

1 **A molecular barcode and online tool to identify and map imported infection**  
2 **with *Plasmodium vivax***

3 Hidayat Trimarsanto<sup>1,2</sup>, Roberto Amato<sup>3,4</sup>, Richard D Pearson<sup>3,4</sup>, Edwin Sutanto<sup>1</sup>, Rintis Noviyanti<sup>1</sup>, Leily  
4 Trianty<sup>1</sup>, Jutta Marfurt<sup>5</sup>, Zuleima Pava<sup>5</sup>, Diego F Echeverry<sup>6,7</sup>, Tatiana M Lopera-Mesa<sup>8</sup>, Lidia Madeline  
5 Montenegro<sup>8</sup>, Alberto Tobón-Castaño<sup>8</sup>, Matthew J Grigg<sup>5,9</sup>, Bridget Barber<sup>5,9</sup>, Timothy William<sup>9,10</sup>,  
6 Nicholas M Anstey<sup>5</sup>, Sisay Getachew<sup>11,12</sup>, Beyene Petros<sup>11</sup>, Abraham Aseffa<sup>12</sup>, Ashenafi Assefa<sup>13</sup>, Awab  
7 Ghulam Rahim<sup>14,15</sup>, Nguyen Hoang Chau<sup>16</sup>, Tran Tinh Hien<sup>16</sup>, Mohammad Shafiul Alam<sup>17</sup>, Wasif A Khan<sup>17</sup>,  
8 Benedikt Ley<sup>5</sup>, Kamala Thriemer<sup>5</sup>, Sonam Wangchuck<sup>18</sup>, Yaghoob Hamed<sup>19</sup>, Ishag Adam<sup>20</sup>, Yaobao  
9 Liu<sup>21,22</sup>, Qi Gao<sup>22</sup>, Kanlaya Sriprawat<sup>23</sup>, Marcelo U Ferreira<sup>24</sup>, Alyssa Barry<sup>25,26, 27,28</sup>, Ivo Mueller<sup>28,29</sup>, Eleanor  
10 Drury<sup>4</sup>, Sonia Goncalves<sup>4</sup>, Victoria Simpson<sup>3,4</sup>, Olivo Miotto<sup>3,4,13</sup>, Alistair Miles<sup>3,4</sup>, Nicholas J White<sup>13,30</sup>,  
11 Francois Nosten<sup>13,23,30</sup>, Dominic P Kwiatkowski<sup>3,4</sup>, Ric N Price<sup>\*5,13,30</sup>, Sarah Auburn<sup>\*5,13,30</sup>

12

- 13 1. Eijkman Institute for Molecular Biology, Jakarta 10430, Indonesia  
14 2. Agency for Assessment and Application of Technology, Jl. MH Thamrin 8, Jakarta 10340,  
15 Indonesia  
16 3. Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, Old Road Campus,  
17 Oxford, OX3 7LF, UK  
18 4. Wellcome Sanger Institute, Hinxton, Cambridge CB10 1SA, UK  
19 5. Global and Tropical Health Division, Menzies School of Health Research and Charles Darwin  
20 University, Darwin, Northern Territory 0811, Australia  
21 6. Centro Internacional de Entrenamiento e Investigaciones Medicas, CIDEIM, Cali, Colombia  
22 7. Universidad Icesi, Cali, Colombia  
23 8. Grupo Malaria, Facultad de Medicina, Universidad de Antioquia, Medellin, Colombia

- 24 9. Infectious Diseases Society Sabah-Menzies School of Health Research Clinical Research Unit,  
25 Kota Kinabalu, Sabah, Malaysia
- 26 10. Clinical Research Centre, Queen Elizabeth Hospital, Sabah, Malaysia
- 27 11. College of Natural Sciences, Addis Ababa University, P.O. Box 52, Addis Ababa, Ethiopia
- 28 12. Armauer Hansen Research Institute, P.O. Box 1005, Jimma Road, Addis Ababa, Ethiopia
- 29 13. Ethiopian Public Health Institute, Addis Ababa, Ethiopia
- 30 14. Mahidol-Oxford Tropical Medicine Research Unit, Mahidol University, Bangkok 10400, Thailand
- 31 15. Nangarhar Medical Faculty, Nangarhar University, Ministry of Higher Education, Afghanistan
- 32 16. Centre for Tropical Medicine, Oxford University Clinical Research Unit, 764 Vo Van Kiet, W.1,  
33 Dist.5, Ho Chi Minh City, Vietnam
- 34 17. Infectious Diseases Division, International Centre for Diarrheal Diseases Research, Bangladesh  
35 Mohakhali, Dhaka, 1212, Bangladesh
- 36 18. Royal Center for Disease Control, Department of Public Health, Ministry of Health, Thimphu,  
37 Bhutan
- 38 19. Infectious and Tropical Diseases Research Center, Hormozgan University of Medical Sciences,  
39 Bandar Abbas, Hormozgan Province, Iran
- 40 20. Faculty of Medicine, University of Khartoum, P.O. Box 102, Khartoum, Sudan
- 41 21. Medical College of Soochow University, Suzhou, Jiangsu, People's Republic of China
- 42 22. Key Laboratory of National Health and Family Planning Commission on Parasitic Disease Control  
43 and Prevention, Jiangsu Provincial Key Laboratory on Parasite and Vector Control Technology,  
44 Jiangsu Institute of Parasitic Diseases, Wuxi, Jiangsu, People's Republic of China
- 45 23. Shoklo Malaria Research Unit, Mae Sot, Tak 63110, Thailand
- 46 24. Department of Parasitology, Institute of Biomedical Sciences, University of São Paulo, São Paulo,  
47 Brazil
- 48 25. School of Medicine, Deakin University, Geelong, VIC, Australia

- 49 26. Burnet Institute, Melbourne, VIC, Australia
- 50 27. Population Health and Immunity Division, The Walter and Eliza Hall Institute of Medical  
51 Research, Parkville, VIC, Australia
- 52 28. Department of Medical Biology, The University of Melbourne, Parkville, VIC, Australia
- 53 29. Department of Parasites and Insect Vectors, Institut Pasteur, Paris, France
- 54 30. Centre for Tropical Medicine and Global Health, Nuffield Department of Medicine, University of  
55 Oxford, OX3 7LJ, UK
- 56
- 57 \*Corresponding author: Dr Sarah Auburn, Menzies School of Health Research, PO Box 41096, Casuarina,  
58 Darwin, NT 0811, Australia; Tel: (+61) 8 8946 8503
- 59

## 60 **Abstract**

61 Imported cases present a considerable challenge to the elimination of malaria. Traditionally, patient  
62 travel history has been used to identify imported cases, but the long-latency liver stages confound this  
63 approach in *Plasmodium vivax*. Molecular tools to identify and map imported cases offer a more robust  
64 approach, that can be combined with drug resistance and other surveillance markers in high-  
65 throughput, population-based genotyping frameworks. Using a machine learning approach  
66 incorporating hierarchical FST (HFST) and decision tree (DT) analysis applied to 831 *P. vivax* genomes  
67 from 20 countries, we identified a 28-Single Nucleotide Polymorphism (SNP) barcode with high capacity  
68 to predict the country of origin. The Matthews correlation coefficient (MCC), which provides a measure  
69 of the quality of the classifications, ranging from -1 (total disagreement) to 1 (perfect prediction),  
70 exceeded 0.9 in 15 countries in cross-validation evaluations. When combined with an existing 37-SNP *P.*  
71 *vivax* barcode, the 65-SNP panel exhibits MCC scores exceeding 0.9 in 17 countries with up to 30%  
72 missing data. As a secondary objective, several genes were identified with moderate MCC scores

73 (median MCC range from 0.54-0.68), amenable as markers for rapid testing using low-throughput  
74 genotyping approaches. A likelihood-based classifier framework was established, that supports analysis  
75 of missing data and polyclonal infections. To facilitate investigator-lead analyses, the likelihood  
76 framework is provided as a web-based, open-access platform (vivaxGEN-geo) to support the analysis  
77 and interpretation of data produced either at the 28-SNP core or full 65-SNP barcode. These tools can  
78 be used by malaria control programs to identify the main reservoirs of infection so that resources can be  
79 focused to where they are needed most.

