# A molecular barcode to inform the geographical origin and transmission dynamics of *Plasmodium vivax* malaria.
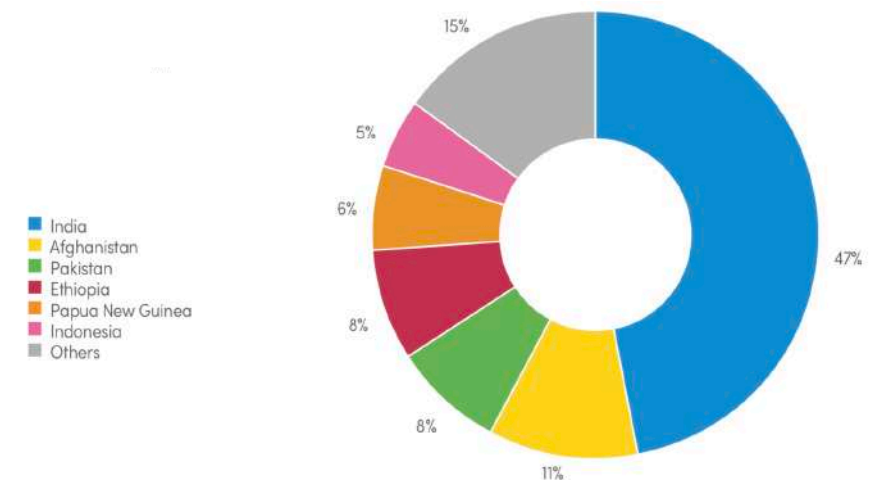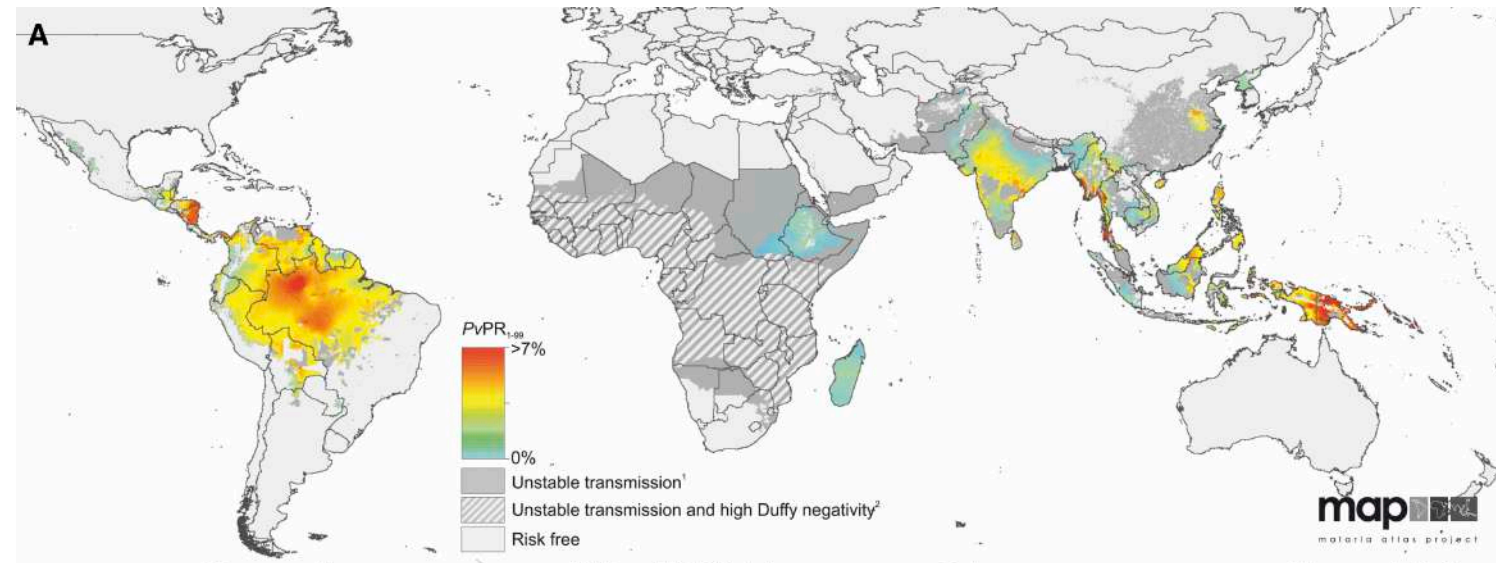
Ernest Diez Benavente
ernest.diezbenavente@lshtm.ac.uk

