
















RESEARCH ARTICLE

**REVISED** **An open dataset of *Plasmodium falciparum* genome variation in 7,000 worldwide samples [version 2; peer review: 2 approved]**

MalariaGEN, Ambroise Ahouidi<sup>1</sup>, Mozam Ali<sup>2</sup>, Jacob Almagro-Garcia<sup>2,3</sup>, Alfred Amambua-Ngwa<sup>2,4</sup>, Chanaki Amaratunga<sup>5</sup>, Roberto Amato<sup>2,3</sup>, Lucas Amenga-Etego <sup>6,7</sup>, Ben Andagalu<sup>8</sup>, Tim J. C. Anderson <sup>9</sup>, Voahangy Andrianarajaka<sup>10</sup>, Tobias Apinjoh<sup>11</sup>, Cristina Ariani<sup>2</sup>, Elizabeth A. Ashley <sup>12</sup>, Sarah Auburn<sup>13,14</sup>, Gordon A. Awandare<sup>7,15</sup>, Hampate Ba <sup>16</sup>, Vito Baraka <sup>17,18</sup>, Alyssa E. Barry<sup>19-21</sup>, Philip Bejon <sup>22</sup>, Gwladys I. Bertin <sup>23</sup>, Maciej F. Boni<sup>14,24</sup>, Steffen Borrmann<sup>25</sup>, Teun Bousema <sup>26,27</sup>, Orale Branch<sup>28</sup>, Peter C. Bull<sup>22,29</sup>, George B. J. Busby<sup>3</sup>, Thanat Chookajorn <sup>30</sup>, Kesinee Chotivanich<sup>30</sup>, Antoine Claessens <sup>4,31</sup>, David Conway <sup>26</sup>, Alister Craig <sup>32,33</sup>, Umberto D'Alessandro <sup>4</sup>, Souleymane Dama<sup>34</sup>, Nicholas PJ Day <sup>12</sup>, Brigitte Denis<sup>33</sup>, Mahamadou Diakite <sup>34</sup>, Abdoulaye Djimdé <sup>34</sup>, Christiane Dolecek<sup>14</sup>, Arjen M Dondorp <sup>12</sup>, Chris Drakeley <sup>26</sup>, Eleanor Drury<sup>2</sup>, Patrick Duffy<sup>5</sup>, Diego F. Echeverry<sup>35,36</sup>, Thomas G. Egwang<sup>37</sup>, Berhanu Erko<sup>38</sup>, Rick M. Fairhurst<sup>39</sup>, Abdul Faiz <sup>40</sup>, Caterina A. Fanello <sup>12</sup>, Mark M. Fukuda<sup>41</sup>, Dionicia Gamboa <sup>42</sup>, Anita Ghansah<sup>43</sup>, Lemu Golassa <sup>38</sup>, Sonia Goncalves<sup>2</sup>, William L. Hamilton <sup>2,44</sup>, G. L. Abby Harrison<sup>21</sup>, Lee Hart <sup>3</sup>, Christa Henrichs<sup>3</sup>, Tran Tinh Hien <sup>24,45</sup>, Catherine A. Hill<sup>46</sup>, Abraham Hodgson<sup>47</sup>, Christina Hubbart <sup>48</sup>, Mallika Imwong<sup>30</sup>, Deus S. Ishengoma <sup>17,49</sup>, Scott A. Jackson <sup>50</sup>, Chris G. Jacob<sup>2</sup>, Ben Jeffery<sup>3</sup>, Anna E. Jeffreys <sup>48</sup>, Kimberly J. Johnson <sup>3</sup>, Dushyanth Jyothi <sup>2</sup>, Claire Kamaliddin <sup>23</sup>, Edwin Kamau <sup>51</sup>, Mihir Kekre<sup>2</sup>, Krzysztof Kluczynski<sup>3</sup>, Theerarat Kochakarn<sup>2,30</sup>, Abibatou Konaté<sup>52</sup>, Dominic P. Kwiatkowski <sup>2,3,48</sup>, Myat Phone Kyaw<sup>53,54</sup>, Pharath Lim<sup>5,55</sup>, Chanthap Lon<sup>41</sup>, Kovana M. Loua <sup>56</sup>, Oumou Maïga-Ascofaré<sup>34,57,58</sup>, Cinzia Malangone <sup>2</sup>, Magnus Manske<sup>2</sup>, Jutta Marfurt<sup>13</sup>, Kevin Marsh <sup>14,59</sup>, Mayfong Mayxay <sup>60,61</sup>, Alistair Miles<sup>2,3</sup>, Olivo Miotto <sup>2,3,12</sup>, Victor Mobegi <sup>62</sup>, Olugbenga A. Mokuolu<sup>63</sup>, Jacqui Montgomery<sup>64</sup>, Ivo Mueller<sup>21,65</sup>, Paul N. Newton<sup>66</sup>, Thuy Nguyen<sup>2</sup>, Thuy-Nhien Nguyen<sup>24</sup>, Harald Noedl<sup>67</sup>, François Nosten <sup>14,68</sup>, Rintis Noviyanti<sup>69</sup>,