80

## 81 **Keywords**

82 *Plasmodium vivax*, malaria, imported malaria, geographic origin, surveillance, genotyping, genomics,  
83 molecular barcode

84

## 85 **Background**

86 The last three World Malaria Reports have revealed a disturbing rise in malaria cases, and, outside  
87 Subsaharan Africa, an increasing proportion of malaria due to *Plasmodium vivax*, undermining the  
88 painstaking efforts to reduce transmission over the past decade<sup>1</sup>. These trends highlight the urgent  
89 need for new surveillance tools, with greater attention to the non-falciparum species. In today's global  
90 climate, human populations are highly mobile, with imported cases of malaria confounding local control  
91 efforts and enhancing the risks of drug resistance spread and outbreaks. There is thus a critical need to  
92 develop tools that can help to determine where patients acquired their infection.

93 The challenge of imported infections is particularly pertinent for *P. vivax*, in view of the parasite's ability  
94 to form dormant liver stages (hypnozoites) that can reactivate weeks to months after the initial  
95 infection, as well as highly persistent, low density blood-stage infections<sup>2,3</sup>. The re-emergence of *P. vivax*  
96 in multiple regions where it was once almost eliminated serves as an important reminder of the need to  
97 maintain diligent surveillance of this species<sup>4</sup>. In low endemic settings where the incidence of local  
98 infections is declining, the relative proportion of imported cases generally rises, emphasizing the  
99 importance for surveillance tools that can identify imported *P. vivax* cases. Traditionally, imported cases  
100 have been identified and mapped using information on patient travel history, but the persistent blood  
101 stage infections and long-latency liver stages constrain the accuracy of this approach in *P. vivax*  
102 infections. Molecular tools to identify and map imported *P. vivax* cases offer an attractive complement  
103 to traditional epidemiological tools.

104 Amplicon-based sequencing has become a favored approach for targeted genotyping of malaria  
105 parasites. Using highly parallel sequencing platforms, such as the latest generation of Illumina  
106 sequencers, amplicon-based sequencing can be applied at moderate to high-throughput, with high  
107 accuracy and sensitivity. These platforms are flexible, allowing iterative enhancement of the Single  
108 Nucleotide Polymorphism (SNP) barcodes, which can provide an affordable genotyping approach,  
109 amenable to population-based molecular surveillance.

110 Previous studies have used mitochondrial and apicoplast markers to resolve imported from local *P. vivax*  
111 isolates, but the resolution of these organelles is constrained<sup>5-7</sup>. In 2015, a panel of 42 SNPs was  
112 identified to facilitate parasite finger-printing and geographic assignment<sup>8</sup>. The proposed 42-SNP Broad  
113 barcode was derived from genomic data available from 13 isolates from 7 countries. In the last 4 years,  
114 the repository of genomic data on *P. vivax* has expanded greatly, allowing further refinement of a  
115 parsimonious and widely applicable genotyping barcode.

116 The primary objective of our study was to develop molecular tools for identifying and characterizing  
117 imported *P. vivax* cases amenable to population-based surveillance frameworks, so that these data can  
118 be used to inform strategic decisions on where and how to deploy malaria control interventions. We  
119 tailored our molecular tools primarily to surveillance frameworks using Illumina or other high-  
120 throughput genotyping platforms. As a secondary objective, we sought to identify single gene regions  
121 permissible to lower throughput approaches for use in settings or situations where high-throughput or  
122 centralized approaches are not feasible. In addition, we provide informatics tools to support users in  
123 analyzing genotyping data produced at the barcode that can accommodate missing data and polyclonal  
124 infections.

125

## 126 **Materials and Methods**

127

### 128 **Overview of the marker selection approach**

129 A flow diagram outlining the steps involved in the marker selection process is provided in Figure 1a. In  
130 accordance with the multiplexing features of the Illumina platform, we sought to identify approximately  
131 50 new SNP-based markers to append to the recently published Broad barcode<sup>8</sup>, to provide a composite  
132 panel with  $\leq 100$  markers for country-level geographic assignment of *P. vivax* infections. The decision to  
133 append markers to the Broad barcode rather than selecting a de novo panel of SNPs was pragmatic,  
134 aimed at promoting consensus and continuity with existing molecular tools available to the vivax  
135 community. A likelihood-based classifier approach was chosen for the respective evaluation of marker  
136 sets and end-user data analysis tasks. This approach was chosen since it allows manual addition of  
137 specific SNP sets, such as the Broad barcode. Two selector algorithms, hierarchical FST (HFST) and  
138 decision-tree (DT), were implemented in the likelihood-based classifier framework to select SNPs with  
139 high country-level prediction values. The primary SNP selection method was the HFST selector, which  
140 leverages on the prior knowledge of the population structure to inform on a relatively parsimonious SNP

141 set with moderately high prediction. The DT selector, the secondary method, was used to select  
142 additional SNPs to append to the Broad barcode and the HFST panel for further enhancement of  
143 geographical prediction. A 10-fold cross-validation strategy was used to assess the performance of the  
144 selectors with the likelihood-based classification framework.

145 To achieve the secondary objective of the study, identifying single gene regions with moderate-to-high  
146 country-level resolution, simulations were run across individual genes using the HFST-0.75 (HFST with  
147 FST threshold of 0.75) selector model with the likelihood classifier. The top 20 genes with the highest  
148 pooled median Matthew Correlation Coefficient (MCC) scores for each selector model were reported  
149 (Figure 1b).

150

## 151 **Overview of the web-based data analysis and sharing platform**

152 To establish accessible informatics tools to support users to analyze and interpret their data, a platform  
153 was created incorporating data classification tools for determining the most likely country of origin of a  
154 sample using genetic data at a given barcode. Existing source code, developed for a microsatellite-based  
155 *P. vivax* data sharing and analysis platform<sup>9</sup>, was modified to create a new web-based platform  
156 (vivaxGEN-geo), to collate SNP data generated at the geographic barcode. A likelihood-based classifier  
157 approach was chosen for geographic assignment within the vivaxGEN-geo platform owing to the ability  
158 to i) incorporate manual SNP sets, ii) evaluate barcodes with missing data, and iii) evaluate heterozygous  
159 genotype calls.

160

## 161 **Likelihood-based classifier framework**

162 The custom classifier was developed to handle bi-allelic heterozygote calls for mix-infection cases by  
163 treating the samples as diploid samples, as well as missing data by treating as heterozygote calls. The

164 classifier was derived from Bernoulli Naive Bayes with modification to the likelihood equation and  
165 elimination of prior probability, since the distribution of our dataset did not reflect the distribution in  
166 nature, but rather the implication from sample and extracted DNA quality, as well as the characteristics  
167 of the original study such as duration and type of the study. Hence the classifier only depends on the  
168 likelihood of the SNP data. The likelihood equation was modified to handle the bi-allelic data as follows:

$$169 \quad L(\mathbf{X} | \mathbf{Ck}) = \prod_i^n p_{ki}^{x_i} \cdot (1 - p_{ki})^{(2-x_i)}$$

170 where  $\mathbf{X}$  is the SNP data set of a sample,  $Ck$  is a group (or a country),  $x_i$  is the number of alternate alleles  
171 at position  $i$  and  $p_{ki}$  is the frequency of the alternate allele at position  $i$  of country  $k$  counted as diploid  
172 samples.