Ernest Diez Benavente
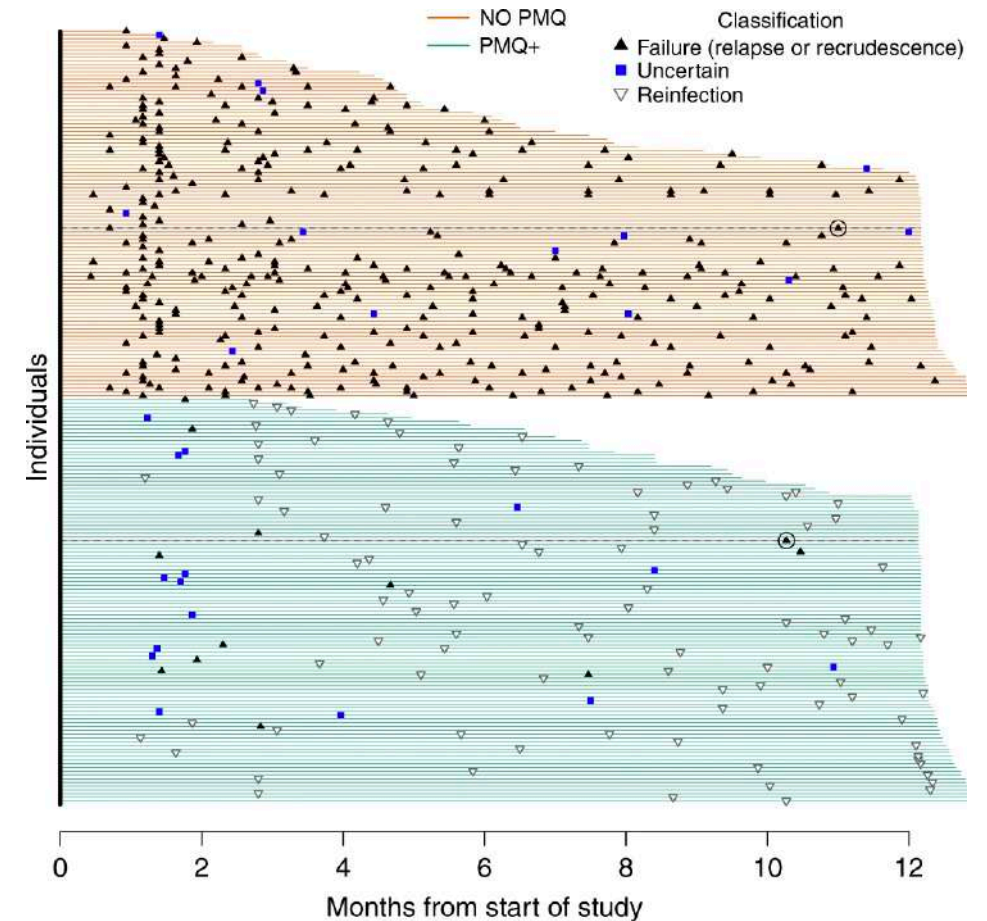ernest.diezbenavente@lshtm.ac.uk

# *Plasmodium vivax,* burden of disease

- *P. vivax* is one of the six protozoan *Plasmodium* parasites that causes **human malaria** and is transmitted by the bite of **female Anopheles mosquitoes**.

- It has the **widest geographical distribution of the six human malarias** (WHO, 2019).

- Estimated **7.5 million symptomatic cases of *P. vivax* malaria in 2018** (WHO, 2019), **53% of these are in the WHO South-East Asia Region**, with the majority being in India (47%). Most predominant parasite in the WHO Region of the Americas, representing 75% of malaria cases (WHO, 2019).

- The reasons for its distribution include amongst others:
  - ~70 species of *Anopheles* capable of transmitting the disease
  - The wide ranges of temperatures that the disease can survive on
  - Host genetics: Protection associated with Duffy Binding Receptor negativity



PvPR₁₋₉₉
>7%
0%
Unstable transmission[1]
Unstable transmission and high Duffy negativity[2]
Risk free



India 47%
Afghanistan 11%
Pakistan 8%
Ethiopia 8%
Papua New Guinea 6%
Indonesia 5%
Others 15%

# Plasmodium vivax, a challenging parasite

- Despite being **less virulent than *P. falciparum***, it can still **cause life-threatening infections** and due to its hypnozoite formation in the liver, causes **relapsing malaria**.

- **Relapsing *Plasmodium vivax*** can be treated using **primaquine** which kills the liver stages of the parasite, but it is **not recommended for G6PD deficient patients** as it can cause **severe anemia**.

