpopulations (Fig. 1 and Supplementary Figs. 1 and 2). The network and DAPC analyses supported this structure but with six main clusters of differentiation (Fig. 2 and Supplementary Fig. 3).

South Africa (DAPC group 4, hpAfrica2 in FS) and the reference genomes from Australia/New Guinea (DAPC group 6, hpSahul in FS) differentiated extensively from the others. A further DAPC analysis not considering these two groups showed a clear separation of two of the remaining clusters from the others: one composed of isolates of African and American origin (DAPC group 2, FS cluster III) and one that includes isolates from Central/East Asia and Indigenous Americans (DAPC group 5, FS cluster IV). The remaining (groups 1 and 3) were more similar and intertwined, representing Southern Europe/Northeast Africa and Eurasia/Central Asia and Americas, respectively (Supplementary Fig. 3). The population assignments according to the respective analyses are summarized in Supplementary Data 3.

The hpEurope subpopulations span from the Atlantic coast to South Asia

In the fineSTRUCTURE analysis, three main European/Eurasian subpopulations emerged (Supplementary Figs. 1 and 2), of which hspNEurope and hspSWEurope have previously been described^{8,10}. The hspEurasia population is proposed in this study, and includes the already reported hspCEurope/hspSEurope^{8,10–12}, and hspMiddleEast¹⁰. Previous studies had limited coverage of Eastern Europe and the Middle East. The *HpGP* strains from Lithuania, Latvia, Russia, Poland, Bulgaria, Türkiye, and Jordan allowed mapping of the Eurasian *H. pylori* relationships with unprecedented detail (Fig. 1). Two northern European populations showed an east-west differentiation, in

Table 1 | Summary of the *HpGP* strain collection

Country	Total number	Non-atrophic gastritis (%)	Intestinal metaplasia (%)	Gastric cancer (%)
Algeria ^a	10	100		
Argentina	10	100		
Bangladesh	10	100		
Brazil	21	48	38	14
Bulgaria ^a	8	100		
Canada	20	35	65	
Chile	46	54	46	
China	10			100
DR Congo ^a	11	91		9
Colombia	45	78	16	7
Costa Rica	8	100		
Dominican Republic ^a	11	91		9
France	21	48		52
Germany	17	59		41
Ghana ^a	2	100		
The Gambia	5	100		
Greece	21	48		52
Guatemala	3	100		
Honduras	26	35	38	27
Indonesia	11	91		9
India	10	100		
Iran	4	100		
Iceland ^a	11	91	9	
Israel	10	70	30	
Italy	29	34	34	31
Japan	29	38	21	41
Jordan ^a	10	100		
Kazakhstan ^a	2	100		
Kyrgyzstan ^a	10	100		
Korea	54	19	19	63
Latvia ^a	34	29	24	47
Lithuania	23	43	35	22
Malaysia	19	47	53	
Mexico	22	45		55
Myanmar ^a	12	83		17
Nepal	13	77		23
Nigeria	4	100		
Peru	33	30	24	45
Poland ^a	20	100		
Portugal	30	57	27	17
Russia	10	60		40
Singapore	21	38	33	29
South Africa	9	100		
Spain	106	72	13	15
Sweden	30	33	33	33
Switzerland	15	60	40	
Taiwan	24	42		58
Türkiye	17	59	18	24
US (continental)	68	96	1	3
US (Puerto Rico) ^a	7	100		
Vietnam	9			100
Total	1011	60	17	23

^aGeographical areas from which no *H. pylori* whole-genome sequences were previously available in GenBank

hspNEurope an east clade with genomes from Latvia, Lithuania, and Russia separated from a north-western clade with genomes from UK, Sweden, Iceland, and Canada (Supplementary Fig. 2). Within hspEurasia, three main clades could be noted of which two spanned from west

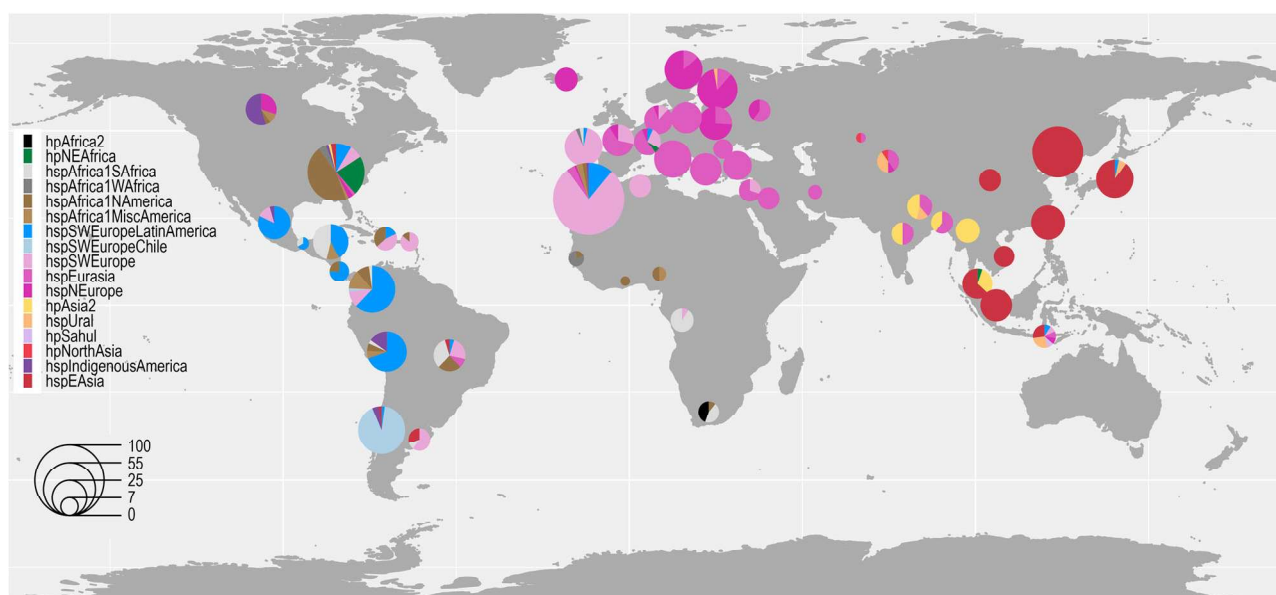


Fig. 1 | World map of *HpGP* strain origins and population assignments. The area of each pie is proportional to the number of *HpGP* genomes from each country and colored by the *H. pylori* population (hp) and subpopulation (hsp) as assigned by fineSTRUCTURE (Supplementary Figs. 1 and 2).

to east (Supplementary Fig. 2). The first, Central-Eastern European hspEurasia1, dominates in Germany, Poland, Lithuania, Latvia, Türkiye, and Russia, while hspEurasia2 is more Southern with representation from France in the west, via Italy and Greece, to Jordan and Iran in the Middle East. Thirdly, hspEurasia3 includes genomes from India and Bangladesh, but also Greece, which separated from the others but were still within the hspEurasia subpopulation.

The European subpopulations have different ancestry proportions

To further investigate the proposed subpopulations we inferred ancestry by comparing genomes within our contemporary dataset in a directed chromosome painting using only the proposed *H. pylori* ancestral populations hpAfrica2, hpNEAfrica, hspAfrica1WAfrica, hpAsia2, hspUral, hpNorthAsia, and hspEAsia as donors (i.e., contributors of genomic ancestry)^{3,10}. We confirmed a gradient in inferred ancestry along both the north-south axis with increasing Asian ancestry and decreasing African ancestry in the hspEurasia1 and hspNEurope populations and the east-west axis with hspSWEurope having a higher proportion of hspAfrica1WAfrica ancestry and with the similar contribution of hpNEAfrica as the Eurasia2 population (Fig. 3 and Supplementary Fig. 3).

The more central Asian hspEurasia3 on the other hand, showed markedly higher hpAsia2 ancestry than the other hpEurope populations, concordant with its geographical co-existence with hpAsia2. Interestingly, hspUral was a more pronounced Asian ancestor for all the hpEurope subpopulations than hpNorthAsia and hspEAsia, the latter two being very even contributors, except for in hspNEurope, where hpNorthAsian ancestry was slightly higher. This relationship was also supported by the network analysis (Fig. 2).

Central Asia can be described with increased resolution but still has underrepresented regions

Apart from the relatively well-investigated hspEAsian subpopulation¹³, the fineSTRUCTURE analysis grouped the central Eurasian strains into three main clades: hpAsia2 and two clades preliminarily termed hpNorthAsia and hspUral, based on their association with reference strains previously described by Moodley et al.¹⁴.

HpAsia2 is one of the main ancestral populations of *H. pylori* but has been comparatively understudied. In the *HpGP* dataset, genomes

belonging to hpAsia2 are mainly from India, Bangladesh, Myanmar, and Nepal, with the Nepalese forming a clade slightly separated from the others (Supplementary Fig. 2a, c). As seen in Fig. 1, hpAsia2 co-exists with the hpEurope hspEurasia3 population in all these countries, except for Myanmar, where only hpAsia2 is present. The DAPC analysis, on the other hand, did not distinguish hpAsia2 from hspNEurope and hspEurasia using $k=6$, while the separation was evident and very consistent using $k=17$ (Supplementary Fig. 4).

HpNorthAsia was previously established as one of the main Siberian populations using Multilocus Sequence Typing (MLST)¹⁴. In our reference panel, hpNorthAsia (including hspAltai) and its subpopulation hspSiberia1 were represented by genomes from central and eastern Siberia. In our analyses, these two populations did not segregate, and *HpGP* genomes from Kazakhstan and Kyrgyzstan were also associated with this cluster (Fig. 1 and Supplementary Fig. 1).

hspUral has been suggested as a southern central Asian subpopulation of hpAsia2. In our dataset, a cluster with a relatively wide geographical representation from Kazakhstan and Kyrgyzstan to Indonesia and Japan (Fig. 1 and Supplementary Figs. 1 and 2) is associated with the hspUral reference genomes. Our main chromosome painting analysis suggested the proposed hspUral population to contain two subclades with very different painting profiles (Supplementary Fig. 2), which was supported by the DAPC and network analysis, and the fineSTRUCTURE principal component analysis (PCA) (see https://hpgp.shinyapps.io/Interactive_figures, Fig. 4). The ancestral contributions to the central Asian genomes confirmed the *HpGP* “hspUral” clade not to have pronounced contribution by the hspUral references but relatively high hpAsia2, hpNorthAsia and hspEAsia painting proportions (Fig. 3). The variability of contributions was also high within the clade, suggesting this may not constitute one pure subpopulation but may consist of representatives of several HpAsia subpopulations (https://hpgp.shinyapps.io/Interactive_figures, Fig. 2). One hpAsia2 reference genome, L7, from Ladakh in northern India grouped with this cluster, especially close to two Nepalese genomes. Several “hspUral” genomes also showed an association with hpSahul in the chromosome painting (Supplementary Fig. 2), which may indicate a relationship between this group and the recently suggested hpRyukyu¹⁵.

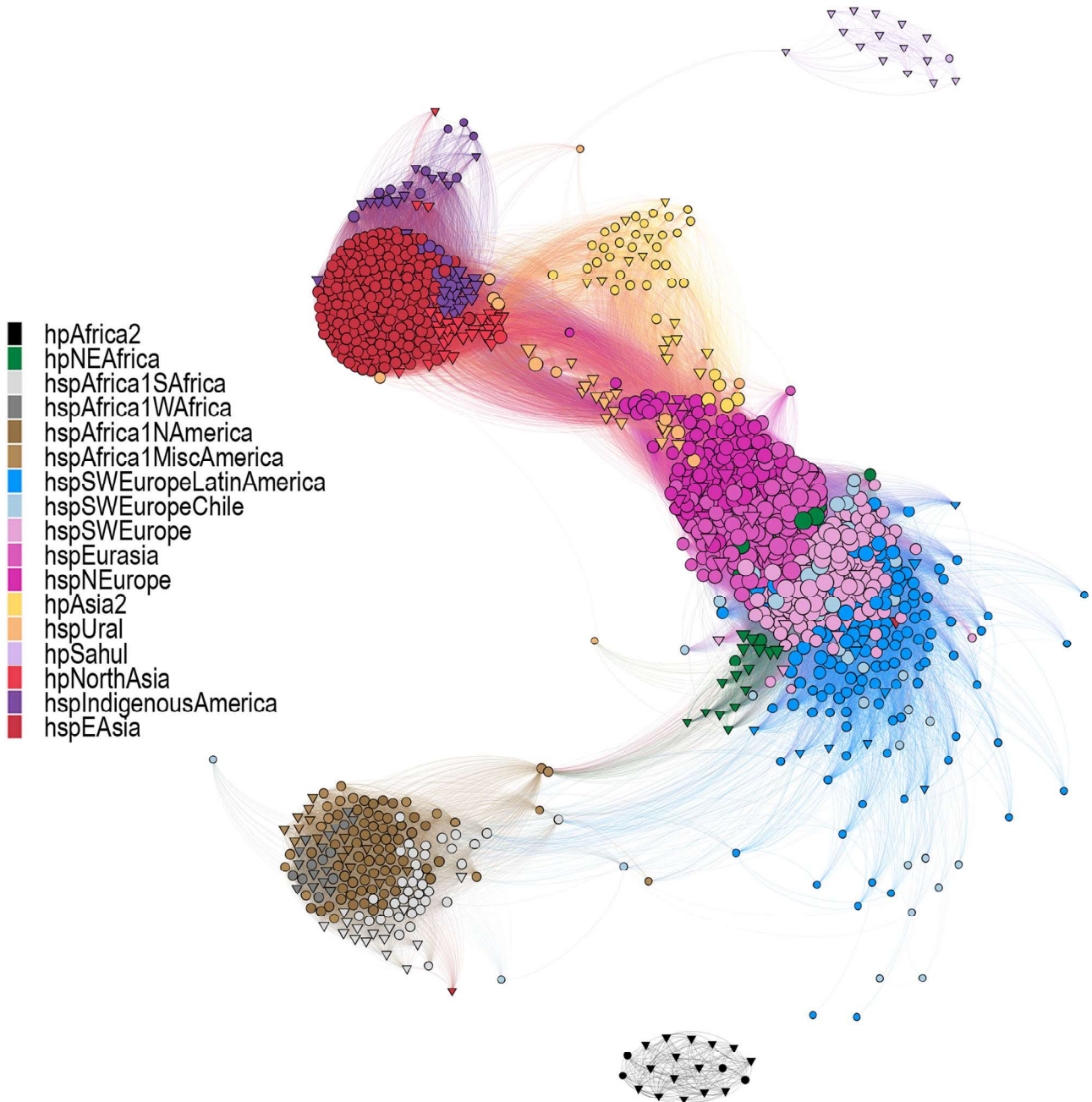


Fig. 2 | Distance network analyses of the core genome of the *H. pylori* strains studied. Fruchterman–Reingold layout of the pruned distance network between *HpGP* genomes (circles) and reference genomes (triangles) (see Methods). Colors indicate the *H. pylori* population (hp) and subpopulation (hsp) as assigned by fineSTRUCTURE (Supplementary Figs. 1 and 2). The length and opacity of each link

are proportional to the genetic distance between genomes (nodes), with higher opacity and shorter length indicating genetic closeness and less opacity and higher length indicating higher genetic distance between strains. The size of each node is proportional to the connectivity (number of links) of that node, indicating that bigger nodes have connections to more other strains than those of lesser sizes.