173

#### 174 **SNP Selection**

175 To select optimal SNPs for country classification, a combination of the HFST and DT selector methods  
176 were employed. The DT selector utilized the Python-based scikit-learn package for the decision tree  
177 implementation, which employed an optimized version of the CART (Classification And Regression Tree)  
178 algorithm and Gini coefficient. To avoid overfitting, a minimum of 3 samples was required for a leaf  
179 node. The Hierarchical FST (HFST) approach worked by traversing across a bifurcating guide tree and  
180 selecting SNP with the highest FST between the two populations represented by the two nodes of the  
181 branch with the assumption that the SNP with the highest FST might differentiate those two  
182 populations. If the highest FST of a certain branch was lower than a given threshold during guide tree  
183 traversal, the DT method was employed to obtain additional SNPs to separate the given branch. To  
184 avoid overfitting, a maximum tree depth of 2 was set for this particular DT step. The HFST analysis in this  
185 study was undertaken using a guide tree constructed using Nei's population distance matrix  
186 implemented with a neighbor-joining algorithm.



187 The classification performance was measured with MCC for each country<sup>10</sup>. In addition, the pooled  
188 median, mean, minimum and first-quartile MCC were collected as additional measurements.

189 Three models, HFST-0.90 (HFST and DT with FST threshold of 0.90), HFST-0.95 (HFST and DT with FST  
190 threshold of 0.95), and pure DT were trained with the full dataset. For each of the three models, 500  
191 repeats were run to allow for different random seeds of the DT analysis, and the top 25 SNP sets with  
192 the highest aggregate minimum MCC score as evaluated by the likelihood classifier were obtained from  
193 each model. A stratified 100 repeat, 10-fold cross-validation was run on each of the 25 SNP sets from  
194 each model, and the best SNP set from each of model, as indicated by highest aggregate minimum MCC  
195 score within a repeat, was selected. To compare the Broad SNP panel to the three new SNP panels  
196 identified by the HFST-0.90, HFST-0.95 and pure DT selectors, a 500 repeat, stratified 10-fold cross-  
197 validation was undertaken on each SNP panel.

198

### 199 **Missing data evaluation**

200 To assess the durability of prediction performance of the SNP sets with missing data, a simulation was  
201 run by removing genotype data randomly. The Likelihood classifier was trained against the selected SNP  
202 sets using all samples. For each country, 25 samples were sampled randomly with replacement and  
203 missing genotype calls were added to the SNP sets in 10%, 20% and 30% proportions. The random  
204 samples were then subjected to the trained classifier. This process was run in 100 repeats and MCC-  
205 score of the prediction for each country was reported.

206

### 207 **Datasets**

208 The analysis included genomic data on *P. vivax* isolates collected from 26 countries. Published data were  
209 included from 19 countries derived from the European Nucleotide Archive<sup>11-15</sup>, and new data from 10  
210 countries (Supplementary Table 1, Supplementary Figure 1). New genomic data were derived from

211 patients recruited to partner studies in Afghanistan, Bangladesh, Bhutan, Colombia, Ethiopia, Indonesia,  
212 Iran, Malaysia, Sudan, and Vietnam. With the exception of Colombia, the patient sampling frameworks  
213 have been described previously<sup>11,12,14,16-20</sup>. The samples from Colombia were collected within the  
214 framework of cross-sectional epidemiological surveys conducted between 2013 and 2017. Whole  
215 genome sequencing, read alignment and variant calling were undertaken within the framework of a *P.*  
216 *vivax* community study in the Malaria Genomic Epidemiology Network (MalariaGEN)<sup>21</sup>. Data was  
217 derived from an open dataset of *Plasmodium vivax* genome variation comprising 2,671,112 discovered  
218 variants across 1,366 isolates (MalariaGEN manuscript in preparation). The data were initially filtered to  
219 derive a set of 670,962 high-quality bi-allelic SNPs with VQSLOD score >0, and minimum read depth and  
220 minimum minor allele count of 2. Individual genotype calls were defined as heterozygotes based on an  
221 arbitrary threshold of a minor allele ratio > 0.1 and a minimum of 2 reads for each allele; all other  
222 genotype calls were defined as homozygous for the major allele. A pair of isolates with distance less  
223 than 0.0005 (0.05%) were considered non-independent. Amongst non-independent sample pairs, the  
224 isolate with the higher percentage of genotype failures was removed from the dataset; this removal  
225 process was iterated until all non-independent isolates had been removed from the dataset. The  
226 samples and SNPs were then subjected to further filtering to eliminate missing data using information  
227 derived from a simulation which calculated the total number of SNPs with no missing data and the  
228 number of consecutive informative SNPs as defined by SNPs with minimum minor allele count (MAC) >2.  
229 The remaining dataset was defined as Dataset 1.

230 The isolates in Dataset 1 were initially assigned to national groups based on the country in which the  
231 patient presented at the clinic with the infection. The national-level groupings were evaluated further  
232 using country-level assignments derived from the genome-wide data classification with the likelihood  
233 classifier. Infections presenting with country classifications differing from the country of presentation  
234 were considered suspected imported infections and removed to produce Dataset 2.

235 Of the 42 Broad barcode SNPs, 37 SNPs were present in the 670K dataset (bi-allelic high-quality SNPs)  
236 and exhibited successful amplicon-based sequencing assays (personal communication, Wellcome Sanger  
237 Institute Core Sequencing Facility); these 37 SNPs were not present in dataset 1 or 2. A new dataset  
238 (Dataset 3) was prepared for evaluation of the Broad barcode comprising of samples with complete data  
239 across the 37 Broad barcode SNPs and partial data across SNPs selected from the HFST and DT  
240 algorithms.

241

## 242 **Software and Web Service Availability**

243 All custom, in-house scripts used for data filtering, simulation, analyses and visualization are available  
244 from <http://github.com/trmznt/vivaxgen-geo>. The VivaxGEN-geo web service provides a user-friendly  
245 online tool for country classification with all SNP sets, and is accessible at <https://geo.vivaxgen.org/>. The  
246 likelihood classifier provided on the online platform has been trained with 809 samples (dataset 4),  
247 consisting of all samples with complete data at all SNP sets. The classifier tool reports the three highest  
248 likelihoods for country of origin and their associated probabilities.

249

## 250 **Ethics**

251 All samples were collected with written informed consent from patients or their legal guardians. Ethical  
252 approvals for the published samples are detailed in the original papers<sup>11-15</sup>. Approvals for the newly  
253 represented studies are outlined in Supplementary Document 1.

254

## 255 **Results**

### 256 **Geographic clustering patterns using the genome-wide dataset**

257 The primary dataset (Dataset 1) was derived using the missing data simulation, to minimise genotype  
258 failures (Supplementary Figure 3), it comprised 854 high-quality samples and 294,628 high-quality

259 informative SNPs, with no missing data. The median percentage of heterozygous calls in each country  
260 ranged from 0.02% to 0.08%. Details on the geographic locations of the samples in dataset 1 are  
261 presented in Supplementary Table 1. Neighbor-joining analysis on dataset 1 revealed distinct geographic  
262 clustering of most countries (Supplementary Figure 4). Exceptions included the isolates from  
263 Afghanistan, Iran, India and Sri Lanka, which appeared to form a single cluster, warranting further  
264 analysis of this geographic region with larger sample sets. Although several isolates in border regions  
265 including Vietnam relative to Cambodia, and Thailand relative to Myanmar, overlapped between  
266 countries, the majority of isolates in these countries could be differentiated by national boundaries.  
267 Visual inspection of the neighbour-joining tree revealed potential imported cases. Using country-level  
268 assignments derived from the genome-wide data classification with the likelihood classifier and manual  
269 confirmation of the neighbor-joining tree, 21 isolates presented country classifications differing from the  
270 country of presentation (Supplementary Table 1). After exclusion of the imported cases, and countries  
271 represented by a single sample, a total of 831 isolates and 20 countries remained, constituting Dataset 2  
272 (Supplementary Table 1).