Alexis Nzila<sup>70</sup>, Lynette I. Ochola-Oyier<sup>22</sup>, Harold Ocholla<sup>71,72</sup>, Abraham Oduro <sup>6</sup>, Irene Omedo <sup>22</sup>, Marie A. Onyamboko<sup>73</sup>, Jean-Bosco Ouedraogo <sup>74</sup>, Kolapo Oyebola <sup>75,76</sup>, Richard D. Pearson <sup>2,3</sup>, Norbert Peshu<sup>22</sup>, Aung Pyae Phyo <sup>12,68</sup>, Chris V. Plowe<sup>77</sup>, Ric N. Price <sup>12,13,45</sup>, Sasithon Pukrittayakamee<sup>30</sup>, Milijaona Randrianarivelosia<sup>78,79</sup>, Julian C. Rayner <sup>2</sup>, Pascal Ringwald<sup>80</sup>, Kirk A. Rockett <sup>2,48</sup>, Katherine Rowlands<sup>48</sup>, Lastenia Ruiz<sup>81</sup>, David Saunders<sup>41</sup>, Alex Shayo <sup>82</sup>, Peter Siba<sup>83</sup>, Victoria J. Simpson<sup>3</sup>, Jim Stalker<sup>2</sup>, Xin-zhuan Su <sup>5</sup>, Colin Sutherland<sup>26</sup>, Shannon Takala-Harrison<sup>84</sup>, Livingstone Tavul<sup>83</sup>, Vandana Thathy<sup>22,85</sup>, Antoinette Tshetu<sup>86</sup>, Federica Verra<sup>87</sup>, Joseph Vinetz<sup>42,88</sup>, Thomas E. Wellems <sup>5</sup>, Jason Wendler<sup>48</sup>, Nicholas J. White<sup>12</sup>, Ian Wright <sup>3</sup>, William Yavo<sup>52,89</sup>, Htut Ye<sup>90</sup>

<sup>1</sup>Hopital Le Dantec, Universite Cheikh Anta Diop, Dakar, Senegal

<sup>2</sup>Wellcome Sanger Institute, Hinxton, UK

<sup>3</sup>MRC Centre for Genomics and Global Health, Big Data Institute, University of Oxford, Oxford, UK

<sup>4</sup>Medical Research Council Unit The Gambia, at the London School of Hygiene and Tropical Medicine, Banjul, The Gambia

<sup>5</sup>National Institute of Allergy and Infectious Diseases (NIAID), NIH, Bethesda, USA

<sup>6</sup>Navrongo Health Research Centre, Ghana Health Service, Navrongo, Ghana

<sup>7</sup>West African Centre for Cell Biology of Infectious Pathogens (WACCBIP), University of Ghana, Accra, Ghana

<sup>8</sup>United States Army Medical Research Directorate-Africa, Kenya Medical Research Institute/Walter Reed Project, Kisumu, Kenya

<sup>9</sup>Texas Biomedical Research Institute, San Antonio, USA

<sup>10</sup>Université d'Antananarivo, Antananarivo, Madagascar

<sup>11</sup>University of Buea, Buea, Cameroon

<sup>12</sup>Mahidol-Oxford Tropical Medicine Research Unit (MORU), Bangkok, Thailand

<sup>13</sup>Menzies School of Health Research, Darwin, Australia

<sup>14</sup>Nuffield Department of Medicine, University of Oxford, Oxford, UK

<sup>15</sup>University of Ghana, Legon, Ghana

<sup>16</sup>Institut National de Recherche en Santé Publique, Nouakchott, Mauritania

<sup>17</sup>National Institute for Medical Research (NIMR), Dar es Salaam, Tanzania

<sup>18</sup>Department of Epidemiology, International Health Unit, University of Antwerp, Antwerp, Belgium

<sup>19</sup>Deakin University, Geelong, Australia

<sup>20</sup>Burnet Institute, Melbourne, Australia

<sup>21</sup>Walter and Eliza Hall Institute, Melbourne, Australia

<sup>22</sup>KEMRI Wellcome Trust Research Programme, Kilifi, Kenya

<sup>23</sup>Institute of Research for Development (IRD), Paris, France

<sup>24</sup>Oxford University Clinical Research Unit (OUCRU), Ho Chi Minh City, Vietnam

<sup>25</sup>Institute for Tropical Medicine, University of Tübingen, Tübingen, Germany

<sup>26</sup>London School of Hygiene and Tropical Medicine, London, UK

<sup>27</sup>Radboud University Medical Center, Nijmegen, The Netherlands

<sup>28</sup>NYU School of Medicine Langone Medical Center, New York, USA

<sup>29</sup>Department of Pathology, University of Cambridge, Cambridge, UK

<sup>30</sup>Mahidol University, Bangkok, Thailand

<sup>31</sup>LPHI, MIVEGEC, INSERM, CNRS, IRD, University of Montpellier, Montpellier, France

<sup>32</sup>Liverpool School of Tropical Medicine, Liverpool, UK

<sup>33</sup>Malawi-Liverpool-Wellcome Trust Clinical Research, Blantyre, Malawi

<sup>34</sup>Malaria Research and Training Centre, University of Science, Techniques and Technologies of Bamako, Bamako, Mali