- Some malaria endemic regions have reported an **increase in the proportion of *P. vivax* cases during effective control of *P. falciparum* malaria**, highlighting the resilience of this parasite.

- *Plasmodium vivax* can't be maintained in in-vitro culture.

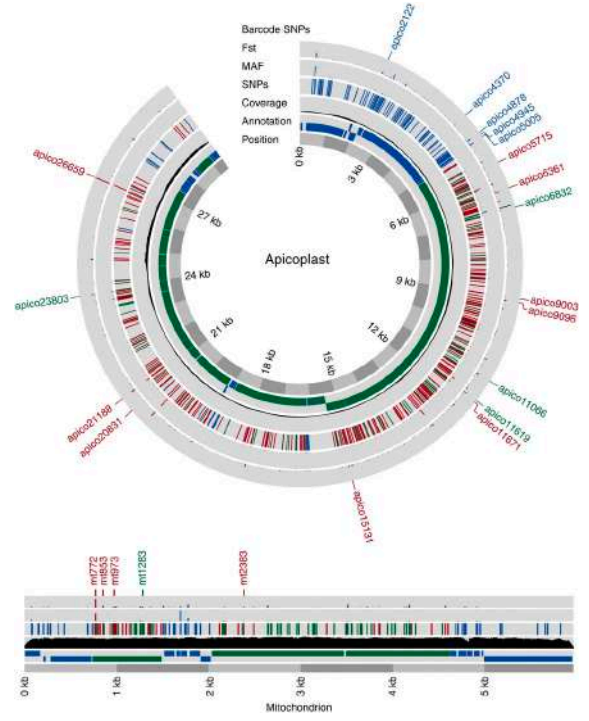- Reports of **chloroquine drug resistance (first line drug)** in parts of SEA.



Taylor et al. 2019

Studying transmission using SNP barcode based on nuclear genome

Geographic inference using SNP barcode based on organellar genomes

Daniels et al. 2009

Preston et al. 2014

Daniels et al. 2015

Using phone tracking data and 101-SNP barcode based on nuclear genome to map imported malaria in Bangladesh