273

#### 274 **Performance of the Broad barcode, HFST and DT SNP selection**

275 When the HFST selector was applied with an FST threshold of 0.90 (HFST-0.90), a set of 28 SNPs (listed in  
276 Supplementary Table 2) were identified. This dataset exhibited median MCC scores exceeding 0.9 in all  
277 countries with the exception of Vietnam (0.75) and Cambodia (0.80). On increasing the FST threshold to  
278 0.95 (HFST-0.95), the HFST model identified 51 SNPs (listed in Supplementary Table 3), which displayed  
279 MCC scores exceeding 0.95 in all countries except for Vietnam (0.85) and Cambodia (0.87). Using the DT  
280 selector alone, 50 SNPs (listed in Supplementary Table 4) displayed comparable performance to the 51-  
281 SNP panel, with a slightly lower aggregate minimum MCC score.

282 The results of cross-validation of the classification performance of the five SNP panels (37-SNP Broad  
283 barcode, 28-SNP, 28-SNP plus Broad barcode (65-SNP), 50-SNP and 51-SNP panels) are illustrated in

284 Figure 2, and the MCC and F scores reflecting the consensus results of the cross-validation are  
285 summarized in Table 1. The performance, ranked from lowest to highest, was: 37-SNP Broad barcode  
286 (median MCC = 0.82), 28-SNP (MCC = 0.99), 50-SNP (MCC = 1.00), 65-SNP (MCC = 1.00), and 51-SNP  
287 (MCC = 1.00).

288

### 289 **Missing data simulations**

290 In the missing data simulations, genotyping failures had the greatest impact on the classification of  
291 samples from Cambodia and Vietnam (Figure 3). With 10% missing data, the median MCCs of the 28-SNP  
292 panel exceeded 0.9 in all countries, with exception of Vietnam (MCC = 0.80) and Cambodia (MCC =  
293 0.77). With this level of missing data, the addition of the 37 Broad SNPs (65-SNP panel) increased the  
294 median MCC to 0.83 in Vietnam and 0.82 in Cambodia. When missing data increased to 30%, the 65-SNP  
295 panel achieved median MCCs above 0.9 in most countries, with exception of Vietnam (MCC = 0.79) and  
296 Cambodia (MCC = 0.75). The 50- and 51-SNP panels both achieved MCC scores exceeding 0.95 for all  
297 countries except Cambodia (0.80-0.82) and Vietnam (0.83-0.85) with 10-30% missing data.

298

### 299 **Evaluation of single gene regions to predict country classification**

300 The suitability of single genes to predict country classifications were assessed by simulations of  
301 individual genes using HFST-0.75 selector model with the likelihood classifier framework. The top 20  
302 genes with the highest pooled median MCC scores for the HFST-0.75 are presented in Supplementary  
303 Table 5. The highest prediction capacity, with median MCC score of 0.68, was PVP01\_0302600, a gene  
304 coding a 11.5 Kb conserved protein with unknown function. The gene list also included three members  
305 of the *cysteine repeat modular protein family* (CRMP): CRMP1 (median MCC = 0.63), CRMP3 (MCC =  
306 0.57) and CRMP4 (MCC = 0.56).

307

## 308 **Discussion**

309 The primary objective of the study was to develop molecular tools amenable to population-based  
310 surveillance frameworks to identify and map imported *P. vivax* infections. Using machine-learning  
311 methods, 3 new SNP panels were identified with high country classification performance, able to  
312 distinguish imported *P. vivax* infections across a range of endemic scenarios. The most parsimonious  
313 panel, the 28-SNP barcode, exhibited high country classification, and can be appended to the 37 bi-  
314 allelic, assayable Broad barcode SNPs for marginal improvement in predictive capacity in samples with  
315 moderate levels of missing data. The combined 65-SNP barcode generated robust country classification  
316 in most endemic areas, even when the proportion of missing data rose to 30%. However, the validity of  
317 the 65-SNP barcode was lower in Cambodia and Vietnam, a likely reflection of the porous border  
318 between these two countries. Although the 50- and 51-SNP panels achieved better resolution in these  
319 areas, characterization of parasite transmission across borders with high levels of gene flow may be  
320 addressed better by the addition of markers suited to an analysis of identity-by-descent<sup>22</sup>. The  
321 application and wider validation of the 65-SNP barcode is underway, with amplicon-based sequencing  
322 assays already established for the 37 Broad barcode SNPs, and under development for the 28 new  
323 markers.

324 The analysis and interpretation of “real-world” genotyping data raises significant challenges from low-  
325 quality samples such as those collected on dried blood spots. In anticipation of these needs we  
326 established a likelihood-based framework with the capacity to deal with polyclonal infections as well as  
327 missing data. This framework has been integrated into the vivaxGEN-geo online platform, so that users  
328 can analyze and interpret their data without needing complex bioinformatics skills and avoiding the  
329 subjective visual inspection of neighbour-joining trees or principal component plots. Whilst the  
330 informatics tools implemented in vivaxGEN-geo are tailored to *P. vivax*, we anticipate that a similar  
331 approach can be adapted to other species. To facilitate wider application the source code will be made  
332 publicly available.

333 The variants in the 28-SNP panel are located in genes representing a range of functions, some of which  
334 may be unstable over time. Although our dataset represents one of the most geographically diverse  
335 panels of *P. vivax* isolates currently available, with representation of all of the major vivax-endemic  
336 regions, the predictive capacity of the derived tools are likely to be constrained by the geographic  
337 representation of the reference panel. In particular, representation from central and south America and  
338 the Indian subcontinent were limited in our data set. Despite this limitation the dataset used comprises  
339 good representation of isolates from areas of public health relevance, including the epicenter of  
340 chloroquine-resistant *P. vivax* in Papua Indonesia<sup>23,24</sup>. The likelihood-based framework is able to re-  
341 evaluate the predictive potential of current marker sets as new genomic data become available, so that  
342 the selected SNP panels can be refined further in an iterative process. Furthermore, as the reference  
343 panel expands with increasing data generated at the barcode SNPs, the accuracy of the likelihood-based  
344 classifications will improve.

345 In addition to the independent selection of SNPs, a number of informative genes were identified, each  
346 of which had moderately high geographic resolution power. Genotyping of these genes or gene regions  
347 are amenable to standard capillary sequencing, offering an alternative approach, albeit with slightly  
348 lower resolution, to define a parasite's geographic origin in settings where high-throughput genotyping  
349 facilities are unavailable. The genes with the greatest geographic resolution, included members of the  
350 *cysteine repeat modulator protein* (CRMP) family (CRMP1, CRMP3 and CRMP4) implicated in essential  
351 roles in parasite transmission from the mosquito to the human host<sup>25</sup>. It is plausible that the CRMP  
352 genes have maintained high geographic differentiation to ensure parasite adaptation to the local vector  
353 species. Although adaptations of these genes are likely to be temporally stable, the resolution of these loci  
354 may be constrained by the distributions of host *Anopheles* vector species.

355 In 2017, up to 100% of all confirmed malaria cases in 17 malaria-endemic countries in the Asia-Pacific  
356 region, the Middle East and the Americas, where *P. vivax* infections predominate, were reported as  
357 being infections<sup>1</sup>. Malaria control programs in these countries can utilize the information derived from

358 the molecular tools provided from our analysis to assess the efficacy of ongoing interventions in  
359 reducing local transmission, and to determine the major routes of infection importation. The tools have  
360 potential to reduce ambiguity for certifying malaria elimination by the World Health Organization,  
361 where one of the key requirements is the demonstration that all malaria cases detected in-country over  
362 at least three consecutive years were imported. For this purpose, countries approaching elimination will  
363 need to maintain archival samples for future molecular comparisons against putatively imported cases.