African and African-descent genomes

The *HpGP* dataset includes African genomes from understudied countries such as Algeria, Democratic Republic of Congo (DRC), Ghana, and Nigeria, adding to previous knowledge from the Gambia and South Africa. The fineSTRUCTURE analysis confirms earlier observations of the presence of four African populations in this continent, hpAfrica2, in the *HpGP* dataset represented in South Africa; hspAfrica1SAfrica, which reaches as far north as DRC; hspAfrica1WAfrica represented in the Gambia, as previously reported, and hpNEAfrica. However, the Ghanaian and Nigerian genomes grouped with the more admixed hspAfrica1NorthAmerica and hspAfrica1MiscAmericas populations, interspersed with, and by chromosome painting indistinguishable from

genomes from the US, Puerto Rico (US territory), Dominican Republic, Colombia, and Brazil, likely a result of the trans-Atlantic slave trade from West Africa into the Americas.

The East African reference genomes from Sudan and Ethiopia grouped within the hspSWEurope umbrella but distinctive from the European SWEurope clade, instead forming a cluster with North American genomes (Fig. 1 and Supplementary Fig. 1). In the fineSTRUCTURE PCA plots, especially pronounced in PC10 and PC11, the reference hpNEAfrican genomes and the US *HpGP* genomes clearly formed two segregated groups except for one Malaysian and one Swiss genome that grouped with the references (https://hpgp.shinyapps.io/Interactive_figures, Fig. 5). Both the DAPC and network analysis

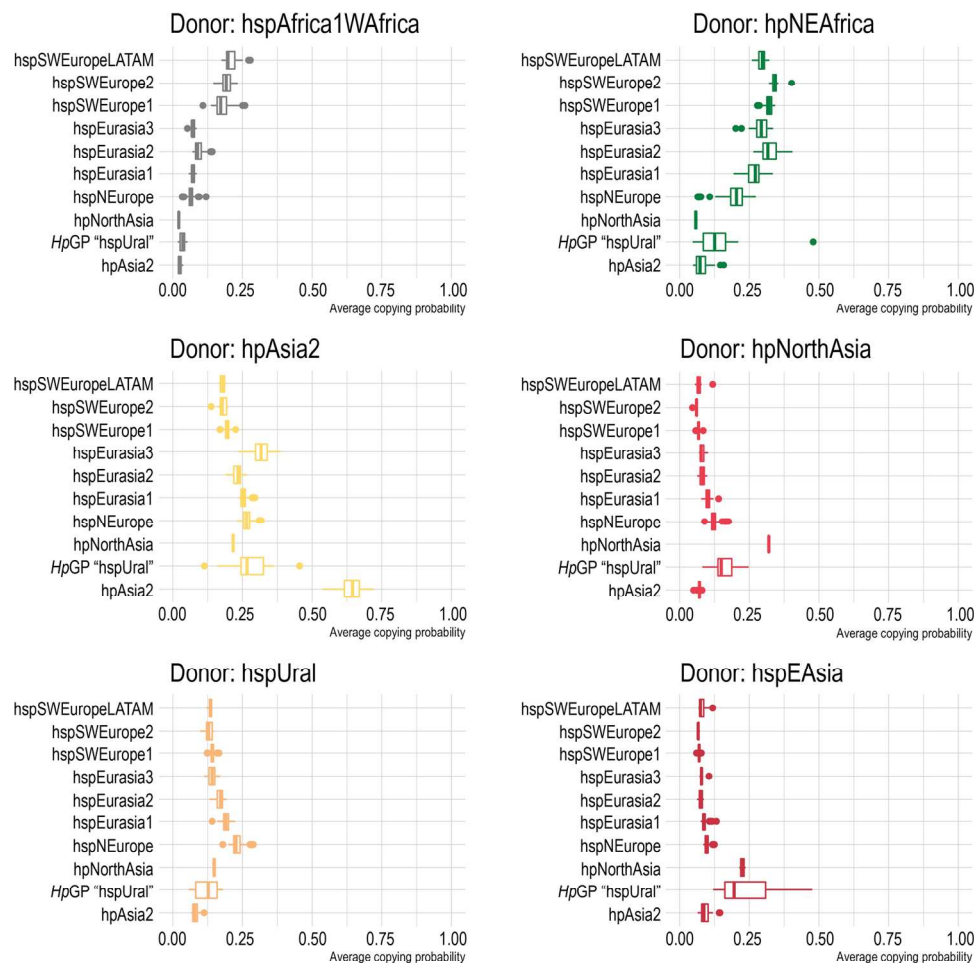


Fig. 3 | Inferred ancestral genomic contributions to the Eurasian *HpGP* genomes. Ancestral chromosome painting proportions by donor and Eurasian subpopulation. Boxplots show the median value per group, and the 25th and 75th percentiles (hinges), with whiskers extending from the hinge to the largest value no further than $1.5 \times$ IQR (inter-quartile range) from the hinge. Data points beyond the

whiskers are plotted individually. The number of genomes in each respective Eurasian population is hspSWEuropeLatinAmerica, $n = 15$; hspSWEurope2, $n = 12$; hspSWEurope1, $n = 129$; hspEurasia3, $n = 18$; hspEurasia2, $n = 76$; hspEurasia1, $n = 103$; hspNEurope, $n = 95$; hpNorthAsia, $n = 2$; *HpGP* "hspUral", $n = 10$; hpAsia2, $n = 27$.

supported the separation of the hpAfrica1 population from hpNEAfrica, the latter being intermingled with genomes from southern Europe and Iberia.

The Algerian strains did not cluster with the other African strains, but within hspSWEurope, together with genomes from Israel and Colombia in a cluster we termed SWEurope2. Despite showing slightly higher West and Northeast African and lower Asian ancestry than SWEurope1 (https://hpgp.shinyapps.io/Interactive_figures, Fig. 2), our analysis confirmed that North African *H. pylori* more closely resemble Iberian and Middle Eastern bacteria than African bacteria.

North America hosts a geographically dispersed deep clone

The *HpGP* dataset contains 68 genomes from the wide geographical representation of the continental US. This feature allowed us to identify a novel subpopulation of 15 US isolates, which showed high similarity and clustered together with genomes of hpNEAfrican ancestry in the fineSTRUCTURE analysis (Supplementary Figs. 1 and 2) and of which none carried the *cagPAI*.

High levels of sequence homogeneity within *H. pylori* are unexpected as unrelated strains differ in their DNA sequence at almost all genes. To further investigate the novel US subpopulation, we performed core genome (cg) MLST of the entire dataset (Fig. 4a). Within the *HpGP*, over 64% of strain pairs differ in sequence at all the 1040 genes. Even amongst strains sampled from the same country, 34%

differ in all the genes. Only 0.15%, 798 pairs, shared similarity at >1% of genes. All but 213 of these pairs are between strains in the same country. Nearly a tenth (66) of these pairs is found between a group of 12 US strains, showing allele distances between 0.83 and 0.94 (17–6% identical alleles, respectively). Thus, this group represents older clonal relationships, a putative "deep clone"; a set of strains that share a recent common ancestor but have diverged via homologous recombination at a large fraction of their genome. Three strains are somewhat less related to these 12, sharing between 1% and 7% of genes, and were conservatively excluded from this clonal group. Other pairs involving more than two samples from the same population also showed deep clonal relationships (e.g., hspSWEuropeChile). However, the amount and pattern of alleles shared between these samples could be better explained by genetic drift and further analysis within this population is needed to define the boundaries of a putative clone.

The *HpGP* strains from the deep clonal group were sampled from California, Wisconsin, Tennessee, Arkansas, Georgia, and Texas and, in total, represented a fifth of the *HpGP* US genomes. Kmer-based clustering analysis showed an additional five public genomes from two other geographical sources, Ohio and Louisiana, associating closely with the proposed deep clonal group. We used ClonalFrameML to estimate the relationships between the genomes. Assuming a previously estimated 1.38×10^{-5} mutation rate per site per year¹⁶, the common ancestor lived an estimated 175 years before the strains were

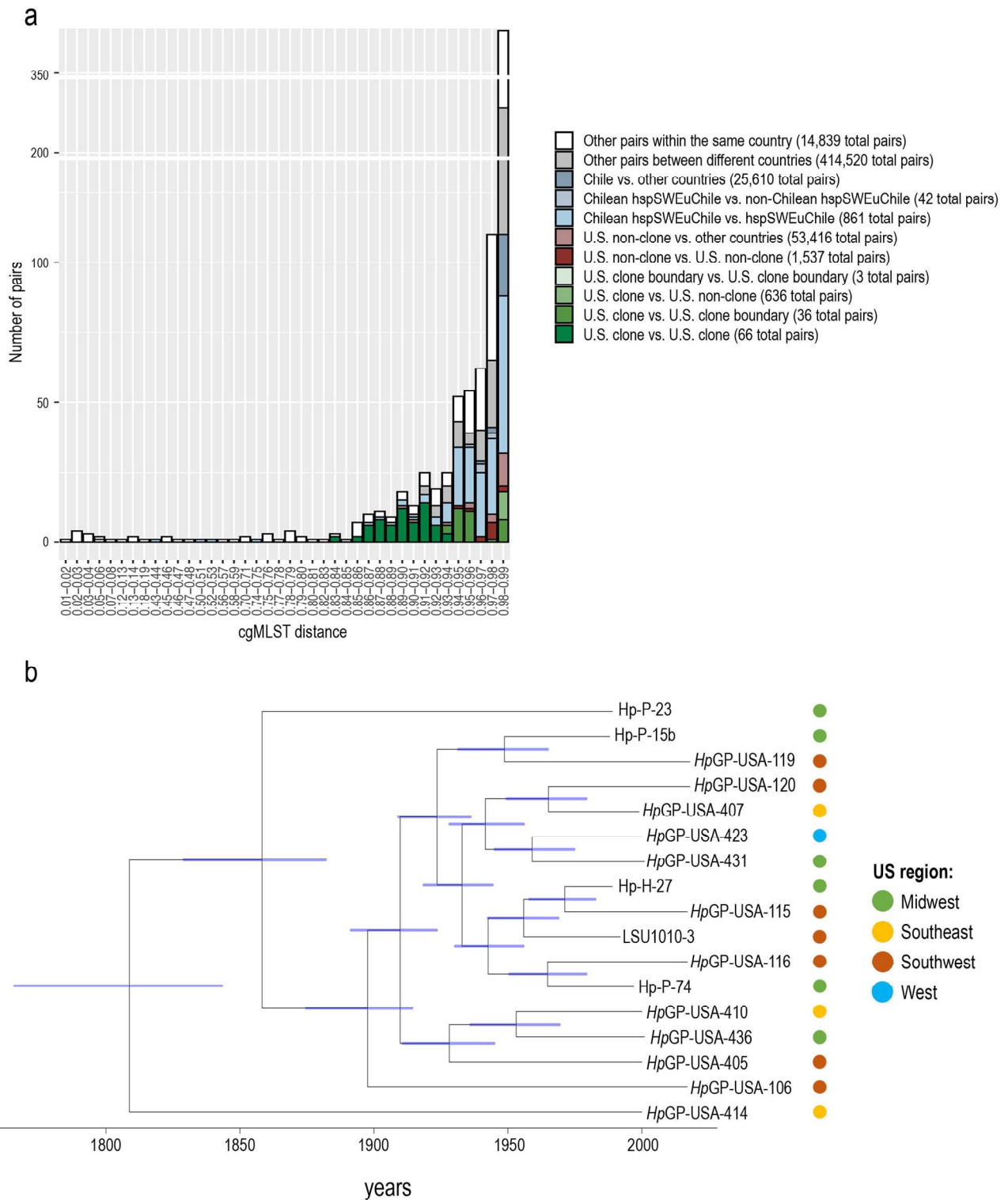


Fig. 4 | In-depth analysis of clonal relationships in the global *H. pylori* dataset. **a** Pairwise core genome MLST (cgMLST) distances of the *HpGP* dataset. Bins illustrate the distribution of core genome allele sharing between pairs of samples. The x-axis ranges from 0.1 to 0.99, with lower values indicating higher number of shared alleles. Every pair is included in a single category of comparison (color bar). Only a small fraction of all possible pairs shares more than 1% of alleles, most of them involving samples from the same country of origin. It is noteworthy that a group of strains from different regions of the US shares between 6% and 17% of alleles

corresponding to 62 and 176 identical genes, suggesting the presence of a deep clone. Other pairs exhibit larger portions of shared alleles (distances <50%), representing recent transmissions between closely related strains. **b** Dated ClonalFrameML tree of the final set of strains considered to belong to the US deep clone *Hp_Clone_US-1*, including five publicly available genomes. Node ages correspond to years based on a previously estimated 1.38×10^{-5} mutation rate per site per year. The colored dots represent the geographical origin of each strain.