Chang et al. 2019
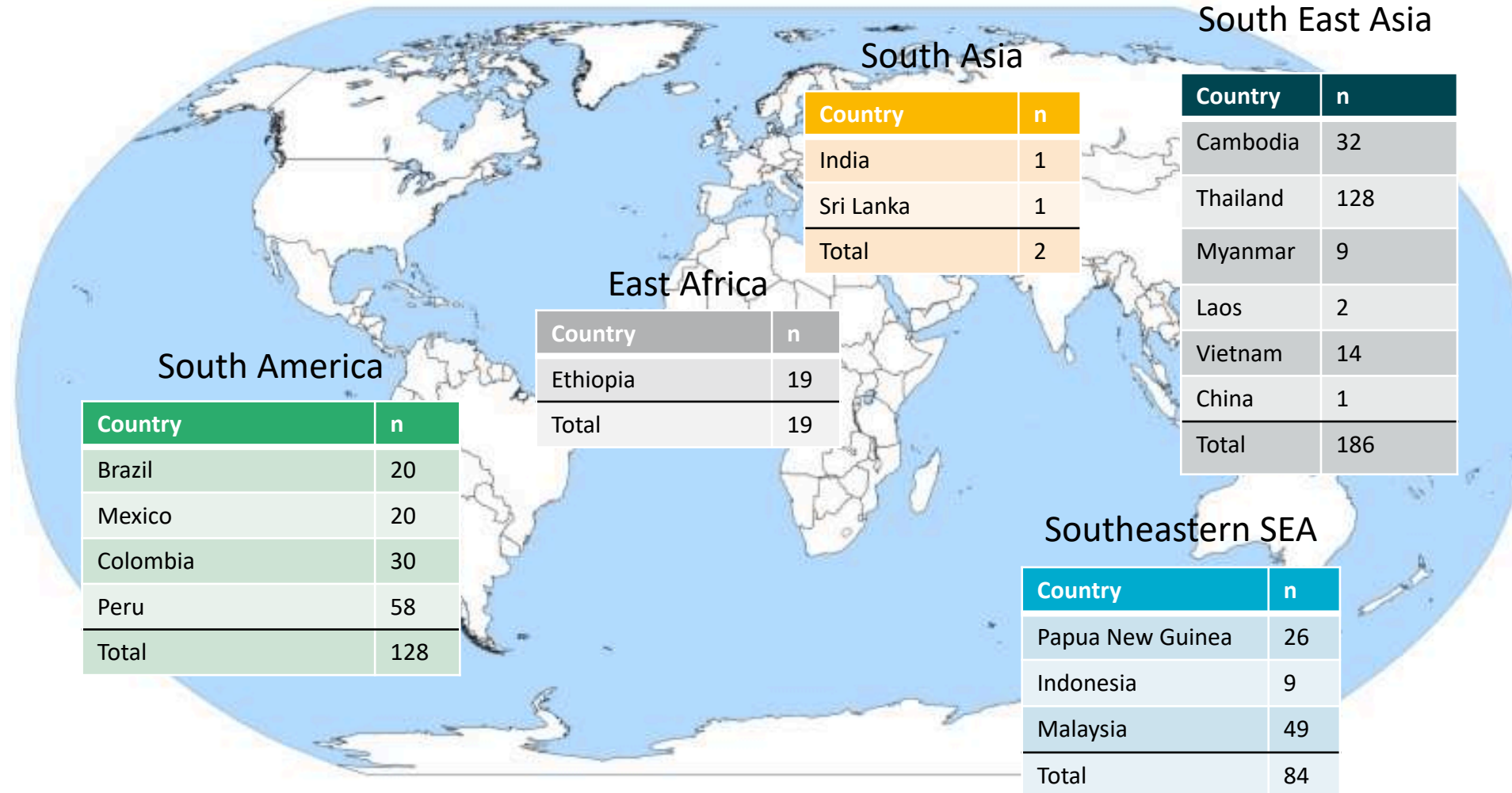
Chang et al. 2019

42-SNP barcode for *P. vivax* population genetics based on data from 13 isolates



Baniecki et al. 2015

Can we create a **SNP barcode** that can **both** be used to