364 The molecular *P. vivax* geographic classification tools presented are designed to empower users in  
365 malaria-endemic countries to analyze and interpret locally produced genotyping data with comparison  
366 to globally available datasets. Amplicon-based sequencing of the full 65-SNP barcode is being developed  
367 and will be combined with other surveillance markers at central laboratories in endemic partner  
368 countries of the Asia Pacific Malaria Elimination Network ([www.apmen.org](http://www.apmen.org)). The data generated from  
369 these centers will inform researchers, National Malaria Control Programs and other key stakeholders on  
370 the incidence, epidemiology and key reservoirs of imported malaria, and, in doing so, help to target  
371 resources to where they are needed most.

372

## 373 **References**

- 374 1. WHO. World Malaria Report 2018. World Health Organization; Geneva 2018. (2018).
- 375 2. White, N.J. & Imwong, M. Relapse. *Adv Parasitol* **80**, 113-50 (2012).
- 376 3. Tripura, R. *et al.* Persistent *Plasmodium falciparum* and *Plasmodium vivax* infections in a  
377 western Cambodian population: implications for prevention, treatment and elimination  
378 strategies. *Malar J* **15**, 181 (2016).
- 379 4. Sattabongkot, J., Tsuboi, T., Zollner, G.E., Sirichaisinthop, J. & Cui, L. *Plasmodium vivax*  
380 transmission: chances for control? *Trends Parasitol* **20**, 192-8 (2004).
- 381 5. Iwagami, M. *et al.* Geographical origin of *Plasmodium vivax* in the Republic of Korea: haplotype  
382 network analysis based on the parasite's mitochondrial genome. *Malar J* **9**, 184 (2010).
- 383 6. Rodrigues, P.T. *et al.* Using mitochondrial genome sequences to track the origin of imported  
384 *Plasmodium vivax* infections diagnosed in the United States. *Am J Trop Med Hyg* **90**, 1102-8  
385 (2014).
- 386 7. Diez Benavente, E. *et al.* Genomic variation in *Plasmodium vivax* malaria reveals regions under  
387 selective pressure. *PLoS One* **12**, e0177134 (2017).



- 388 8. Baniecki, M.L. *et al.* Development of a single nucleotide polymorphism barcode to genotype  
389 *Plasmodium vivax* infections. *PLoS Negl Trop Dis* **9**, e0003539 (2015).
- 390 9. Trimarsanto, H. *et al.* VivaxGEN: An open access platform for comparative analysis of short  
391 tandem repeat genotyping data in *Plasmodium vivax* Populations. *PLoS Negl Trop Dis* **11**,  
392 e0005465 (2017).
- 393 10. Jurman, G., Riccadonna, S. & Furlanello, C. A comparison of MCC and CEN error measures in  
394 multi-class prediction. *PLoS One* **7**, e41882 (2012).
- 395 11. Auburn, S. *et al.* Genomic analysis of a pre-elimination Malaysian *Plasmodium vivax* population  
396 reveals selective pressures and changing transmission dynamics. *Nat Commun* **9**, 2585 (2018).
- 397 12. Auburn, S. *et al.* Genomic analysis of *Plasmodium vivax* in southern Ethiopia reveals selective  
398 pressures in multiple parasite mechanisms. *J Infect Dis* (2019).
- 399 13. Hupalo, D.N. *et al.* Population genomics studies identify signatures of global dispersal and drug  
400 resistance in *Plasmodium vivax*. *Nat Genet* **48**, 953-8 (2016).
- 401 14. Pearson, R.D. *et al.* Genomic analysis of local variation and recent evolution in *Plasmodium*  
402 *vivax*. *Nat Genet* **48**, 959-64 (2016).
- 403 15. Parobek, C.M. *et al.* Selective sweep suggests transcriptional regulation may underlie  
404 *Plasmodium vivax* resilience to malaria control measures in Cambodia. *Proc Natl Acad Sci U S A*  
405 **113**, E8096-E8105 (2016).
- 406 16. Wangchuk, S. *et al.* Where chloroquine still works: the genetic make-up and susceptibility of  
407 *Plasmodium vivax* to chloroquine plus primaquine in Bhutan. *Malar J* **15**, 277 (2016).
- 408 17. Ley, B. *et al.* G6PD Deficiency and Antimalarial Efficacy for Uncomplicated Malaria in  
409 Bangladesh: A Prospective Observational Study. *PLoS One* **11**, e0154015 (2016).
- 410 18. Hamedi, Y. *et al.* Molecular Epidemiology of *P. vivax* in Iran: High Diversity and Complex Sub-  
411 Structure Using Neutral Markers, but No Evidence of Y976F Mutation at *pvm-dr1*. *PLoS One* **11**,  
412 e0166124 (2016).
- 413 19. Getachew, S. *et al.* Variation in Complexity of Infection and Transmission Stability between  
414 Neighbouring Populations of *Plasmodium vivax* in Southern Ethiopia. *PLoS One* **10**, e0140780  
415 (2015).
- 416 20. Taylor, W.R.J. *et al.* Short-course primaquine for the radical cure of *Plasmodium vivax* malaria: a  
417 multicentre, randomised, placebo-controlled non-inferiority trial. *Lancet* (2019).
- 418 21. Malaria Genomic Epidemiology, N. A global network for investigating the genomic epidemiology  
419 of malaria. *Nature* **456**, 732-7 (2008).
- 420 22. Taylor, A.R., Jacob, P.E., Neafsey, D.E. & Buckee, C.O. Estimating Relatedness Between Malaria  
421 Parasites. *Genetics* **212**, 1337-1351 (2019).
- 422 23. Ratcliff, A. *et al.* Therapeutic response of multidrug-resistant *Plasmodium falciparum* and *P.*  
423 *vivax* to chloroquine and sulfadoxine-pyrimethamine in southern Papua, Indonesia. *Trans R Soc*  
424 *Trop Med Hyg* **101**, 351-9 (2007).
- 425 24. Price, R.N. *et al.* Global extent of chloroquine-resistant *Plasmodium vivax*: a systematic review  
426 and meta-analysis. *Lancet Infect Dis* **14**, 982-91 (2014).
- 427 25. Douradinha, B. *et al.* *Plasmodium* Cysteine Repeat Modular Proteins 3 and 4 are essential for  
428 malaria parasite transmission from the mosquito to the host. *Malar J* **10**, 71 (2011).
- 429 26. Battle, K.E. *et al.* Mapping the global endemicity and clinical burden of *Plasmodium vivax*, 2000-  
430 17: a spatial and temporal modelling study. *Lancet* **394**, 332-343 (2019).

431

## 432 Acknowledgements

433 We would like to thank the patients who contributed their samples to the study, local health facilities,  
434 and the health workers and field teams who assisted with the sample collections. We also thank the  
435 staff of the Wellcome Sanger Institute Sample Logistics, Sequencing, and Informatics facilities for their  
436 contributions.