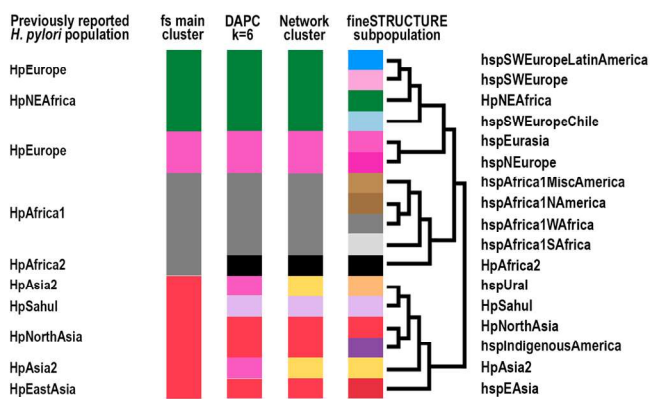


Fig. 5 | Summary of population classifications. Summary of the clustering results using the respective analyses in relation to previously reported MLST and whole genome-based *H. pylori* populations (Hp) and subpopulations (hsp). Colors are based on classifications from the fineSTRUCTURE (fs) analyses visualized in Supplementary Fig. 1, on the $K = 6$ discriminant analysis of principal components, DAPC (Supplementary Fig. 3), and the network clusters (Fig. 2). The topology of the dendrogram to the left is based on the fineSTRUCTURE hierarchical clustering of Supplementary Fig. 1.

collected (95% confidence interval, 107–227 years), while the majority of internal nodes are estimated to be less than 50 years old (Fig. 4b). Thus, the sampled strains are not epidemiologically associated with each other, and instead represent independent strains from a circulating population of clonally related bacteria, which we suggest calling Hp_Clone_US-1.

Latin American subpopulations are more admixed than others

A total of 238 strains from different regions of Latin America were included in the HpGP (Table 1). In the fineSTRUCTURE analysis, most Latin American strains clustered into two previously described populations, hspAfrica1MiscAmerica and hspSWEuropeLatinAmerica^{8,11}, and in hspSWEuropeChile (Supplementary Figs. 1 and 2). Around one-third of the Latin American genomes clustered in non-Latin American populations, the majority in hspAfrica1SAfrica, and hspSWEurope. However, there were also hspEAsia genomes in Argentina, Brazil, and Chile and two hspEurasia genomes from Brazil. Generally, the Latin American genomes were more admixed than their European and African counterparts, with a higher African proportion in hspSWEurope Latin American genomes and a higher European proportion in genomes grouping with hspAfrica1 (https://hpgp.shinyapps.io/Interactive_figures, Fig. 3)

Notably, most Chilean isolates clustered in a separate group, hspSWEuropeChile, (Supplementary Fig. 1), similar to Colombian isolates (hspSWEuropeColombia) previously described^{8,11}. This population is close to hspSWEuropeLatinAmerica and hspSWEurope, as can be seen in the fineSTRUCTURE PCA, particularly in components PC1 and PC7 (https://hpgp.shinyapps.io/Interactive_figures, Fig. 5). However, in the DAPC and network analyses, these strains are dispersed but still near hspSWEurope (Fig. 2 and Supplementary Fig. 3), which is supported by very high self-painting proportions in the chromosome painting analyses (Supplementary Fig. 2), and high pairwise similarities between the genomes of this subpopulation in the cgMLST analysis (Fig. 4a).

Indigenous American *H. pylori* have different ancestral contributions

The fineSTRUCTURE analysis confirmed the hspIndigenousAmerica group^{8,14}. This population is made up of isolates from urban areas of mixed human ancestry, as well as Indigenous communities. HspIndigenousAmerica can be subdivided into two groups called

hspIndigenousNAmerica and hspIndigenousSAmerica (Supplementary Figs. 1 and 2). While hspIndigenousNAmerica is composed of strains from Indigenous communities in North America (Canada and US), the hspIndigenousSAmerica group mostly contains isolates from Latin American regions. In this dataset, we added observations of this subpopulation in Chile, Mexico, Peru, Spain, and the US.

According to the ancestral chromosome painting, and corroborated by the network results, hspIndigenousSAmerica shows a higher proximity to hspEAsia, while hspIndigenousNAmerica has a higher Indigenous-ancestral proportion and is closer to hpNorthAsia in the network analysis, even relatively distanced from hspIndigenousSAmerica (Fig. 2, https://hpgp.shinyapps.io/Interactive_figures, Fig. 3).

Discussion

The intimate association between humans and *H. pylori* started at the beginning of our species and represents a unique story of co-evolution between kingdoms that has fascinated researchers and the public and contributed to understanding human migration dynamics^{14,17}. However, the challenge is to understand the consequences of this thousands-of-years of co-evolution for human health, and on the whole-genome level, bacterial population structure has mostly been studied in the setting of specific geographical areas^{8,10,13–15,18–20}. Ongoing analyses by the HpGP Research Network are comparing between strains from patients with different gastric diseases in order to identify genetic and epigenetic bacterial features that determine human pathogenicity. The HpGP provides a publicly available worldwide collection of complete genomes and epigenomes with high-quality metadata for future investigations of *H. pylori* pathobiology.

Here we present a phylogeographic characterization of the HpGP genomes and outline the global population structure of this bacterium. We used three complementing comparative genomics approaches, fineSTRUCTURE/Chromosome Painting analysis, DAPC, and network analysis of pairwise distances, including interactive visualization of the data, which allowed us to study different aspects of the genomic relationships. A summary of the classifications using the different methods, including their relation to previously reported populations, is presented in Fig. 5, with details in Supplementary Data 3. The higher dynamic range of the DAPC and network analysis clearly showed that hpAfrica2 and hpSahul were very distant from all other populations (Fig. 2, Supplementary Fig. 3b and https://hpgp.shinyapps.io/Interactive_figures, Fig. 1), and the DAPC presented another four main clusters of similarity: a South/West African cluster and a NEAfrica/SWEurope cluster, of which both also had a high presence in the Americas, a North-Central Eurasian cluster, and a North/East Asian cluster, which also included hspIndigenousAmerica. All analyses, however, additionally provided evidence for strong interactions between the hspEurasia and hspSWEurope genomes, and in the 3D plots of ancestry contribution, these populations form a continuum of different ancestry levels, rather than being discrete populations (https://hpgp.shinyapps.io/Interactive_figures, Figs. 2 and 4). Iterating the DAPC analysis to test the consistency of classifications showed, for example, that northeast European genomes from Latvia, Lithuania, Poland, and Russia interchangeably were classified to the clusters corresponding to hspNEurope and hspEurasia. Similarly, some Spanish and Latin American genomes jumped between clusters corresponding to different subpopulations of hspSWEurope (Supplementary Figs. 3d and 4d). However, it was infrequent that genomes were reclassified across the main populations, which supported the relative stability of categories. A few genomes, especially from Indonesia, showed chimeric chromosome painting patterns, for example, a hspUral/hspEAsia combination and a hspUral/hpNEAfrica combination, which constitute rare and exciting intersects between distant populations.

The finding of a highly homogenous group of geographically dispersed genomes in the US motivated us to search for evidence of

distant clonal relationships amongst all *HpGP* strains. The exceptional recombination rate of *H. pylori* means that strains with a common ancestor a few hundred years ago will have recombined most of their genomes, eliminating evidence of a shared clonal frame. Furthermore, an estimated 3.5 billion humans are infected with *H. pylori*²¹ meaning that the current bacterial population size is enormous. As a result, it has been rare to find evidence of clonal relationships between strains collected from distant geographic locations. However, the availability of complete genomes makes it possible to detect deep clones that have recombined in a large fraction of their genome but still share some signal of clonal descent, and the probability of sampling clonally related strains increases quadratically with sample size, meaning that clones will become increasingly common as database sizes increase.

The frequency of the deep clone Hp_Clone_US-1 in the US population is likely somewhere between 3% (proportion in non-*HpGP* US samples) and 18% (proportion in *HpGP*), while it has not yet been found outside the US. The US population in the year 1830 was less than 13 million individuals and has increased to over 330 million through natural population growth and immigration. Assuming the lineage was introduced into the US by a single individual around 1830 and infected 10-fold or more people in each human generation, it would be present in around 3 million individuals today, or about 4% of sampled individuals. These calculations ignore factors such as mixed infection and are subject to many uncertainties but demonstrate that a high level of non-vertical transmission and a significant fitness advantage over other *H. pylori* is necessary to explain the current frequency of Hp_Clone_US-1 in US individuals.

The relative frequency of different transmission routes in the spread of *H. pylori* remains unclear, and while there is evidence of frequent vertical transmission in some populations, other evidence suggests the infection spreads more readily among children^{22,23}. Recent work has emphasized the role of transmission within communities, especially in locations without modern sanitary infrastructure. Our results imply that Hp_Clone_US-1 has been expanding continuously, with several pairs of strains isolated from patients in different states having estimated common ancestors within the last 70 years, which suggests the possibility of occasional mass transmission events in the 20th-century USA. Identification of further clones worldwide should provide additional information to understand when and how some lineages of *H. pylori* can spread fast through human populations. Interestingly, all members of the clone lack the *cag* pathogenicity island, suggesting that also Cag negative strains can be highly competitive under modern conditions.

We note that several geographical regions and human populations remain understudied. Acquiring a better coverage of *H. pylori* whole genomes from South and Central Asia, and a broader representation from the Russian Federation is pivotal. These additional samples would not only offer deeper insights into the hspUral subclades but might also illuminate the possibility of uncovering novel subpopulations stemming from the main ancestral group, HpAsia2. Also, the African continent is still poorly studied in terms of *H. pylori* genomics, which severely limits our understanding of not only population structure but important aspects of bacterial virulence and pathophysiology.

This *HpGP* manuscript was designed as a landmark paper, detailing *Helicobacter pylori* population structure in a global, high-quality dataset. Our intention is for the manuscript to serve as a launching point for individual researchers to deepen the exploration of the detailed data generated by our network. We hope the material (i.e., data and strains) generated by the *HpGP*, including shared resources, codes, and interactive visualizations, together with our main results, will be widely used and will facilitate secondary analyses with the ultimate goal of reducing the burden of the pathologies associated with this bacterial carcinogen.

Methods

Sample acquisition

The *HpGP* samples represent a convenient set. Contributors of samples were identified through advertisements at international scientific meetings, direct invitations to known colleagues and investigators with published sets of *H. pylori* strains, as well as referrals. A limited number of *H. pylori* genomes was publicly available from Spain, one of the main countries responsible for colonial activities in the Americas. Thus, in collaboration of members of the Spanish Association of Gastroenterology, we oversampled this country to better understand the admixed genomes from individuals from Latin America and the Caribbean.

We obtained gastric tissues (fresh frozen with and without culture media; $n = 351$) and cultures (pooled or single colonies; $n = 660$) of *H. pylori* from patients with non-atrophic gastritis ($n = 606$), advanced intestinal metaplasia ($n = 172$, with extension to gastric corpus or incomplete type restricted to antrum), and gastric cancer ($n = 233$). Samples were collected between 1995 and 2020. Biospecimens were shipped to the Division of Gastroenterology, Hepatology, and Nutrition at Vanderbilt University for processing. Before shipment, clinical information and sample descriptions were submitted to the coordinating center at the US National Cancer Institute to confirm eligibility. Biospecimens from the 72 collaborating centers were shipped frozen on dry ice. All individuals provided informed consent, and local Institutional Review Boards approved sample collection. The *HpGP* was exempted from institutional review board evaluation by the National Institutes of Health Office of Human Subjects Research Protection. The summary statistics of 1011 included strains are presented in Table 1, and corresponding NCBI accession numbers and genome statistics are presented in Supplementary Data 1.