identify **geographic origin** and help characterize **transmission dynamics** of *Plasmodium vivax* infections?
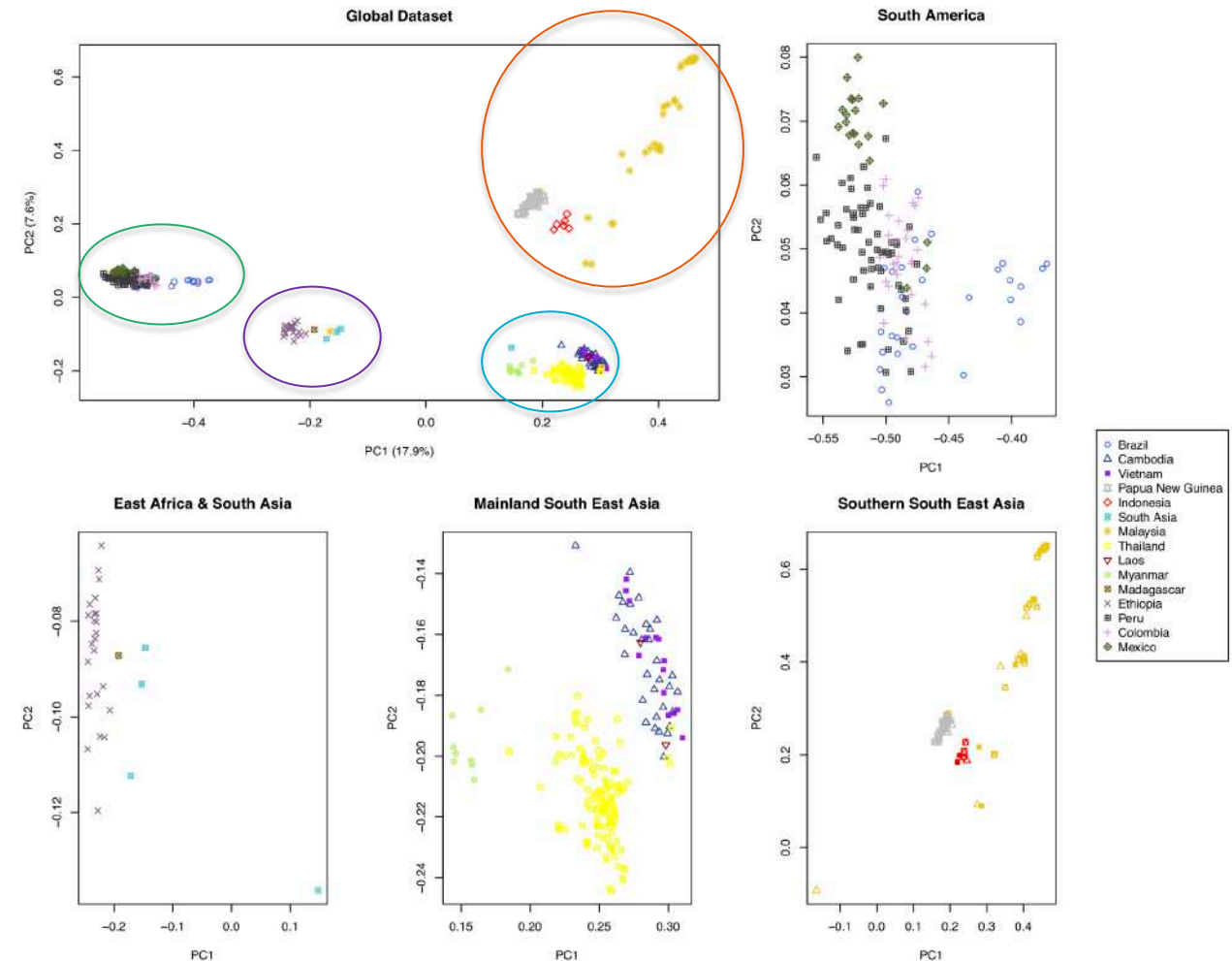
# A SNP barcode to inform geographic origin and transmission dynamics in *Plasmodium vivax*.