437

## 438 **Financial Support**

439 Financial support for the study was provided by the Wellcome Trust (Senior Fellowship in Clinical  
440 Science awarded to R.N.P., 200909), the Australian Department of Foreign Affairs and Trade (TDCRRI  
441 72904), the Australian National Health and Medical Research Council (NHMRC) ('Improving Health  
442 Outcomes in the Tropical North: A multidisciplinary collaboration 'HOT North' Career Development  
443 Fellowship awarded to S.A.), and the Bill and Melinda Gates Foundation (OPP1164105). D.F.E received  
444 financial support from Colciencias -Colombia, call 656-2014 "Es Tiempo de Volver" award FP44842-503-  
445 2014. The patient sampling and metadata collection was funded by the Asia-Pacific Malaria Elimination  
446 Network (108-07), the Malaysian Ministry of Health (BP00500420), and the NHMRC (1037304 and  
447 1045156; Fellowships to N.M.A. [1042072 and 1135820], B.E.B. [1088738] and M.J.G. [1074795]). M.J.G  
448 was also supported by a 'Hot North' Earth Career Fellowship (1131932). M.U.F is supported by a senior  
449 researcher scholarship from the Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq),  
450 Brazil. The whole genome sequencing component of the study was supported by grants from the  
451 Medical Research Council and UK Department for International Development (M006212) and the  
452 Wellcome Trust (206194, 204911) awarded to D.P.K. This work was supported by the Australian  
453 Centre for Research Excellence on Malaria Elimination (ACREME), funded by the NHMRC (APP  
454 1134989).

455

## 456 **Figures**

### 457 **Figure 1. Overview of the marker selection approaches**

458 Flow diagrams illustrating the datasets, selector models and classification approaches used to identify  
459 and evaluate independent SNP panels (A) and single gene regions (B). Decision Tree (DT), HFST-0.90  
460 (HFST and DT with FST threshold of 0.90), HFST-0.95 (HFST and DT with FST threshold of 0.95) and HFST-  
461 0.75 (HFST and DT with FST threshold of 0.75) represent the SNP selector models. The DT, HFST-0.9 and  
462 HFST-0.95 SNP selector models were run in 500 repeats for SNP selection (A), and the HFST-0.75 model  
463 was run in 5 repeats for gene selection (B). For SNP selection (A), the top-25 SNP sets were selected  
464 from each model for a further 100 repeats of stratified cross-validation from which one SNP set was  
465 selected from each of the DT, HFST-0.9 and HFST-0.95 SNP selector models.

466

### 467 **Figure 2. Comparison between the 37-SNP Broad barcode, new marker panels and combined SNP sets**

468 The Broad-37 SNP set reflects 37 of the 42 Broad SNPs represented amongst the 294K high-quality SNPs.  
469 The SNP-28 SNP set reflects 28 high-performance SNPs derived from the HFST selector with FST  
470 threshold of 0.9. The SNP-28+Broad SNP set reflects the combined Broad-37 and SNP-28 SNP sets for a  
471 total of 65 SNPs. The SNP-50 set reflects the 50 SNPs selected by the Decision Tree selector. The SNP-51  
472 set reflects 51 high-performance SNPs from the HFST selector with threshold FST of 0.95. The boxplots  
473 present the MCC scores from 500 repeats with stratified 10-fold cross validation for each SNP set.  
474 Country labels are provided on the y-axis; MEDIAN, MEAN, Q1 (1st percentile) and MIN reflect the  
475 respective summary statistics for the pooled MCC scores across all countries.

476

### 477 **Figure 3. Simulation of missing data in the 28-SNP, 65-SNP, 50-SNP and 51-SNP panels**

478 Result of 200 repeats, 25 samples per country simulation of missing data (genotype fails) of 10%, 20%  
479 and 30% against the 37-SNP Broad barcode, 28-SNP set, 65-SNP set (28-SNP + Broad panel), 51-SNP and  
480 50-SNP set. The 65-SNP set demonstrated marginally better performance relative to the 28-SNP set with  
481 missing data. However, both the 50-SNP and 51-SNP panels outperformed the 65-SNP panel with  
482 missing data.

483

484 **Supplementary Figure 1. *P. vivax* prevalence map pinpointing the countries included in the study**

485 *P. vivax* prevalence map from the Malaria Atlas Project (*Plasmodium vivax* parasite rate in all ages  
486 globally (1-99) from (2000-2017)<sup>26</sup>, with counties included in dataset 2 demarked by stars.

487

488 **Supplementary Figure 2. Overview of the datasets**

489

490 **Supplementary Figure 3. Output from the data quality simulation**

491 The upper panel shows the number of complete SNPs (green), complete informative SNPs with minor  
492 allele count (MAC) = 1 (orange) and complete informative SNPs with MAC = 2 (red) against the number  
493 of samples. The lower panel shows the number of differences in SNPs between consecutive number of  
494 samples, with informative SNPs with MAC = 1 (blue) and informative SNPs with MAC = 2 (orange). The  
495 maximum of both MAC=1 and MAC=2 were 958 samples.

496

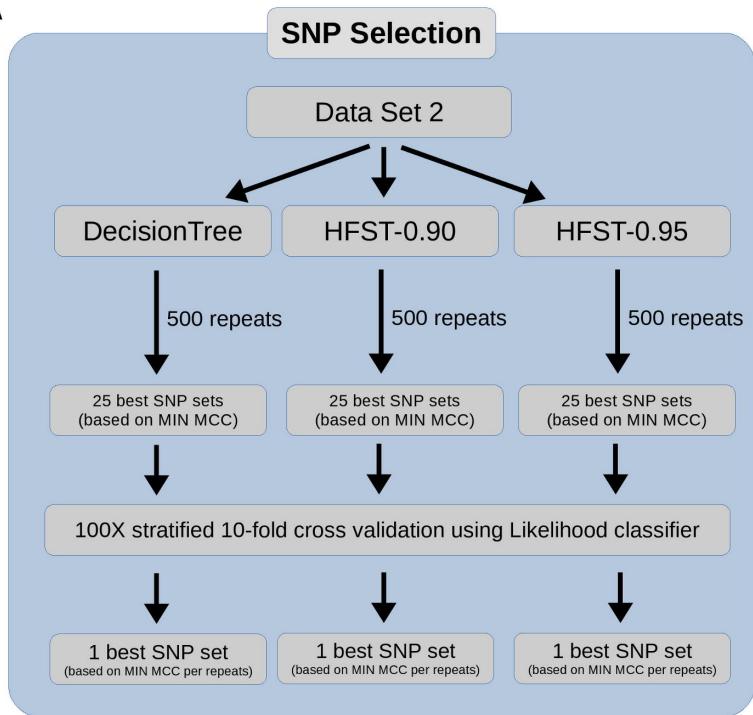
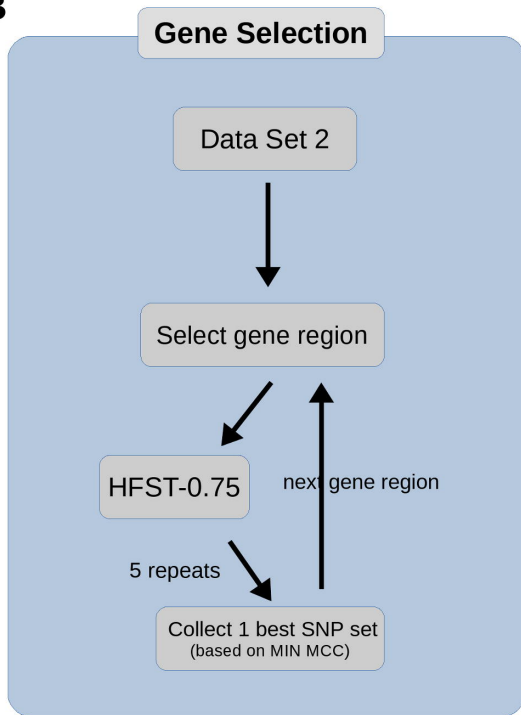
497 **Supplementary Figure 4. Neighbour-joining tree of the global dataset**