Isolation and expansion of *H. pylori* strains and DNA extraction

Gastric tissues (biopsies or fragments from resections) were homogenized under sterile conditions in 100 μ L of sterile phosphate-buffered saline (PBS, pH 7.4) using a homogenizer (Kimble-Kontes, Vineland, NJ, US). Then, 300 μ L of sterile PBS was added to each sample, mixed, and plated onto two selective Trypticase soy agar (TSA) plates with 5% sheep blood containing vancomycin (20 mg/L), bacitracin (200 mg/L), nalidixic acid (10 mg/L) and amphotericin B (2 mg/L) (Sigma, St Louis, MO, US). In addition, a 1:10 dilution was plated on a no-antibiotic TSA plate (BBL; LABSCO, Nashville, TN, US). Agar plates were incubated under microaerobic conditions (Campy Pak Plus envelope, BBL) at 37 °C for 4–6 days until small gray translucent colonies appeared. Gram stains and assays for oxidase and urease were performed. Colony morphology was consistent with the characteristic shape of *H. pylori* colonies. A pool and one single colony of *H. pylori* were expanded and frozen into 1 mL of freezing media (Brucella broth plus 15% glycerol). The single colony was also expanded and used for DNA extraction using Qiagen, QIAamp DNA Mini kit (Qiagen, Catalog number 51306), following the protocol and using the EB buffer to elute the DNA. Original cultures (pooled or single colonies) were processed using the same protocol.

PacBio whole-genome library preparation and sequencing

DNA samples were sequenced at the Cancer Genomics Research Laboratory at the US National Cancer Institute. The manufacturer's protocol was performed for constructing whole-genome libraries from microbial DNA using the SMRTbell Template Prep Kit. Briefly, 1000 ng of genomic DNA, as determined by Quant-iTTM PicoGreen dsDNA Reagent (Thermo Fisher Scientific, Waltham, MA, US), was sheared using the g-TUBE (Covaris, Inc., Woburn, MA, US) to an average fragment size of 10 kb. Following fragmentation and purification, DNA damage, and end repair, hairpin adapters were ligated to the fragment ends to generate SMRTbell libraries. For sequencing on the PacBio RSII instrument, standard hairpin adapters were used. For sequencing on

the PacBio Sequel and Sequel II instruments, barcoded hairpin adapters were used. For the Sequel and Sequel II, barcoded SMRTbell libraries were pooled (up to 8 for the Sequel and up to 48 for the Sequel II), and stringent purification was performed using AMPure PB beads to remove small fragments. Following purification, sequencing primer annealing and DNA polymerase binding of the pooled SMRTbell libraries was performed according to the manufacturer's protocols. SMRT sequencing of the libraries proceeded on the PacBio instrument using 1 SMRT Cell per isolate (RSII) or pool of libraries (Sequel, Sequel II). Genome coverage ranged from 111 to 5678× (median, 949×).

Genome assembly

Since samples were sequenced from multiple generations of PacBio instruments (RSII, Sequel, Sequel II), raw data from the RSII in h5 format were converted to Sequel's subreads XML format so that the same analytical pipeline could be applied to data from all three instruments sequencer. Thus, RSII data were reanalyzed by the same analysis pipeline as Sequel and Sequel II after initial assembly by HGAP3. The original assemblies (from HGAP3) and the new assemblies from the newer SMRTlink assembly tools (HGAP4 or Microbial Assembly) were compared and were highly consistent. Newly reassembled RSII data that generated a single contiguous chromosomal contig were kept.

Whole-genome assembly using raw subreads was performed using SMRTLink's HGAP4/Microbial Assembly, as well as Hifiasm v0.13-r308²⁴ on the HiFi (circular consensus sequencing, CCS) reads. Prior to Hifiasm, raw subreads were converted to circular consensus reads, filtering for CCS reads with minimum predicted read quality higher than 0.99. Chromosomes and plasmids were assembled, which proved to be the most accurate and efficient to achieve complete assembly of all contigs in silico. Circularization with circlator v1.5.3²⁵ was performed on every contig in each strain. Bacterial chromosomal contig start points were all shifted to NusB gene with an additional 12 nt at the 3'-end. MUMmer v3.23²⁶ was used to perform a self-alignment to screen for assembly issues or artifactual contigs, and prokka v1.14.6²⁷ annotation of conserved *H. pylori* genes was run on both raw subread and HiFi read assemblies. The assemblies generated from raw subreads were generally used for methylation calling and downstream analysis unless they failed to be circularized, or prokka annotation suggested a pseudogene percentage higher than 5%, or they contained an unexpected number of tRNA/rRNAs. In those samples, the Hifiasm assembly was used. Candidate chromosomal contigs were also aligned to a published 26695 *H. pylori* strain (NC_000915.1) as a sanity check to verify the contig shared high homology with known *H. pylori*. The detailed analyses of *HpGP* plasmid sequences and their geographic and chromosomal contexts will be reported in full elsewhere.

Assembly quality control

To address the assembly quality, we applied the 3Cs protocol suggested by PacBio (<https://www.pacb.com/blog/beyond-contiguity/>). First, we assessed sequence contiguity and determined that the *HpGP* de novo assemblies all have a contig N50 over 1 Mb. As expected, single chromosomal contigs range from -1.5 to 1.7 Mb. Second, we measured the completeness of our assemblies using BUSCO (Benchmarking Universal Single-Copy Orthologs) scores²⁸ v5.1.3. BUSCO checks the presence or absence of highly conserved genes, and a score >95% is considered a good assembly. For the assemblies that did not achieve a BUSCO score as high as 95%, we either discarded that sequence, or a second attempt was carried out either in silico or in the laboratory. All 1011 *HpGP* assemblies have BUSCO scores above 95%. To further measure correctness, we checked the ratio of pseudogenes, including frameshifted, incomplete, internal stop, ambiguous residues, and multiple problems against the total number of genes. Any assembly with a ratio of pseudogene of more than 5% of the total was discarded. Although the *HpGP* set includes assemblies from three different PacBio

sequencing instruments (15 RSII, 832 Sequel, and 164 Sequel II), the measures of assembly quality (contiguity, BUSCO scores, genomic sequence length, number of total genes, and number of pseudogenes) were similar for the 1011 assemblies from these instruments.

Finally, a consolidated QC report was generated to summarize contig lengths, BUSCO score, and coverage depth (Supplementary Data 1). The minimal chromosomal contig average confidence QV score among most strains was as high as 90.

National Center for Biotechnology Information (NCBI) annotation

The *HpGP* chromosomal sequences were submitted to NCBI, including annotation with the NCBI Prokaryotic Genome Annotation Pipeline, PGAP (https://www.ncbi.nlm.nih.gov/genome/annotation_prok/)^{29–31}. For sequences that could not be circularized ($n=7$), 100 Ns were added to mark breakpoint locations in the genomic sequences. The individual accession numbers and genome statistics are presented in Supplementary Data 1.

Representative genome dataset

To relate the *HpGP* dataset to previous knowledge about *H. pylori* population structure, we used a reference dataset representing the 17 global *H. pylori* subpopulations ($n=255$ genomes, see Supplementary Data 2) described prior to January 2022. To acquire a balanced dataset, we selected 15 genomes per subpopulation, and for the subpopulations with more reported genomes, we selected representatives based on (1) consistency of population assignments in previous publications, (2) assembly quality (contig number, genome size 1.7 ± 0.2 Mbp, and (3) as wide geographical representation within the subpopulation as possible to try to encompass the full breadth of each subpopulation. For the African continent, hpAfrica2, hspAfrica1SAfrica, hspAfrica1WAfrica, and hpNEAfrica were represented, and from Europe hspNEurope, hspSEurope, and hspSWEurope. From Asia, we complemented hpAsia2 and hspEAsia with newly published genomes from hpNorthAsia and the proposed subpopulations hspSiberia, and hspUral¹⁴. For the Americas hspSWEuropeLatinAmerica, hspAfrica1MiscAmericas, hspAfrica1NorthAmerica, hspIndigenousAmericaN, and hspIndigenousAmericaS were represented⁸, and lastly, for Oceania we included hpSahul. The genomes were annotated using prokka v1.14.6²⁷ as previously described^{8,11}.

Core genome analysis

All population structure analyses (fineSTRUCTURE, Chromosome Painting, Network analysis, and DAPC), were based on the same core gene alignment. This was generated using the prokka-annotated 1011 *HpGP* genomes plus a resequenced ATCC reference strain 26695 (*HpGP*-26695) and the 255 representative genomes, a total of 1267 genomes. The analysis was performed using the panaroo pipeline v1.2.10³² using 90% protein sequence identity and 75% gene length coverage cut-off.

Population structure analysis

The genome-wide haplotype data was calculated as described previously³³; we conducted SNP calling for each alignment, and imputation for polymorphic sites with missing frequency <1% using BEAGLE v3.3.2³⁴. This genome-wide haplotype contained 387,927 SNPs in 1227 genes and was used to define isolate populations and subpopulations based on the similarity of the haplotype copying profiles obtained by fineSTRUCTURE v4. Then, fineSTRUCTURE³⁵ analysis was performed with 200,000 iterations of both the burn-in and Markov chain Monte Carlo (MCMC) method to cluster individuals based on the coancestry matrix as described³⁶. The results were visualized as a heat map with each cell indicating the proportion of DNA "chunks" a recipient receives from each donor. Furthermore, the posterior distribution of the clusters was visualized using

fineSTRUCTURE's tree-building algorithm to define the populations and subpopulations produced. With the previously obtained coancestry matrix, multiple principal component analysis (PCA) was calculated to analyze the population structure in detail. Principal components (PCs) 1 to 11 were calculated and visualized using R.

DAPC analysis

We employed discriminant analysis of principal components (DAPC) to further investigate the genetic structure of our data. DAPC describes clusters in genetic data by creating synthetic variables (discriminant functions) that maximize variance among groups while minimizing variance within groups⁹. DAPC is a multivariate approach, not model based; hence, it makes no assumptions about Hardy-Weinberg or linkage equilibrium on genetic loci. Before running the DAPC, we assessed the number of clusters most supported for our *H. pylori* dataset by employing the *find.clusters* function in adegenet R package, comparing the results of 100 independent runs using a custom-made R script and selecting the optimal number of clusters according to Bayesian information criteria (BIC). In order to assess the uncertainty of the group assignments of each individual, we visualized posterior group membership probabilities based on the DAPC analysis using the function *compplot*.

Using SNP-sites v2.5.1³⁷, we extracted 601,000 SNPs from the 1267 genome panaroo core gene alignment. We employed the function *optim.a.score* of the adegenet R package to identify the optimal number of principal components to consider for the analyses, as too many could lead to overfitting, while a low number of components could decrease discriminatory power between groups.

We first ran a DAPC employing the most supported number of clusters/groups as estimated by the *find.cluster* analysis: $K = 6$. Initially, we ran the DAPC considering all the sequences in our alignment so as to visualize the entire genetic variability of our data. Given the outlier position of the two groups, we subsequently ran another DAPC analysis excluding these outliers to better emphasize the differences among the other clusters. Finally, we computed posterior group membership probability for each individual. This parameter is based on the retained discriminant functions of the DAPC analysis and represents the probability of each sample to be assigned to a group, which can be interpreted in order to assess how clear-cut or admixed the clusters are. We also ran the DAPC procedure considering $K = 17$, the same number of clusters identified by the fineSTRUCTURE analysis.

Network analysis of core genes

The core gene alignment obtained with panaroo was used to estimate distances with PAUP³⁸ v4.0a166, using maximum likelihood criteria. Each distance was normalized between 0 and 1 as previously described⁸. With this normalization, 0 means the highest genetic similarity, and 1 signifies the highest dissimilarity between two strains. Next, a complete network is created, where all pairs of strains have a measure of genetic distance based on this previous normalization. Strains are represented as vertices, and their distances are represented as edges. In the beginning, this network is fully connected and has no perceptible structure. A process of edge and node pruning is carried out to reveal the underlying structure of the genetic similarity between strains. This process consists of ranking the values of the edges and removing them subsequently, starting with the most dissimilar (equal or close to one). This process is continued until the network is subdivided into a determined number of Connected Components (CC). We consider a CC as a set of more than two nodes connected between them but isolated from other groups of nodes. If a single node is stripped of all its edges (singleton), we discard this node from the set of nodes of the resulting network. Figure 2 was created following this pruning process with a CC threshold of 2. This means that edges and singletons were removed until the full network was separated into two groups of nodes, with the separated group being the hpAfrica2 group.

Chromosome painting

Full dataset analysis. To identify the patterns of shared genomic content of *H. pylori* isolates, we conducted chromosome painting using ChromoPainterV2³⁵, designating all genomes as recipients (1011 HpGP genomes), and randomly selected ~20 isolates per population as donors (335 genomes; Supplementary Data 3). Each strain was painted using all the other donor samples and the result is visualized in a bar plot built with R.

HpGP only ancestral chromosome painting. Since we have no genomes from the true historically ancestral populations, genomic ancestry is commonly inferred from contemporary representatives of these populations. A second chromosome painting was thus performed using hpAfrica2, hspAfrica1WAfrica, hpNEAfrica, hpAsia2, hpNorthAsia, hspUral, and hspEAsia as donors to infer ancestral contributions to the populations in the HpGP dataset^{3,10}. We also only selected donors among the reference collection for which *H. pylori* population assignment and geographical origin were concordant (Supplementary Data 3).

Core gene multilocus sequence typing (cgMLST)

To investigate the existence of clonal relationships in *H. pylori*, we estimated the total number of identical loci shared among strains from the HpGP dataset by performing a cgMLST as implemented by chewBBACA³⁹ software v2.8.5. chewBBACA uses a gene-by-gene method to compare coding sequences and assign alleles based on a BLAST Score Ratio (BSR)⁴⁰. We first used Prodigal⁴¹ v2.6.3, including the option *-t* to create a training file from the assembled version of the 26695 *H. pylori* reference strain resequenced as part of the HpGP dataset. Then, the "CreateSchema" module of chewBBACA was applied to the 1011 HpGP genomes and the Prodigal training file to estimate a whole-genome MLST (wgMLST) scheme. The 3943 wgMLST genes were then compared with the "AlleleCall" module, using the default BSR threshold of 0.6. A total of 867 genes identified as paralogs were removed from the wgMLST using the "RemoveGenes" module, reducing the scheme to 3076 loci. We then used the "ExtractCgMLST" module to create a cgMLST with all loci present in more than 95 percent of strains (*-t* 0.95), obtaining a total of 981,110 alleles for 1040 loci, an average of 943 different alleles per locus.