- Starting with **837 isolates WGS** from the different endemic areas of *P. vivax*.

- Identified 1,522,046 variants using a bioinformatics pipeline using ***trimmomatic, bwa* and *samtools***.

- After quality filtering using different thresholds for **coverage, missing data, mixed calls and low quality base calls** the final high quality data consisted of **433 isolates and 720,340 SNPs**.

### South Asia

| Country | n |
|---|---|
| India | 1 |
| Sri Lanka | 1 |
| Total | 2 |

### South East Asia

| Country | n |
|---|---|
| Cambodia | 32 |
| Thailand | 128 |
| Myanmar | 9 |
| Laos | 2 |
| Vietnam | 14 |
| China | 1 |
| Total | 186 |

### East Africa

| Country | n |
|---|---|
| Ethiopia | 19 |
| Total | 19 |

### South America

| Country | n |
|---|---|
| Brazil | 20 |
| Mexico | 20 |
| Colombia | 30 |
| Peru | 58 |
| Total | 128 |

### Southeastern SEA

| Country | n |
|---|---|
| Papua New Guinea | 26 |
| Indonesia | 9 |
| Malaysia | 49 |
| Total | 84 |

- PCA generated with 720K SNPs **shows good clustering of isolates by region** and to **relatively good separation of isolates by country of origin**.

- The **whole genome pairwise genetic distance between isolates** obtained using the Manhattan method (number of **SNPs differences** in pairwise comparisons) was used as the **gold standard for assessing the performance in the measurement of relatedness of the different subsets of SNPs**.

- **NOTE: The genetic distance used to generate the PCA plots has been divided by the sum of the MAF of all the SNPs considered in order to make the distances comparable across subset of SNPs.**



Diez Benavente et al. PLOS Genetics 2020

SNPs partitioned in **three subsets of equal number of SNPs sorted by MAF**.



SNP Subset 1 (n=240113)
min MAF= 0, max MAF= 0.002

SNP Subset 2 (n=240113)
min MAF= 0.002, max MAF= 0.007

SNP Subset 3 (n=240114)
min MAF= 0.007, max MAF= 0.499

Country
- Brazil
- Cambodia
- Vietnam
- Papua New Guinea
- Indonesia
- South Asia
- Malaysia
- Thailand
- Laos
- Myanmar
- Madagascar
- Ethiopia
- Peru
- Colombia
- Mexico

$r^2 = 0.25$

$r^2 = 0.27$

$r^2 = 0.94$

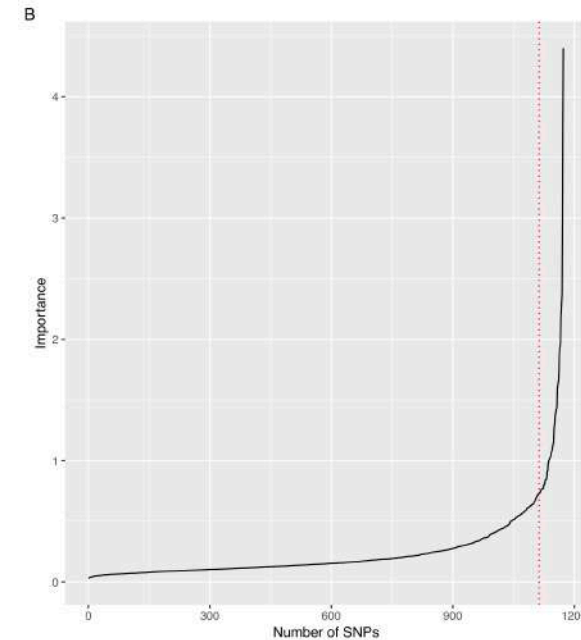Diez Benavente et al. PLOS Genetics 2020

- After MAF selection **16,110 SNPs** were used for further selection.