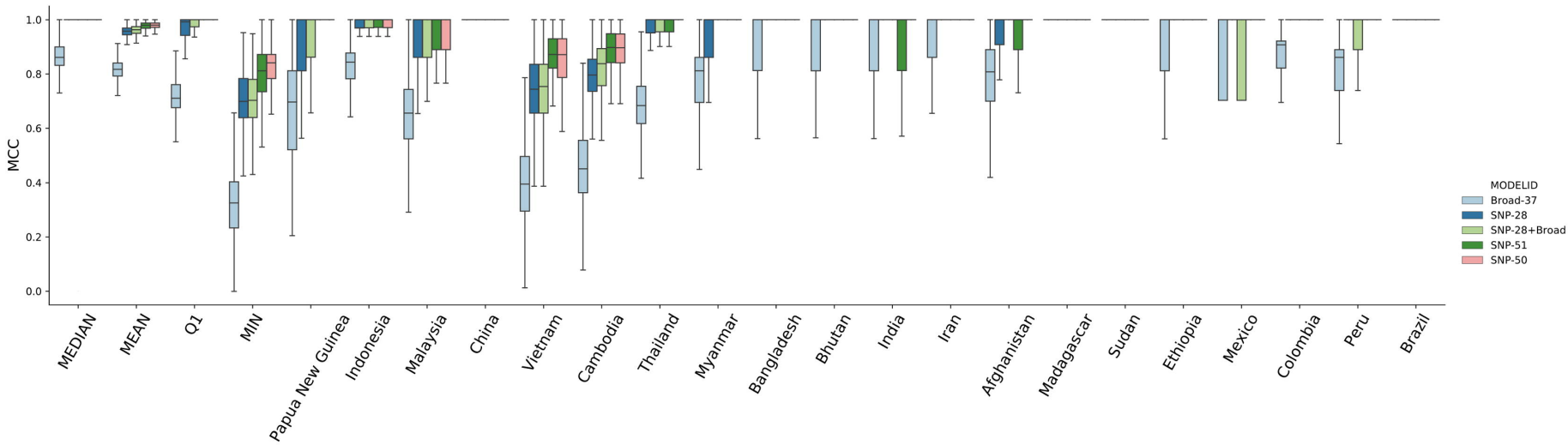
498 The tree was constructed using genotyping data from 854 samples at 294K SNPs.

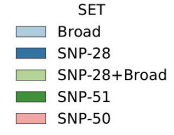
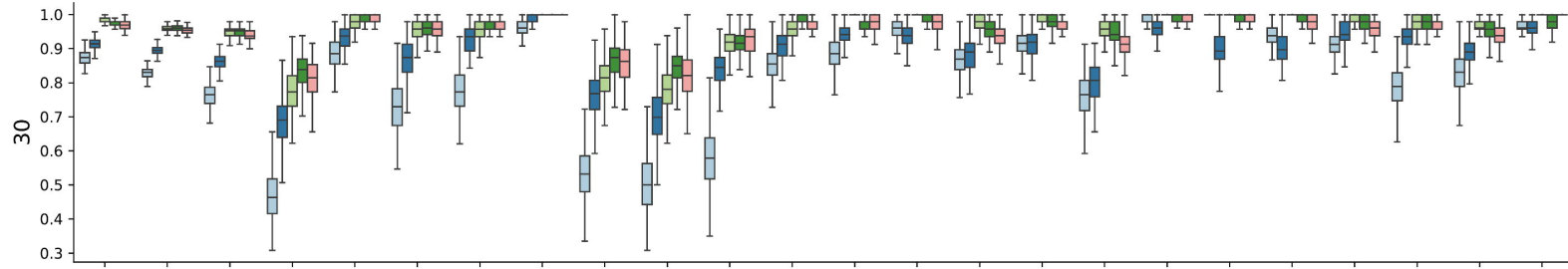
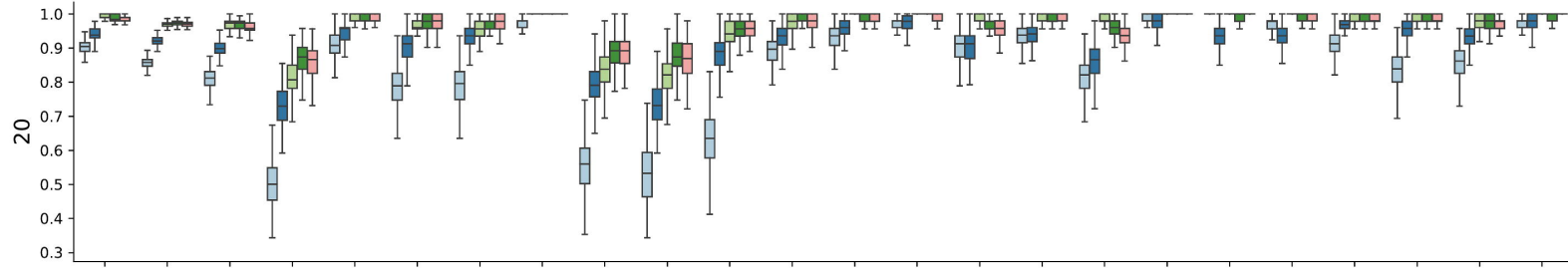
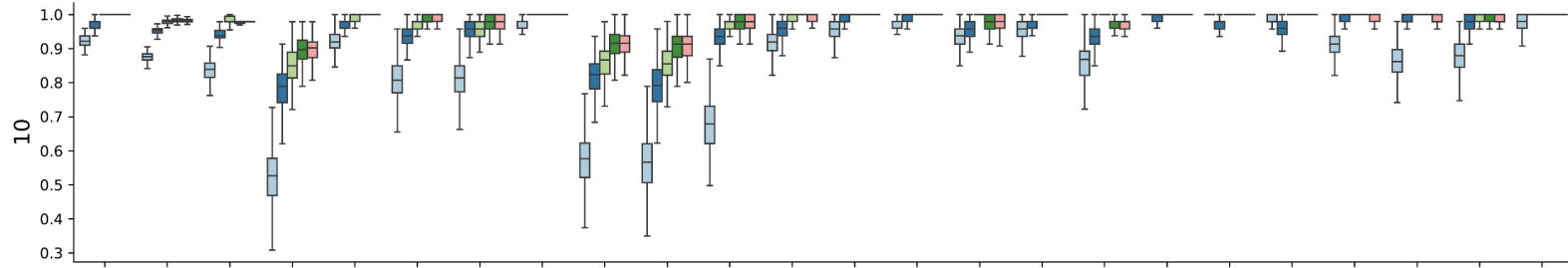
## Tables

**Table 1. Summary of MCC and F-scores from the consensus results of 500 repeats of the stratified 10-fold cross-validation of the SNP panels**