We last used the cgMLST allelic profile to calculate pairwise distances with GrapeTree⁴² v1.5.0, running it in "--wgMLST" mode with the "distance" method (*-method distance*) while ignoring missing data (*--missing 0*). We analyzed the distribution of cgMLST distances between pairs of strains in categories such as "US clone", "US clone boundary", "US non-clone", "Chile", "Chilean hspSWEuropeChile", "non-Chilean hspSWEuropeChile", "within the same country", and "between different countries", as depicted in Fig. 4a.

Analysis of public US genomes

We downloaded all whole-genome sequences publicly available in the Enterobase *H. pylori* database (<https://enterobase.warwick.ac.uk/species/index/helicobacter>) with the US as the country of isolation as of September 18, 2022 ($n = 226$). Sixty-seven sequences were either isolated from non-human hosts, results of experimental infections, repeated samplings from the same individual or overlapping the HpGP set, thus were excluded. The remaining 151 genomes (Supplementary Data 4) were combined with the HpGP US genomes and the 255 references in a kmer-based genomic distance analysis using mash v2.3⁴³. The five genomes clustering with the US deep clone were added to the dataset used for in-depth analysis.

Dating of the US deep clone

A core gene alignment of the highly clonal US genomes, including the five public ones, was generated with panaroo using the settings described above. Three genomes, HpGP-USA-401, HpGP-USA-404, and

HpGP-USA-414 had diverged from the clone both by phylogeny and chromosome painting profile and were excluded from further analysis. A phylogenetic tree was computed using PhyML v3.1⁴⁴ and input to ClonalFrameML v1.11-3-g4f13f23⁴⁵, executed using default parameters. Node ages were determined using the R BactDating package⁴⁶, using 10,000 Markov chain Monte Carlo iterations and a mutation rate of 1.38×10^{-5} per site per year, as has previously been estimated¹⁶.

Data visualization

The map figures of the dataset's geographical distribution, including the gray background map, were plotted using the ggplot2⁴⁷ and ggmaps⁴⁸ package in R. The painting profiles were summarized as described above, and plotting and statistical analysis was performed in R using the ggplot2 and plotly⁴⁹ packages.

Strain availability

The *HpGP* set of *H. pylori* strains is available from the US National Cancer Institute for scientific purposes upon a reasonable request. However, restrictions apply to its availability as some samples require authorization from contributing centers to be distributed to third parties.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

The whole-genome sequences generated within the *HpGP* have been deposited in the NCBI GenBank database under BioProject accession code PRJNA529500 [<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA529500>] (Supplementary Data 1). NCBI or equivalent public accessions for the reference set are listed in Supplementary Data 2. The whole *HpGP* genome dataset and the 255 reference genomes are also deposited to Zenodo, DOI: 10.5281/zenodo.10048320. Source Data for the individual figures are available with this paper.

Code availability

The computational scripts to process the data and plot figures are available at <https://github.com/HpGP/Code-and-Data> v1.0. This code is also archived on Zenodo under <https://doi.org/10.5281/zenodo.8381170>.

References

- Fox, J. G. & Wang, T. C. Inflammation, atrophy, and gastric cancer. *J. Clin. Investig.* **117**, 60–69 (2007).
- Conteduca, V. et al. *H. pylori* infection and gastric cancer: state of the art (review). *Int. J. Oncol.* **42**, 5–18 (2013).
- Falush, D. et al. Traces of human migrations in *Helicobacter pylori* populations. *Science* **299**, 1582–1585 (2003).
- Linz, B. et al. An African origin for the intimate association between humans and *Helicobacter pylori*. *Nature* **445**, 915–918 (2007).
- Moodley, Y. et al. Age of the association between *Helicobacter pylori* and man. *PLoS Pathog.* **8**, e1002693 (2012).
- Yamaoka, Y. *Helicobacter pylori* typing as a tool for tracking human migration. *Clin. Microbiol. Infect.* **15**, 829–834 (2009).
- Sung, H. et al. Global Cancer Statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J. Clin.* **71**, 209–249 (2021).
- Munoz-Ramirez, Z. Y. et al. A 500-year tale of co-evolution, adaptation, and virulence: *Helicobacter pylori* in the Americas. *ISME J.* **15**, 78–92 (2021).
- Jombart, T., Devillard, S. & Balloux, F. Discriminant analysis of principal components: a new method for the analysis of genetically structured populations. *BMC Genet.* **11**, 94 (2010).
- Thorpe, H. A. et al. Repeated out-of-Africa expansions of *Helicobacter pylori* driven by replacement of deleterious mutations. *Nat. Commun.* **13**, 6842 (2022).
- Thorell, K. et al. Rapid evolution of distinct *Helicobacter pylori* subpopulations in the Americas. *PLoS Genet.* **13**, e1006546 (2017).
- Berthenet, E. et al. A GWAS on *Helicobacter pylori* strains points to genetic variants associated with gastric cancer risk. *BMC Biol.* **16**, 84 (2018).
- You, Y. et al. Genomic differentiation within East Asian *Helicobacter pylori*. *Microb. Genom.* <https://doi.org/10.1099/mgen.0.000676> (2022).
- Moodley, Y. et al. *Helicobacter pylori*'s historical journey through Siberia and the Americas. *Proc. Natl Acad. Sci. USA.* <https://doi.org/10.1073/pnas.2015523118> (2021).
- Suzuki, R. et al. *Helicobacter pylori* genomes reveal Paleolithic human migration to the east end of Asia. *iScience* **25**, 104477 (2022).
- Didelot, X. et al. Genomic evolution and transmission of *Helicobacter pylori* in two South African families. *Proc. Natl Acad. Sci. USA.* **110**, 13880–13885 (2013).
- Moodley, Y. & Linz, B. *Helicobacter pylori* sequences reflect past human migrations. *Genome Dyn.* **6**, 62–74 (2009).
- Kumar, N., Albert, M. J., Al Abkal, H., Siddique, I. & Ahmed, N. What constitutes an Arabian *Helicobacter pylori*? Lessons from comparative genomics. *Helicobacter.* <https://doi.org/10.1111/hel.12323> (2017).
- Kumar, N. et al. Comparative genomic analysis of *Helicobacter pylori* from Malaysia identifies three distinct lineages suggestive of differential evolution. *Nucleic Acids Res.* **43**, 324–335 (2015).
- Oleastro, M., Rocha, R. & Vale, F. F. Population genetic structure of *Helicobacter pylori* strains from Portuguese-speaking countries. *Helicobacter.* <https://doi.org/10.1111/hel.12382> (2017).
- Li, Y. et al. Global prevalence of *Helicobacter pylori* infection between 1980 and 2022: a systematic review and meta-analysis. *Lancet Gastroenterol. Hepatol.* **8**, 553–564 (2023).
- Ford, A. C. et al. Effect of sibling number in the household and birth order on prevalence of *Helicobacter pylori*: a cross-sectional study. *Int. J. Epidemiol.* **36**, 1327–1333 (2007).
- Goodman, K. J. & Correa, P. Transmission of *Helicobacter pylori* among siblings. *Lancet* **355**, 358–362 (2000).
- Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat. Methods* **18**, 170–175 (2021).
- Hunt, M. et al. Circlator: automated circularization of genome assemblies using long sequencing reads. *Genome Biol.* **16**, 294 (2015).
- Kurtz, S. et al. Versatile and open software for comparing large genomes. *Genome Biol.* **5**, R12 (2004).
- Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).
- Manni, M., Berkeley, M. R., Seppey, M., Simao, F. A. & Zdobnov, E. M. BUSCO Update: novel and streamlined workflows along with broader and deeper phylogenetic coverage for scoring of eukaryotic, prokaryotic, and viral genomes. *Mol. Biol. Evol.* **38**, 4647–4654 (2021).
- Li, W. et al. RefSeq: expanding the Prokaryotic Genome Annotation Pipeline reach with protein family model curation. *Nucleic Acids Res.* **49**, D1020–D1028 (2021).
- Haft, D. H. et al. RefSeq: an update on prokaryotic genome annotation and curation. *Nucleic Acids Res.* **46**, D851–D860 (2018).
- Tatusova, T. et al. NCBI prokaryotic genome annotation pipeline. *Nucleic Acids Res.* **44**, 6614–6624 (2016).
- Tonkin-Hill, G. et al. Producing polished prokaryotic pangenomes with the Panaroo pipeline. *Genome Biol.* **21**, 180 (2020).

33. Yahara, K., Didelot, X., Ansari, M. A., Sheppard, S. K. & Falush, D. Efficient inference of recombination hot regions in bacterial genomes. *Mol. Biol. Evol.* **31**, 1593–1605 (2014).
34. Browning, B. L. & Browning, S. R. A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *Am. J. Hum. Genet.* **84**, 210–223 (2009).
35. Lawson, D. J., Hellenthal, G., Myers, S. & Falush, D. Inference of population structure using dense haplotype data. *PLoS Genet.* **8**, e1002453 (2012).
36. Yahara, K. et al. Chromosome painting in silico in a bacterial species reveals fine population structure. *Mol. Biol. Evol.* **30**, 1454–1464 (2013).
37. Page, A. J. et al. SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments. *Microbial. Genomics.* <https://doi.org/10.1099/mgen.0.000056> (2016).
38. Wilgenbusch, J. C. & Swofford, D. Inferring evolutionary trees with PAUP*. *Curr. Protoc. Bioinformatics* Chapter 6, Unit 6.4. <https://doi.org/10.1002/0471250953.bi0604s00> (2003).
39. Silva, M. et al. chewBBACA: a complete suite for gene-by-gene schema creation and strain identification. *Microb. Genom.* <https://doi.org/10.1099/mgen.0.000166> (2018).
40. Rasko, D. A., Myers, G. S. & Ravel, J. Visualization of comparative genomic analyses by BLAST score ratio. *BMC Bioinforma.* **6**, 2 (2005).
41. Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinforma.* **11**, 119 (2010).
42. Zhou, Z. et al. GrapeTree: visualization of core genomic relationships among 100,000 bacterial pathogens. *Genome Res.* **28**, 1395–1404 (2018).
43. Ondov, B. D. et al. Mash: fast genome and metagenome distance estimation using MinHash. *Genome Biol.* **17**, 132 (2016).
44. Guindon, S., Delsuc, F., Dufayard, J. F. & Gascuel, O. Estimating maximum likelihood phylogenies with PhyML. *Methods Mol. Biol.* **537**, 113–137 (2009).
45. Didelot, X. & Wilson, D. J. ClonalFrameML: efficient inference of recombination in whole bacterial genomes. *PLoS Comput. Biol.* **11**, e1004041 (2015).
46. Didelot, X., Croucher, N. J., Bentley, S. D., Harris, S. R. & Wilson, D. J. Bayesian inference of ancestral dates on bacterial phylogenetic trees. *Nucleic Acids Res.* **46**, e134 (2018).
47. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis* (Springer-Verlag New York). <https://ggplot2.tidyverse.org> (2016).
48. Kahle, D. W. H. ggmap: spatial visualization with ggplot2. *R. J.* **5**, 144–161 (2013).
49. *Collaborative Data Science* (Plotly Technologies Inc., Montréal, QC, 2015).

Acknowledgements

Our special thanks are extended to all the individuals who were the hosts of the strains of the *HpGP* collection, who represent the human populations that bear the burden of *H. pylori*-associated disease, and whose biological samples serve to advance research aimed at reducing this disease burden. We immensely thank Lisa D. Finkelstein, Ramona Bhatnagary, Jillian M. Varonin, Mary Jane Williams, Karen Williams Kinney, and Melissa A. Raymond from the US National Institutes of Health (National Cancer Institute's Technology Transfer Center, National Cancer Institute's Division of Cancer Epidemiology, and Genetics, and Division of Logistic Services) for their administrative and logistic support in establishing the collaboration agreements and importing the multiple sets of biospecimens. We also thank the US Centers for Disease Control and Prevention's Import Permit Program. We dedicate this work to our deceased colleagues Pablo Luna, Radislav Nakov, Bongani Kaimila, and Khean Lee Goh who passed away in recent years.