- We are interested in obtaining independent markers in order to avoid oversampling SNPs that might be high MAF due to sampling differences in our population of isolates.

- *Tagster*, a software originally designed for SNP selection for genotyping chip design was used.

- Tagster identifies SNPs that "tag" other SNPs in vicinity, so we used this with relatively high threshold for distance between SNPs (500 Kb) and la LD cutoff of 0.7.

- After selection we identified **1,173 independent SNPs** that were able to capture the variability of 40% of the SNPs in our dataset.

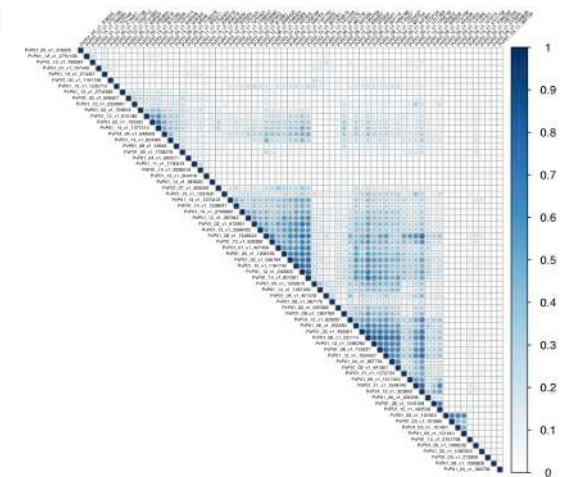Diez Benavente et al. PLOS Genetics 2020

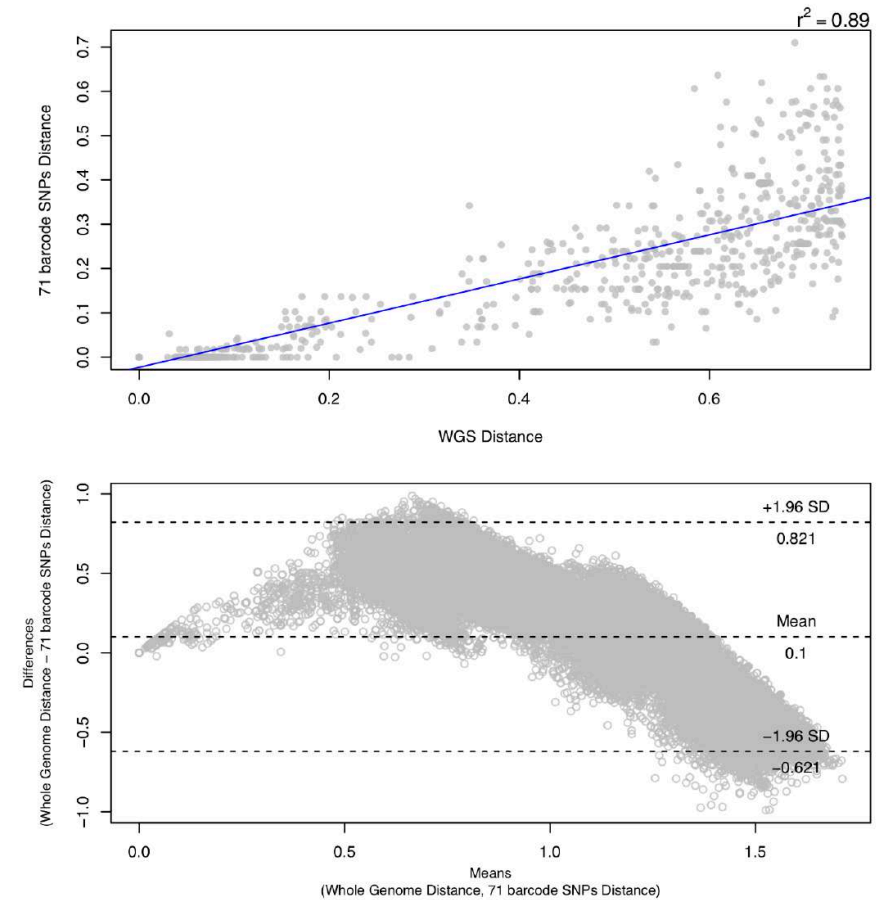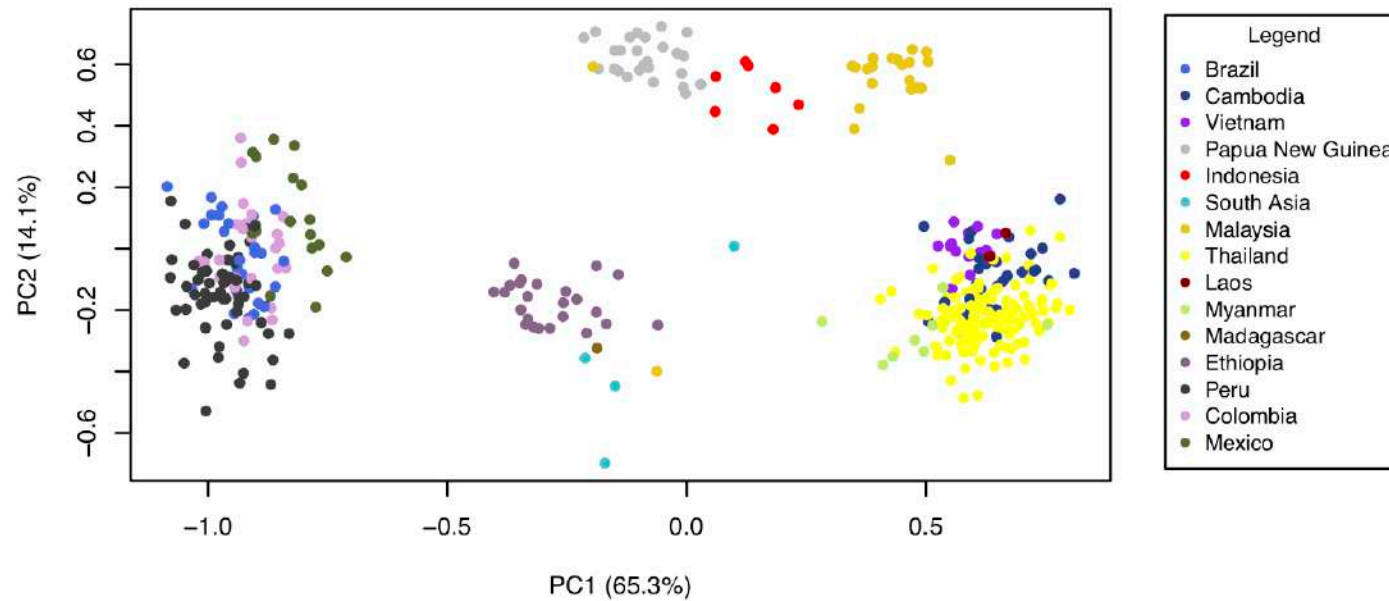# Using Random Forests to predict geographic origin and identify most relevant SNPs for geographic origin.

- Random forest classification using country as the outcome to be predicted, we partitioned the dataset randomly into 80% training (n=346) and 20% for validation (n=87), obtaining a out-of-bag error rate of 17.1%.

- We identify the **60 SNPs with the highest importance** for classification and combine this with 11 SNPs with $F_{ST}$ >0.7 (highly differentiating) for the populations with the highest misclassification rates (Thailand/Myanamar).