Population	37-SNP (Broad)		28-SNP		65-SNP		50-SNP		51-SNP	
	MCC	F	MCC	F	MCC	F	MCC	F	MCC	F
Afghanistan	0.852	0.857	0.987	0.988	1	1	0.974	0.975	0.987	0.988
Bangladesh	0.796	0.778	0.881	0.875	1	1	1	1	1	1
Bhutan	0.665	0.615	1	1	1	1	0.894	0.889	1	1
Brazil	0.815	0.8	1	1	1	1	1	1	1	1
Cambodia	0.456	0.518	0.78	0.809	0.813	0.837	0.893	0.907	0.898	0.911
China	0.912	0.909	1	1	1	1	1	1	1	1
Colombia	0.895	0.903	1	1	1	1	1	1	1	1
Ethiopia	0.929	0.931	1	1	1	1	1	1	1	1
India	0.714	0.706	1	1	1	1	0.935	0.933	0.925	0.923
Indonesia	0.819	0.857	0.971	0.978	0.984	0.988	0.987	0.99	0.99	0.993
Iran	0.865	0.857	0.925	0.923	1	1	1	1	1	1
Madagascar	0.894	0.889	1	1	1	1	1	1	1	1
Malaysia	0.614	0.627	0.923	0.927	0.949	0.95	0.962	0.963	0.962	0.963
Mexico	0.92	0.919	1	1	0.92	0.919	1	1	1	1
Myanmar	0.529	0.48	0.835	0.824	0.881	0.875	0.935	0.933	1	1
Papua New Guinea	0.616	0.6	0.888	0.889	0.934	0.933	0.976	0.977	0.975	0.976
Peru	0.839	0.846	1	1	0.961	0.962	1	1	1	1
Sudan	1	1	1	1	1	1	1	1	1	1
Thailand	0.681	0.725	0.971	0.975	0.985	0.988	0.99	0.992	0.985	0.987
Vietnam	0.389	0.446	0.714	0.74	0.733	0.757	0.861	0.872	0.866	0.876
<b>Mean</b>	<b>0.76</b>	<b>0.763</b>	<b>0.944</b>	<b>0.946</b>	<b>0.958</b>	<b>0.96</b>	<b>0.97</b>	<b>0.972</b>	<b>0.979</b>	<b>0.981</b>
<b>Median</b>	<b>0.817</b>	<b>0.823</b>	<b>0.994</b>	<b>0.994</b>	<b>1</b>	<b>1</b>	<b>0.995</b>	<b>0.996</b>	<b>1</b>	<b>1</b>
<b>Q1</b>	<b>0.653</b>	<b>0.624</b>	<b>0.914</b>	<b>0.915</b>	<b>0.945</b>	<b>0.946</b>	<b>0.955</b>	<b>0.956</b>	<b>0.983</b>	<b>0.984</b>
<b>Min</b>	<b>0.389</b>	<b>0.446</b>	<b>0.714</b>	<b>0.74</b>	<b>0.733</b>	<b>0.757</b>	<b>0.861</b>	<b>0.872</b>	<b>0.866</b>	<b>0.876</b>

**A****B**

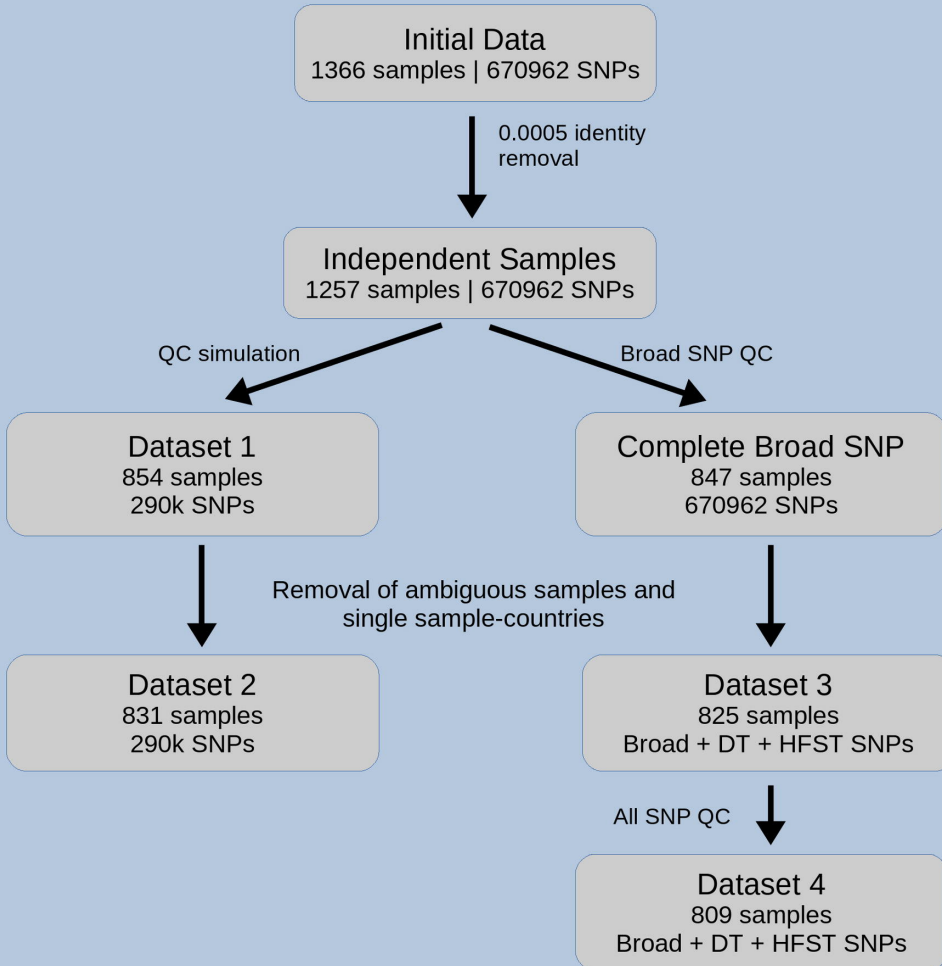


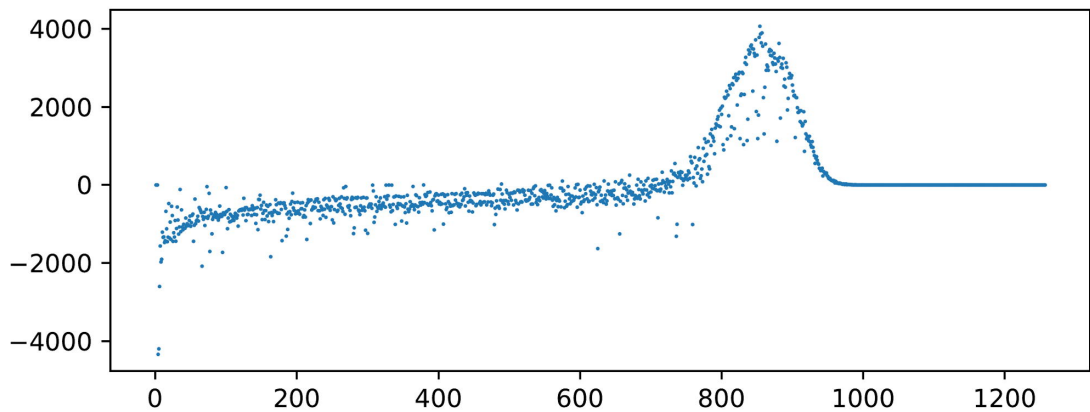
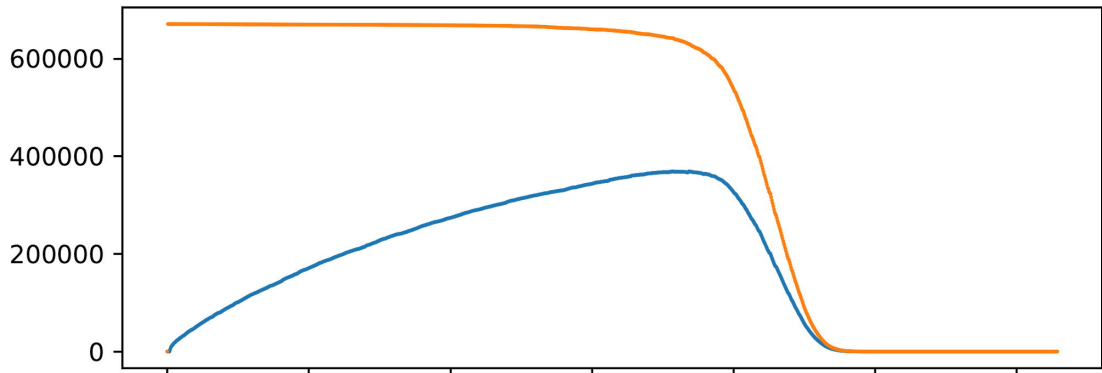






# Sample Processing





- Papua New Guinea
- Indonesia
- Malaysia
- Thailand
- Cambodia
- Vietnam
- China
- Myanmar
- Bangladesh
- Bhutan
- India
- Afghanistan
- Iran
- Ethiopia
- Sudan
- Madagascar
- Brazil
- Colombia
- Mexico
- Peru