The *HpGP* was mainly supported by the Intramural Research Program from the US National Cancer Institute (NCI), National Institutes of Health (NIH). This work was supported in part by the intramural research programs of the US National Library of Medicine, the US National Institute on Minority Health and Health Disparities, and the US National Institute of Allergy and Infectious Diseases. The members of the bioinformatics group received support from the Swedish Society for Medical Research (K.T.), Assar Gabrielsson Foundation (K.T.), and Magnus Bergvall Foundation (K.T.). The collaborating centers for sample collection received grant support from the US NIH (P01CA116087, R01CA077955, R01DK058587 and P30DK058404 to R.M.P.; P01CA028842 and R01CA190612 to K.T.W.; P01CA028842, R01CA190612, K07CA125588, R03CA167773 and P30CA068485 to D.R.M.; K08CA252635 to R.J.H., K22CA226395 to M.G.-P.; and U54GM133807 to M.C.-C.), the German Federal Ministry of Education and Research (BMBF-0315905D, ERA-NET PathoGenoMics to P.M.), the French Association pour la Recherche Contre le Cancer (8412 to F.M.), the French Institut National du Cancer (07/3D1616/IABC-23-12/NC-NG and 2014-152 to F.M.), the Canceropole Grand Sud-Ouest (2010-08-canceropole GSO-Universite Bordeaux 2 to F.M.), the Japanese National Institutes of Health (DK62813 to Y.Y.), the Japanese Ministry of Education, Culture, Sports, Science, and Technology (18KK0266, 19H03473, 21H00346 and 22H02871 to Y.Y.), the National Fund for Innovation and Development of Science and Technology from the Ministry of Higher Education Science and Technology of the Dominican Republic (2012-2013-2A1-65 and 2015-3A1-182 to M.C.), the National Cancer Center of South Korea (2210630, I.J.C.), ArcticNet (RES0010178 to K.J.G.), the Network of Centres of Excellence of Canada, the Canadian Institutes for Health Research (MOP115031 to K.J.G.), Alberta Innovates Health Solutions (201201159 to K.J.G.), the University of Malaya-Ministry of Higher Education (UM.C/625/1/HIR/MOHE/CHAN-02 to J.V.), the Ministry of Science and Technology of Vietnam, the Kyrgyz State Medical Academy, the Italian Ministry of Health for Institutional Research, the Chilean National Fund for Health Research and Development (FONIS A19/0188, FONDECYT 1230504 and ANID-FONDAP 152220002 to A.R.; CONICYT-FONDAP 15130011 and FONDECYT 1231773 to A.H.C.), the Chilean Cancer Prevention and Control Center, the Horizon 2020 Programme of European Union (825832, "CeLac and European consortium for a personalized medicine approach to Gastric Cancer," LEGACy, to T.F.-K. and A.R.), the Fundação de Amparo à Pesquisa do Estado de São Paulo (FAPESP; 2014/26847-0, 2018/14267-2, 2018/02972-3 to E.D.-N.), the Departamento de Ciência e Tecnologia (DECIT), Ministry of Health, Brazil (PRONON, SIPAR 2500.035-167/2015-23 to E.D.-N.), the Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq, 314344/2020-9 to E.T.-S.), the Universidad de Costa Rica (742-B9-310 and 742-90912-19 to V.R.-M.), LABGIPAT (S.D.-B.), the Hospital Clínica Bíblica (C.C.-N.), the Greek Ministry of Culture and Education (InfeNeutra Project, NSRF 2007-2013, MIS450598, D.N.S.), the National Strategic Reference Framework Operational Program "Competitiveness, Entrepreneurship and Innovation" (NSRF 2014-2020, MIS5002486, D.N.S.), the Hellenic *Helicobacter pylori* Study Group (2012-2016, B.M.-G.), the Hellenic Society of Gastroenterology (National Multicenter Laboratory Surveillance Studies, 2018-2019, B.M.-G.), the Ministry of Science and Technology, Executive Yuan, Taiwan (109-2314-B-002-096; MOST 111-2314-B-002-012; MOST 109-2314-B-002-090-MY3 to J.-M.L. and M.-S.W.), the National Research Foundation of Singapore, the Singapore Ministry of Health's National Medical Research Council (Open Fund-Large Collaborative Grant, MOH-OFLCG18May-0003), the University of Puerto Rico Comprehensive Cancer Center, the Fondo Nacional de Desarrollo Científico y Tecnológico (196-2015-FONDECYT to C.C.), Universidad Científica del Sur, and Instituto Nacional de Enfermedades Neoplásicas (INEN, Peru).

The computations and data storage required for the analyses presented were enabled by resources in projects snic-2021/22-229 and snic-2021/23-234 provided by the National Academic Infrastructure for Supercomputing in Sweden (NAISS) and the Swedish National

Infrastructure for Computing (SNIC) at the UPPMAX HPC, partially funded by the Swedish Research Council through grant agreements 2022-06725, and 2018-05973.

Author contributions

Manuscript conception and design: K.T. and Z.Y.M.R. Analysis supervision: K.T., S.G., and D.F. Data analysis: K.T., Z.Y.M.R., D.W., S.S.M., R.B.A., S.G., and R.C.T. Interpretation of results: K.T., Z.Y.M.R., S.S.M., R.B.A., S.G., R.C.T., D.F., M.C.C., and C.S.R. Manuscript writing: K.T., Z.Y.M.R., D.W., R.C.T., D.F., and M.C.C. Data coordinator: D.W. Editing of the manuscript: S.S.M., R.B.A., S.G., HpGP Research Network and C.S.R. Sample acquisition: HpGP Research Network. Conception and design of the HpGP initiative: M.C.C. and C.S.R. HpGP study coordinators: M.C.C. and C.S.R.

Funding

Open access funding provided by University of Gothenburg.

Competing interests

J.P.G. has served as a speaker, consultant, and advisory member for or has received research funding from Mayoly, Allergan, Diasorin, Gebro Pharma, and Richen. E.B.-M. has served as a speaker and consultant for Janssen, Chiesi, Kern and Takeda. R.M.F., J.C.M., and C.F. own patent WO/2018/169423 on microbiome markers for gastric cancer, and R.J.R. works for New England Biolabs, a company that sells research reagents, including restriction enzymes and DNA methyltransferases, to the scientific community. The remaining authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-023-43562-y>.

Correspondence and requests for materials should be addressed to Kaisa Thorell.

Peer review information *Nature Communications* thanks the anonymous reviewers for their contribution to the peer review of this work. A peer review file is available.

Reprints and permissions information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023

¹Department of Chemistry and Molecular Biology, University of Gothenburg, Gothenburg, Sweden. ²Facultad de Ciencias Químicas, Universidad Autónoma de Chihuahua, Chihuahua, Chihuahua, México. ³Cancer Genomics Research Laboratory, Frederick National Laboratory for Cancer Research, Frederick, MD, USA. ⁴Division of Cancer Epidemiology and Genetics, National Cancer Institute, Rockville, MD, USA. ⁵Instituto Nacional de Medicina Genómica, Ciudad de México, México. ⁶Consejo Nacional de Ciencia y Tecnología, Cátedras CONACYT, Ciudad de México, México. ⁷Centro de Ciencias de la Complejidad, Universidad Nacional Autónoma de México, Ciudad de México, México. ⁸Department of Life Sciences and Biotechnology, University of Ferrara, Ferrara, Italy. ⁹Centre for Microbes Development and Health, Institute Pasteur Shanghai, Shanghai, China. ²⁰⁴These authors contributed equally: Kaisa Thorell, Zilia Y. Muñoz-Ramírez. ²⁰⁵These authors jointly supervised this work: M. Constanza Camargo, Charles S. Rabkin. ✉ e-mail: kaisa.thorell@gu.se

HpGP Research Network

Kaisa Thorell^{1,204} ✉, Zilia Y. Muñoz-Ramírez^{2,204}, Difei Wang^{3,4}, Santiago Sandoval-Motta^{5,6,7}, Rajiv Boscolo Agostini⁸, Silvia Ghirotto⁸, Roberto C. Torres⁹, Judith Romero-Gallo¹⁰, Uma Krishna¹⁰, Richard M. Peek Jr¹⁰, M. Blanca Piazuelo¹⁰, Naïma Raaf¹¹, Federico Bentolila¹², Hafeza Aftab¹³, Junko Akada¹⁴, Takashi Matsumoto¹⁴, Freddy Haesebrouck¹⁵, Rony P. Colanzi¹⁶, Thais F. Bartelli¹⁷, Diana Noronha Nunes¹⁷, Adriane Pelosof¹⁷, Claudia Zitron Sztokfisz¹⁷, Emmanuel Dias-Neto¹⁷, Paulo Pimentel Assumpção¹⁸, Ivan Tishkov¹⁹, Laure Brigitte Kouitcheu Mabeku²⁰, Karen J. Goodman²¹, Janis Geary²¹, Taylor J. Cromarty²¹, Nancy L. Price²¹, Douglas Quilty²², Alejandro H. Corvalan²³, Carolina A. Serrano²⁴, Robinson Gonzalez²⁵, Arnoldo Riquelme²⁵, Apolinaria García-Cancino²⁶, Cristian Parra-Sepúlveda²⁶, Giuliano Bernal²⁷, Francisco Castillo²⁸, Alisa M. Goldstein⁴, Nan Hu⁴, Philip R. Taylor⁴, Maria Mercedes Bravo²⁹, Alvaro Pazos³⁰, Luis E. Bravo³¹, Keith T. Wilson¹⁰, James G. Fox³², Vanessa Ramírez-Mayorga³³, Silvia Molina-Castro³³, Sundry Durán-Bermúdez³⁴, Christian Campos-Núñez³⁵, Manuel Chaves-Cervantes³⁵, Evariste Tshibangu-Kabamba^{36,37}, Ghislain Disashi Tumba³⁷, Antoine Tshimpi-Wola³⁸, Patrick de Jesus Ngoma-Kisoko³⁸, Dieudonné Mumba Ngoyi^{38,39}, Modesto Cruz⁴⁰, Celso Hosking⁴⁰, José Jiménez Abreu⁴¹, Christine Varon⁴², Lucie Benejat⁴², Ousman Secka⁴³, Alexander Link⁴⁴, Peter Malfertheiner⁴⁴, Michael Buenor Adinortey⁴⁵, Ansumana Sandy Bockarie⁴⁶, Cynthia Ayefoumi Adinortey⁴⁷, Eric Gyamerah Ofori⁴⁸, Dionyssios N. Sgouras⁴⁹, Beatriz Martinez-Gonzalez⁴⁹, Spyridon Michopoulos⁵⁰, Sotirios Georgopoulos⁵¹, Elisa Hernandez⁵², Braulio Volga Tacatic⁵³, Mynor Aguilar⁵⁴, Ricardo L. Dominguez⁵⁵, Douglas R. Morgan⁵⁶, Hjördís Harðardóttir⁵⁷, Anna Ingibjörg Gunnarsdóttir⁵⁷, Hallgrímur Guðjónsson⁵⁷, Jón Gunnlaugur Jónasson^{57,58}, Einar S. Björnsson^{57,58}, Mamatha Ballal⁵⁹, Vignesh Shetty^{59,60}