- We retrain the model using the training set but in this case using only the 71 SNP barcode positions, obtaining an overall accuracy **of 91.4% when predicting in the validation set**.
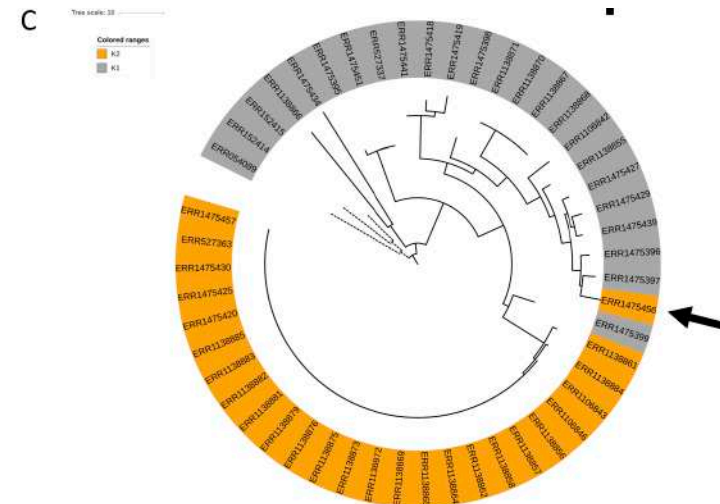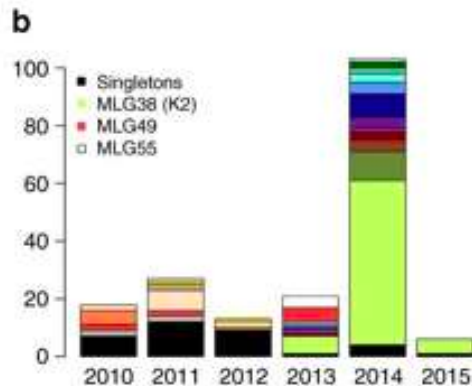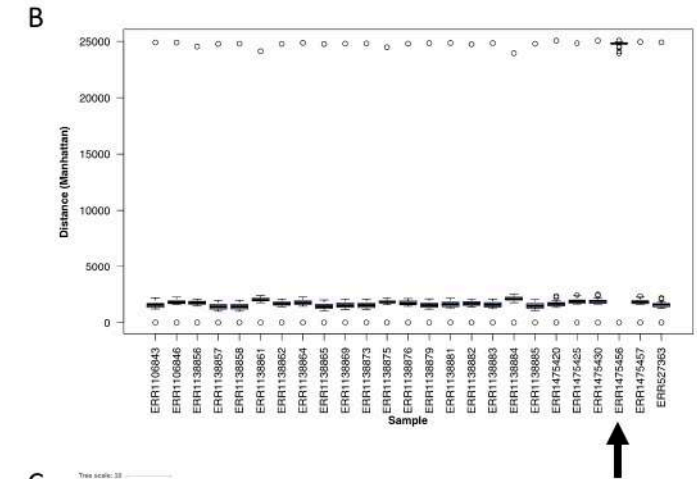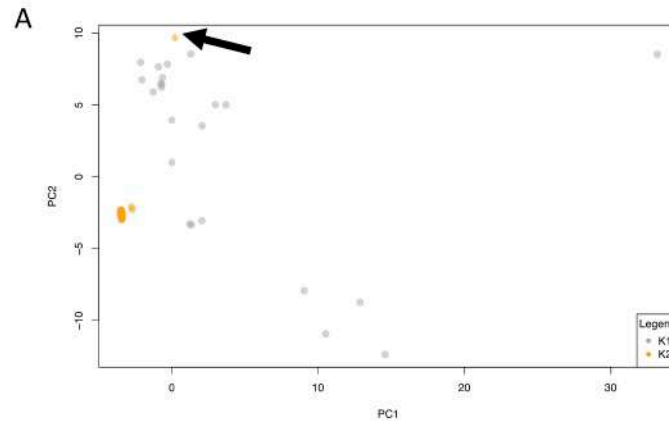


Mean 0.15, IQR = 0.02 - 0.24

Diez Benavente et al. PLOS Genetics 2020

# How does the 71 SNP barcode perform?



Diez Benavente et al. PLOS Genetics 2020

Auburn et al. Nat Comms 2018

Diez Benavente et al. PLOS Genetics 2020

# Use case example for the 71 SNPs barcode, origin classification of isolates from returning travelers to the UK (MRL).



| Country | n |
|---------|---|
| Afghanistan | 31 |
| Bangladesh | 2 |
| India | 37 |
| Pakistan | 35 |
| Eritrea | 17 |
| Ethiopia | 7 |
| Sudan | 10 |
| Uganda | 5 |
| The Philippines | 1 |
| Guyana | 3 |

Diez Benavente et al. PLOS Genetics 2020

# Conclusions

- SNPs barcodes can be designed to be used for geographic origin prediction and to study transmission dynamics of *P. vivax* simultaneously.

- Whole genome sequencing can be used to inform barcode design for *Plasmodium sp*.

- Our in-silico 71-SNP barcode for *P. vivax* showed strong regional classification and promising country level geographical origin prediction.

- Some caveats:
  - We have not assessed the impact of mixed infections, this work assumes only mono infections are considered.
  - Random Forests might not be the best method to perform variable selection, other models could be used.
  - We used genetic distance as our measure of relatedness, but other measures have shown stronger power to understand transmission dynamics (i.e. identity by descent) although it is estimated that at least 200 SNPs would be needed to assess this.
  - Our training dataset is missing some important regions (i.e. South Asia).

# Acknowledgements