Muhammad Miftahussurur⁶¹, Titong Sugihartono⁶¹, Ricky Indra Alfaray⁶¹, Langgeng Agung Waskito⁶¹, Kartika Afrida Fauzia⁶¹, Ari Fahrial Syam⁶², Hasan Maulahela⁶², Reza Malekzadeh⁶³, Masoud Sotoudeh⁶³, Avi Peretz^{64,65}, Maya Azrad^{64,65}, Avi On^{65,66}, Valli De Re⁶⁷, Stefania Zanussi⁶⁷, Renato Cannizzaro⁶⁸, Vincenzo Canzonieri⁶⁹, Takaya Shimura⁷⁰, Kengo Tokunaga⁷¹, Takako Osaki⁷¹, Shigeru Kamiya⁷¹, Khaled Jadallah⁷², Ismail Matalka⁷², Nurbek Igissinov⁷³, Mariia Satarovna Moldobaeva⁷⁴, Attokurova Rakhat⁷⁴, Il Ju Choi⁷⁵, Jae Gyu Kim⁷⁶, Nayoung Kim⁷⁷, Minkyong Song⁷⁸, Mārcis Leja^{78,79,80}, Reinis Vangravs⁷⁸, Ģirts Šķenders^{78,80}, Dace Rudzīte⁸⁰, Aiga Rūdile⁷⁸, Aigars Vanags⁷⁹, Ilze Kikuste⁷⁸, Juozas Kupcinskas⁸¹, Jurgita Skieceviciene⁸¹, Laimas Jonaitis⁸¹, Gediminas Kiudelis⁸¹, Paulius Jonaitis⁸¹, Vytautas Kiudelis⁸¹, Greta Varkalaite⁸¹, Jamuna Vadivelu^{82,83}, Mun Fai Loke⁸², Kumutha Malar Vellasamy⁸², Roberto Herrera-Goepfert⁸⁴, Juan Octavio Alonso-Larraga⁸⁵, Than Than Yee⁸⁶, Kyaw Htet⁸⁶, Takeshi Matsuhisa⁸⁷, Pradeep Krishna Shrestha⁸⁸, Shamsul Ansari⁸⁹, Olumide Abiodun⁹⁰, Christopher Jemilohun⁹¹, Kolawole Oluseyi Akande⁹², Oluwatosin Olu-Abiodun⁹³, Francis Ajang Magaji⁹⁴, Ayodele Omotoso⁹⁵, Chukwuemeka Chukwunwendu Osuagwu⁹⁶, Uchenna Okonkwo⁹⁵, Opeyemi O. Owoseni⁹⁷, Carlos Castaneda⁹⁸, Miluska Castillo⁹⁹, Billie Velapatino¹⁰⁰, Robert H. Gilman¹⁰¹, Paweł Krzyżek¹⁰², Grażyna Gościńskiak¹⁰², Dorota Pawełka¹⁰³, Izabela Korona-Głowniak¹⁰⁴, Halina Cichoz-Lach¹⁰⁵, Monica Oleastro¹⁰⁶, Ceu Figueiredo^{107,108,109}, Jose C. Machado^{107,108,109}, Rui M. Ferreira^{107,108}, Dmitry S. Bordin^{110,111,112}, Maria A. Livzan¹¹³, Vladislav V. Tsukanov¹¹⁴, Patrick Tan^{115,116,117}, Khay Guan Yeoh^{118,119}, Feng Zhu¹¹⁸, Reid Ally^{120,121}, Rainer Haas¹²², Milagrosa Montes¹²³, María Fernández-Reyes¹²³, Esther Tamayo¹²³, Jacobo Lizasoain¹²³, Luis Bujanda¹²⁴, Sergio Lario^{125,126}, María José Ramírez-Lázaro^{125,126}, Xavier Calvet^{125,126}, Eduard Brunet-Mas^{125,126}, María José Domper-Arnal^{127,128}, Sandra García-Mateo^{127,128}, Daniel Abad-Baroja¹²⁹, Pedro Delgado-Guillena¹³⁰, Leticia Moreira¹³¹, Josep Botargues¹³², Isabel Pérez-Martínez^{133,134}, Eva Barreiro-Alonso^{133,135,136}, Virginia Flores¹³⁷, Javier P. Gisbert^{138,139}, Edurne Amorena Muro¹⁴⁰, Pedro Linares¹⁴¹, Vicente Martín¹⁴², Laura Alcoba¹⁴¹, Tania Fleitas-Kanonnikoff¹⁴³, Hisham N. Altayeb^{144,145}, Lars Engstrand¹⁴⁶, Helena Enroth¹⁴⁷, Peter M. Keller^{148,149}, Karoline Wagner¹⁵⁰, Daniel Pohl¹⁵¹, Yi-Chia Lee¹⁵², Jyh-Ming Liou¹⁵², Ming-Shiang Wu¹⁵², Bekir Kocazeybek¹⁵³, Suat Saribas¹⁵³, İhsan Taşçı¹⁵⁴, Süleyman Demiryas¹⁵⁴, Nuray Kepil¹⁵⁵, Luis Quiel¹⁵⁶, Miguel Villagra¹⁵⁷, Morgan Norton¹⁵⁸, Deborah Johnson¹⁵⁸, Robert J. Huang¹⁵⁹, Joo Ha Hwang¹⁵⁹, Wendy Szymczak¹⁶⁰, Saranathan Rajagopalan¹⁶⁰, Emmanuel Asare¹⁶⁰, William R. Jacobs Jr.¹⁶⁰, Haejin In^{160,161}, Roni Bollag¹⁶², Aileen Lopez¹⁶², Edward J. Kruse¹⁶³, Joseph White¹⁶³, David Y. Graham¹⁶⁴, Charlotte Lane¹⁶⁵, Yang Gao¹⁶⁵, Patricia I. Fields¹⁶⁵, Benjamin D. Gold¹⁶⁶, Marcia Cruz-Correa^{167,168}, María González-Pons¹⁶⁷, Luz M. Rodríguez¹⁶⁹, Vo Phuoc Tuan¹⁷⁰, Ho Dang Quy Dung¹⁷⁰, Tran Thanh Binh¹⁷⁰, Tran Thi Huyen Trang¹⁷¹, Vu Van Khien¹⁷¹, Xiongfong Chen¹⁷², Castle Raley¹⁷³, Bailey Kessing¹⁷³, Yongmei Zhao¹⁷², Bao Tran¹⁷³, Andrés J. Gutiérrez-Escobar⁴, Yunhu Wan³, Belynda Hicks³, Bin Zhu³, Kai Yu⁴, Bin Zhu⁴, Meredith Yeager³, Amy Hutchinson³, Kedest Teshome³, Kristie Jones³, Wen Luo³, Quentin Jehanne⁴², Yukako Katsura¹⁷⁴, Patricio Gonzalez-Hormazabal¹⁷⁵, Xavier Didelot¹⁷⁶, Sam Sheppard¹⁷⁷, Eduardo Tarazona-Santos¹⁷⁸, Leonardo Mariño-Ramírez¹⁷⁹, John T. Loh¹⁰, Steffen Backert¹⁸⁰, Michael Naumann¹⁸¹, Christian C. Abnet⁴, Annemieke Smet¹⁸², Douglas E. Berg¹⁸³, Álvaro Chiner-Oms^{184,185}, Iñaki Comas^{185,186}, Francisco José Martínez-Martínez¹⁸⁶, Roxana Zamudio^{178,187}, Philippe Lehours⁴², Francis Megraud⁴², Koji Yahara¹⁸⁸, Martin J. Blaser¹⁸⁹, Tamas Vincze¹⁹⁰, Richard D. Morgan¹⁹⁰, Richard J. Roberts¹⁹⁰, Stephen J. Chanock⁴, John P. Dekker¹⁹¹, Javier Torres¹⁹², Timothy L. Cover^{10,193}, Mehwish Noureen¹⁹⁴, Wolfgang Fischer¹²², Filipa F. Vale^{195,196}, Joshua L. Cherry^{197,198}, Naoki Osada¹⁹⁹, Masaki Fukuyo²⁰⁰, Masanori Arita²⁰¹, Yoshio Yamaoka^{14,164}, Ichizo Kobayashi²⁰², Ikuo Uchiyama²⁰³, Daniel Falush⁹, M. Constanza Camargo^{4,205} & Charles S. Rabkin^{4,205}

¹⁰Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA. ¹¹Department of Natural and Life Sciences, Faculty of Sciences, University of Algiers 1 Benyoucef Benkhedda, Algiers, Algeria. ¹²Departamento de Medicina Interna, Hospital Alemán, Buenos Aires, Argentina. ¹³Department of Gastroenterology, Dhaka Medical College and Hospital, Dhaka, Bangladesh. ¹⁴Department of Environmental and Preventive Medicine, Oita University Faculty of Medicine, Yufu, Japan. ¹⁵Department of Pathobiology, Pharmacology and Zoological Medicine, Faculty of Veterinary Medicine, Ghent University, Ghent, Belgium. ¹⁶Hospital Universitario Japones, Santa Cruz de la Sierra, Bolivia. ¹⁷A.C. Camargo Cancer Center, São Paulo, São Paulo, Brazil. ¹⁸Núcleo de Pesquisas em Oncologia, Universidade Federal do Pará, Belém, Pará, Brazil. ¹⁹Medical University of Sofia, Sofia, Bulgaria. ²⁰Department of Biochemistry, University of Dschang, Dschang, Cameroon. ²¹Faculty of Medicine and Dentistry, Department of Medicine, University of Alberta, Edmonton, AB, Canada. ²²Queen's University, Kingston, ON, Canada. ²³Department of Hematology and Oncology, Faculty of Medicine, Pontificia Universidad Católica de Chile, Santiago, Chile. ²⁴Department of Pediatric Gastroenterology and Nutrition, Faculty of Medicine, Pontificia Universidad Católica de Chile, Santiago, Chile. ²⁵Department of Gastroenterology, Faculty of Medicine, Pontificia Universidad Católica de Chile, Santiago, Chile. ²⁶Facultad de Ciencias Biológicas, Universidad de Concepción, Concepción, Chile. ²⁷Cáncer Lab, Departamento de Ciencias Biomédicas, Facultad de Medicina, Universidad Católica del Norte (Coquimbo), Chile, Coquimbo, Chile. ²⁸Hospital Hanga Roa, Easter Island, Chile. ²⁹Grupo de Investigación en Biología del Cáncer, Instituto Nacional de Cancerología, Bogotá DC, Colombia. ³⁰Departamento de Biología, Universidad de Nariño, Pasto, Nariño, Colombia. ³¹Escuela de Medicina, Universidad del Valle, Cali, Valle, Colombia. ³²Division of Comparative Medicine, Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA. ³³Instituto de Investigaciones en Salud, Universidad de Costa Rica, San Jose, Costa Rica. ³⁴Laboratorio de Patología General y Gastrointestinal (LABGIPAT), San Jose, Costa Rica. ³⁵Servicio de Gastroenterología y Endoscopia Digestiva, Hospital Clínica Bíblica, San Jose, Costa Rica. ³⁶Faculty of Medicine, Osaka Metropolitan University, Osaka, Japan. ³⁷Faculty of Medicine, University of Mbuji-Mayi, Mbuji-Mayi, Kasai-Oriental, Democratic Republic of the

Congo. ³⁸Faculty of Medicine, University of Kinshasa, Kinshasa, Democratic Republic of the Congo. ³⁹Department of Parasitology, National Institute of Biomedical Research, Kinshasa, Democratic Republic of the Congo. ⁴⁰Instituto de Microbiología y Parasitología, Universidad Autónoma de Santo Domingo, Santo Domingo, Dominican Republic. ⁴¹Dominican-Japanese Digestive Disease Center, Dr Luis E. Aybar Health and Hygiene City, Santo Domingo, Dominican Republic. ⁴²Bordeaux Institute of Oncology, BRIC U1312, INSERM, Bordeaux, and National Reference Center for Campylobacters & Helicobacters, CHU de Bordeaux, Bordeaux, France. ⁴³Medical Research Council Unit, The Gambia at the London School of Hygiene & Tropical Medicine, Banjul, The Gambia. ⁴⁴Department of Gastroenterology, Hepatology and Infectious Diseases, Otto-von-Guericke University Magdeburg, Magdeburg, Germany. ⁴⁵Department of Biochemistry, School of Biological Sciences, University of Cape Coast, Cape Coast, Central Region, Ghana. ⁴⁶Department of Internal Medicine and Therapeutics, School of Medical Sciences, University of Cape Coast, Cape Coast, Central Region, Ghana. ⁴⁷Department of Molecular Biology and Biotechnology, School of Biological Sciences, University of Cape Coast, Cape Coast, Central Region, Ghana. ⁴⁸Department of Biology Education, Faculty of Science Education, University of Education, Winneba, Ghana. ⁴⁹Laboratory of Medical Microbiology, Hellenic Pasteur Institute, Athens, Greece. ⁵⁰Department of Gastroenterology, Alexandra Hospital, Athens, Greece. ⁵¹Department of Gastroenterology, Athens Medical, P. Faliron Hospital, Athens, Greece. ⁵²Facultad de Ciencias Médicas, University of San Carlos of Guatemala, Guatemala City, Guatemala. ⁵³Unidad de Gastroenterología, Hospital Roosevelt, Guatemala City, Guatemala. ⁵⁴Gastrocentro, S.A., Guatemala City, Guatemala. ⁵⁵Departamento de Medicina Interna, Hospital de Occidente, Santa Rosa de Copán, Honduras. ⁵⁶School of Medicine, University of Alabama at Birmingham (UAB), Birmingham, AL, USA. ⁵⁷Landspítali – The National University Hospital of Iceland, Reykjavík, Iceland. ⁵⁸University of Iceland, Reykjavík, Iceland. ⁵⁹Department of Microbiology, Kasturba Medical College, Manipal Academy of Higher Education, Manipal, Karnataka, India. ⁶⁰Department of Medicine, University of Cambridge, Cambridge, UK. ⁶¹Universitas Airlangga, Surabaya, East Java, Indonesia. ⁶²University of Indonesia, Jakarta, Indonesia. ⁶³Digestive Disease Research Institute, Tehran University of Medical Sciences, Tehran, Iran. ⁶⁴Clinical Microbiology Laboratory and Research Institute, Tzafon Medical Center, affiliated with Azrieli Faculty of Medicine, Bar Ilan University, Poriya, Israel. ⁶⁵Azrieli Faculty of Medicine, Bar Ilan University, Safed, Israel. ⁶⁶Pediatric Gastroenterology and Nutrition Unit, Tzafon Medical Center, Poriya, Israel. ⁶⁷Unit of Immunopathology and Oncological Biomarkers, Centro di Riferimento Oncologico di Aviano, Aviano, Italy. ⁶⁸Unit of Oncological Gastroenterology, Centro di Riferimento Oncologico di Aviano, Aviano, Italy. ⁶⁹Unit of Pathology, Centro di Riferimento Oncologico di Aviano, Aviano, Italy. ⁷⁰Department of Gastroenterology and Metabolism, Nagoya City University Graduate School of Medical Sciences, Nagoya, Japan. ⁷¹School of Medicine, Kyorin University, Mitaka, Tokyo, Japan. ⁷²Jordan University of Science and Technology, Ar-Ramtha, Jordan. ⁷³Department of Surgical Diseases Internship, Astana Medical University, Nur-Sultan, Kazakhstan. ⁷⁴Kyrgyz State Medical Academy, Bishkek, Kyrgyzstan. ⁷⁵Center for Gastric Cancer, National Cancer Center, Goyang, South Korea. ⁷⁶Department of Internal Medicine, Chung-Ang University Hospital, Seoul, South Korea. ⁷⁷College of Medicine, Seoul National University, Seoul, South Korea. ⁷⁸Institute of Clinical and Preventive Medicine, University of Latvia, Riga, Latvia. ⁷⁹Digestive Diseases Centre GASTRO, Riga, Latvia. ⁸⁰Riga East University Hospital, Riga, Latvia. ⁸¹Department of Gastroenterology, Institute for Digestive Research, Medical Academy, Lithuanian University of Health Sciences, Kaunas, Lithuania. ⁸²Department of Medical Microbiology, Faculty of Medicine, Universiti Malaya, Kuala Lumpur, Malaysia. ⁸³Medical Education Research and Development Unit, Faculty of Medicine, Universiti Malaya, Kuala Lumpur, Malaysia. ⁸⁴Departamento de Patología, Instituto Nacional de Cancerología, Mexico City, Mexico. ⁸⁵Servicio de Endoscopia, Instituto Nacional de Cancerología, Mexico City, Mexico. ⁸⁶Defence Services General Hospital, Yangon, Yangon, Yangon Region, Myanmar. ⁸⁷Nippon Medical School, Tokyo, Japan. ⁸⁸Department of Gastroenterology, Maharajgunj Medical Campus, Tribhuvan University Teaching Hospital, Kathmandu, Nepal. ⁸⁹Division of Health Sciences, Abu Dhabi Women's Campus, Higher Colleges of Technology, Abu Dhabi, United Arab Emirates. ⁹⁰Department of Community Medicine, Babcock University, Ilishan, Ogun State, Nigeria. ⁹¹Department of Medicine, Babcock University, Ilishan, Ogun State, Nigeria. ⁹²Department of Medicine, College of Medicine, University of Ibadan, Ibadan, Oyo, Nigeria. ⁹³Department of Nursing, Crescent University, Abeokuta, Ogun State, Nigeria. ⁹⁴University Teaching Hospital, University of Jos, Jos, Plateau, Nigeria. ⁹⁵University of Calabar Teaching Hospital, Calabar, Cross River, Nigeria. ⁹⁶University of Nigeria Teaching Hospital, Ituku-Ozalla, Enugu State, Nigeria. ⁹⁷Department of Internal Medicine, Federal Medical Center Abeokuta, Abeokuta, Ogun State, Nigeria. ⁹⁸Faculty of Health Sciences, Universidad Científica del Sur, Lima, Peru. ⁹⁹Departamento de investigación, Instituto Nacional de Enfermedades Neoplásicas, Lima, Peru. ¹⁰⁰Universidad Peruana Cayetano Heredia, Lima, Peru. ¹⁰¹Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, USA. ¹⁰²Department of Microbiology, Wrocław Medical University, Wrocław, Poland. ¹⁰³Department of Surgery Teaching, Wrocław Medical University, Wrocław, Poland. ¹⁰⁴Department of Pharmaceutical Microbiology, Medical University of Lublin, Lublin, Poland. ¹⁰⁵Department of Gastroenterology with Endoscopic Unit, Medical University of Lublin, Lublin, Poland. ¹⁰⁶Instituto Nacional de Saúde Dr. Ricardo Jorge, Lisboa, Portugal. ¹⁰⁷Instituto de Patologia e Imunologia Molecular da Universidade do Porto, Porto, Portugal. ¹⁰⁸Instituto de Investigação e Inovação em Saúde, Universidade do Porto, Porto, Portugal. ¹⁰⁹Faculdade de Medicina da Universidade do Porto, Porto, Portugal. ¹¹⁰Department of Pancreatic, Biliary and Upper Digestive Tract Disorders, A. S. Loginov Moscow Clinical Scientific Center, Moscow, Russia. ¹¹¹Department of General Medical Practice and Family Medicine, Tver State Medical University, Moscow, Russia. ¹¹²Department of Propaedeutic of Internal diseases and Gastroenterology A.I. Yevdokimov Moscow State University of Medicine and Dentistry, Moscow, Russia. ¹¹³Department of Faculty Therapy and Gastroenterology, Omsk State Medical University, Omsk, Russia. ¹¹⁴Scientific Research Institute of Medical Problems of the North, Federal Research Centre “Krasnoyarsk Science Centre” of the Siberian Branch of Russian Academy of Science, Krasnoyarsk, Russia. ¹¹⁵Cancer Science Institute of Singapore, National University of Singapore, Singapore, Singapore. ¹¹⁶Cancer and Stem Cell Biology Program, Duke NUS Medical School, Singapore, Singapore. ¹¹⁷Genome Institute of Singapore, Agency for Science, Technology and Research, Singapore, Singapore. ¹¹⁸Department of Medicine, Yong Loo Lin School of Medicine, National University of Singapore, Singapore, Singapore. ¹¹⁹Department of Gastroenterology and Hepatology, National University Health System, Singapore, Singapore. ¹²⁰Chris Hani Baragwanath Academic Hospital, Johannesburg, South Africa. ¹²¹University of the Witwatersrand, Johannesburg, South Africa. ¹²²Max von Pettenkofer Institute of Hygiene and Medical Microbiology, Faculty of Medicine, LMU Munich, Munich, Germany. ¹²³Hospital Universitario Donostia, San Sebastian, Spain. ¹²⁴Department of Gastroenterology, Biodonostia Health Research Institute - Donostia University Hospital, Universidad del País Vasco (UPV/EHU), Centro de Investigación Biomédica en Red de Enfermedades Hepáticas y Digestivas (CIBERehd), San Sebastian, Spain. ¹²⁵Digestive Diseases Unit, Parc Taulí Hospital Universitari. Institut d'Investigació i Innovació Parc Taulí (I3PT-CERCA), Universitat Autònoma de Barcelona, Barcelona, Spain. ¹²⁶Centro de Investigación Biomédica en Red de Enfermedades Hepáticas y Digestivas, Instituto de Salud Carlos III, Madrid, Spain. ¹²⁷Lozano Blesa University Clinic Hospital, Zaragoza, Spain. ¹²⁸Aragon Health Research Institute, Zaragoza, Spain. ¹²⁹Miguel Servet University Hospital, Zaragoza, Spain. ¹³⁰Hospital General de Granollers, Barcelona, Spain. ¹³¹Gastroenterology Department, Hospital Clínic de Barcelona, Centro de Investigación Biomédica en Red de Enfermedades Hepáticas y Digestivas, Facultat de Medicina i Ciències de la Salut, Universitat de Barcelona, Barcelona, Spain. ¹³²Hospital Universitari de Bellvitge, L'Hospitalet de Llobregat, Barcelona, Spain. ¹³³Department of Gastroenterology, Hospital Universitario Central de Asturias, Oviedo, Asturias, Spain. ¹³⁴Diet, Microbiota and Health Group, Instituto de Investigación Sanitaria del Principado de Asturias, Oviedo, Asturias, Spain. ¹³⁵Pharmacology Group, Instituto de Investigación Sanitaria del Principado de Asturias, Oviedo, Asturias, Spain. ¹³⁶Instituto Universitario de Oncología del Principado de Asturias, Oviedo, Asturias, Spain. ¹³⁷Hospital General Universitario Gregorio Marañón, Madrid, Spain. ¹³⁸Gastroenterology Unit, Hospital Universitario de La Princesa, Instituto de Investigación Sanitaria Princesa, Madrid, Spain. ¹³⁹Universidad Autónoma de Madrid, Centro de Investigación Biomédica en Red de Enfermedades Hepáticas y Digestivas, Madrid, Spain. ¹⁴⁰Gastroenterology

Department, Hospital Universitario de Navarra, Pamplona, Navarra, Spain. ¹⁴¹Hospital de Leon, Leon, Spain. ¹⁴²Institut of Biomedicine, University of León, Consortium for Biomedical Research in Epidemiology and Public Health, Leon, Spain. ¹⁴³Instituto de Investigación Sanitaria INCLIVA, Hospital Clínico Universitario de Valencia, Valencia, Spain. ¹⁴⁴Biochemistry Department, Faculty of Science, King Abdulaziz University, Jeddah, Saudi Arabia. ¹⁴⁵Faculty of Medical laboratory Science, Sudan University of Science and Technology, Khartoum, Sudan. ¹⁴⁶Centre for Translational Microbiome Research, Department of Microbiology, Tumor and Cell Biology, Karolinska Institutet, Solna, Sweden. ¹⁴⁷University of Skövde, Skövde, Sweden. ¹⁴⁸Institute for Infectious Diseases, University of Bern, Bern, Switzerland. ¹⁴⁹Clinical Bacteriology/Mycology Unit, University Hospital Basel, Basel, Switzerland. ¹⁵⁰Institute of Medical Microbiology, University of Zurich, Zürich, Switzerland. ¹⁵¹Clinic for Gastroenterology and Hepatology, University Hospital Zurich, Zürich, Switzerland. ¹⁵²College of Medicine, National Taiwan University, Taipei City, Taiwan. ¹⁵³Medical Microbiology Department, Cerrahpasa Medical Faculty, Istanbul University-Cerrahpasa, İstanbul, Türkiye. ¹⁵⁴General Surgery Department, Cerrahpasa Medical Faculty, Istanbul University-Cerrahpasa, İstanbul, Türkiye. ¹⁵⁵Medical Pathology Department, Cerrahpasa Medical Faculty, Istanbul University-Cerrahpasa, İstanbul, Türkiye. ¹⁵⁶Lawrence General Hospital, Lawrence, MA, USA. ¹⁵⁷Carson Tahoe Regional Medical Center, Carson City, NV, USA. ¹⁵⁸White River Medical Center, Batesville, AR, USA. ¹⁵⁹Department of Medicine, Stanford University, Stanford, CA, USA. ¹⁶⁰Albert Einstein College of Medicine, Bronx, NY, USA. ¹⁶¹Rutgers Cancer Institute, New Brunswick, NJ, USA. ¹⁶²Georgia Cancer Center's Biorepository, Augusta University, Augusta, Georgia. ¹⁶³Augusta University Medical Center, Augusta, Georgia. ¹⁶⁴Section of Gastroenterology and Hepatology, Department of Medicine, Baylor College of Medicine, Houston, TX, USA. ¹⁶⁵Enteric Diseases Laboratory Branch, Division of Foodborne, Waterborne, and Environmental Diseases, Centers for Disease Control and Prevention, Atlanta, GA, USA. ¹⁶⁶Gi Care for Kids, LLC, Children's Center for Digestive Healthcare, LLC, Atlanta, GA, USA. ¹⁶⁷University of Puerto Rico Comprehensive Cancer Center, San Juan, Puerto Rico. ¹⁶⁸University of Puerto Rico Medical Sciences Campus, San Juan, Puerto Rico. ¹⁶⁹Gastrointestinal and Other Cancers Research Group, Division of Cancer Prevention, National Cancer Institute, Rockville, MD, USA. ¹⁷⁰Department of Endoscopy, Cho Ray Hospital, Ho Chi Minh City, Vietnam. ¹⁷¹Department of Hepatogastroenterology, 108 Military Central Hospital, Hanoi, Vietnam. ¹⁷²Sequencing Facility Bioinformatics Group, Bioinformatics and Computational Science Directorate, Frederick National Laboratory for Cancer Research, Frederick, MD, USA. ¹⁷³Cancer Research Technology Program, Frederick National Laboratory for Cancer Research, Frederick, MD, USA. ¹⁷⁴Center for the Evolutionary Origins of Human Behavior, Kyoto University, Inuyama, Japan. ¹⁷⁵Instituto de Ciencias Biomédicas (ICBM), Facultad de Medicina, Universidad de Chile, Santiago, Chile. ¹⁷⁶School of Life Sciences, University of Warwick, Coventry, UK. ¹⁷⁷Department of Biology & Biochemistry, University of Bath, Bath, UK. ¹⁷⁸Departamento de Genética, Ecología e Evolução, Universidade Federal de Minas Gerais, Belo Horizonte, Brazil. ¹⁷⁹National Institute on Minority Health and Health Disparities, Bethesda, MD, USA. ¹⁸⁰Division of Microbiology, Department of Biology, Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, Germany. ¹⁸¹Institute of Experimental Internal Medicine, Otto-von-Guericke-Universität Magdeburg, Magdeburg, Germany. ¹⁸²Laboratory of Experimental Medicine and Pediatrics, Faculty of Medicine and Health Sciences, University of Antwerp, Antwerp, Belgium. ¹⁸³Department of Molecular Microbiology, Washington University School of Medicine, St. Louis, MO, USA. ¹⁸⁴Genomics and Health Area, FISABIO – Public Health, Valencia, Spain. ¹⁸⁵CIBER in Epidemiology and Public Health, Madrid, Spain. ¹⁸⁶Tuberculosis Genomics Unit, Instituto de Biomedicina de Valencia, Consejo Superior de Investigaciones Científicas, Valencia, Spain. ¹⁸⁷Quadram Institute Bioscience, Norwich, UK. ¹⁸⁸National Institute of Infectious Diseases, Tokyo, Japan. ¹⁸⁹Center for Advanced Biotechnology and Medicine, Rutgers University, New Brunswick, NJ, USA. ¹⁹⁰New England Biolabs, Ipswich, MA, USA. ¹⁹¹Bacterial Pathogenesis and Antimicrobial Resistance Unit, National Institute of Allergy and Infectious Diseases, Bethesda, MD, USA. ¹⁹²Unidad de Investigación en Enfermedades Infecciosas y Parasitarias, UMAE Pediatría, Instituto de Seguro Social, Mexico City, Mexico. ¹⁹³Veterans Affairs Tennessee Valley Healthcare System, Nashville, TN, USA. ¹⁹⁴Department of Genetics, SOKENDAI University, Mishima, Shizuoka, Japan. ¹⁹⁵Pathogen Genome Bioinformatics and Computational Biology, Research Institute for Medicines, Faculty of Pharmacy, Universidade de Lisboa, Lisboa, Portugal. ¹⁹⁶Instituto de Biosistemas e Ciências Integrativas, Faculdade de Ciências, Universidade de Lisboa, Lisboa, Portugal. ¹⁹⁷National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA. ¹⁹⁸Division of International Epidemiology and Population Studies, Fogarty International Center, National Institutes of Health, Bethesda, MD, USA. ¹⁹⁹Faculty of Information Science and Technology, Hokkaido University, Sapporo, Japan. ²⁰⁰Department of Molecular Oncology, Chiba University, Chiba, Japan. ²⁰¹Bioinformation and DDBJ Center, National Institute of Genetics, Mishima, Shizuoka, Japan. ²⁰²Research Center for Micro-Nano Technology, Hosei University, Tokyo, Japan. ²⁰³National Institute for Basic Biology, National Institutes of Natural Sciences, Aichi, Japan.