<sup>35</sup>Centro Internacional de Entrenamiento e Investigaciones Médicas - CIDEIM, Cali, Colombia

<sup>36</sup>Universidad Icesi, Cali, Colombia

- <sup>37</sup>Biotech Laboratories, Kampala, Uganda
- <sup>38</sup>Aklilu Lemma Institute of Pathobiology, Addis Ababa University, Addis Ababa, Ethiopia
- <sup>39</sup>National Institutes of Health (NIH), Bethesda, USA
- <sup>40</sup>Dev Care Foundation, Dhaka, Bangladesh
- <sup>41</sup>Department of Immunology and Medicine, US Army Medical Component, Armed Forces Research Institute of Medical Sciences (USAMC-AFRIMS), Bangkok, Thailand
- <sup>42</sup>Laboratorio ICMR-Amazonia, Laboratorios de Investigacion y Desarrollo, Facultad de Ciencias y Filosofia, Universidad Peruana Cayetano Heredia, Lima, Peru
- <sup>43</sup>Nogouchi Memorial Institute for Medical Research, Legon-Accra, Ghana
- <sup>44</sup>Cambridge University Hospitals NHS Foundation Trust, Cambridge, UK
- <sup>45</sup>Centre for Tropical Medicine and Global Health, University of Oxford, Oxford, UK
- <sup>46</sup>Department of Entomology, Purdue University, West Lafayette, USA
- <sup>47</sup>Ghana Health Service, Ministry of Health, Accra, Ghana
- <sup>48</sup>Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK
- <sup>49</sup>East African Consortium for Clinical Research (EACCR), Dar es Salaam, Tanzania
- <sup>50</sup>Center for Applied Genetic Technologies, University of Georgia, Athens, GA, USA
- <sup>51</sup>Walter Reed Army Institute of Research, U.S. Military HIV Research Program, Silver Spring, MD, USA
- <sup>52</sup>University Félix Houphouët-Boigny, Abidjan, Cote d'Ivoire
- <sup>53</sup>The Myanmar Oxford Clinical Research Unit, University of Oxford, Yangon, Myanmar
- <sup>54</sup>University of Public Health, Yangon, Myanmar
- <sup>55</sup>Medical Care Development International, Maryland, USA
- <sup>56</sup>Institut Nationale de Santé Publique, Conakry, Guinea
- <sup>57</sup>Bernhard Nocht Institute for Tropical Medicine, Hamburg, Germany
- <sup>58</sup>Research in Tropical Medicine, Kwame Nkrumah University of Sciences and Technology, Kumasi, Ghana
- <sup>59</sup>African Academy of Sciences, Nairobi, Kenya
- <sup>60</sup>Lao-Oxford-Mahosot Hospital-Wellcome Trust Research Unit (LOMWRU), Vientiane, Lao People's Democratic Republic
- <sup>61</sup>Institute of Research and Education Development (IRED), University of Health Sciences, Ministry of Health, Vientiane, Lao People's Democratic Republic
- <sup>62</sup>School of Medicine, University of Nairobi, Nairobi, Kenya
- <sup>63</sup>Department of Paediatrics and Child Health, University of Ilorin, Ilorin, Nigeria
- <sup>64</sup>Institute of Vector-Borne Disease, Monash University, Clayton, Victoria, 3800, Australia
- <sup>65</sup>Barcelona Centre for International Health Research, Barcelona, Spain
- <sup>66</sup>Wellcome Trust-Mahosot Hospital-Oxford Tropical Medicine Research Collaboration, Vientiane, Lao People's Democratic Republic
- <sup>67</sup>MARIB - Malaria Research Initiative Bandarban, Bandarban, Bangladesh
- <sup>68</sup>Shoklo Malaria Research Unit, Bangkok, Thailand
- <sup>69</sup>Eijkman Institute for Molecular Biology, Jakarta, Indonesia
- <sup>70</sup>King Fahid University of Petroleum and Minerals (KFUMP), Dharhran, Saudi Arabia
- <sup>71</sup>KEMRI - Centres for Disease Control and Prevention (CDC) Research Program, Kisumu, Kenya
- <sup>72</sup>Centre for Bioinformatics and Biotechnology, University of Nairobi, Nairobi, Kenya
- <sup>73</sup>Kinshasa School of Public Health, University of Kinshasa, Kinshasa, Congo, Democratic Republic
- <sup>74</sup>Institut de Recherche en Sciences de la Santé, Ouagadougou, Burkina Faso
- <sup>75</sup>Nigerian Institute of Medical Research, Lagos, Nigeria
- <sup>76</sup>Parasitology and Bioinformatics Unit, Faculty of Science, University of Lagos, Lagos, Nigeria
- <sup>77</sup>School of Medicine, University of Maryland, Baltimore, MD, USA
- <sup>78</sup>Institut Pasteur de Madagascar, Antananarivo, Madagascar
- <sup>79</sup>Universités d'Antananarivo et de Mahajanga, Antananarivo, Madagascar
- <sup>80</sup>World Health Organization (WHO), Geneva, Switzerland
- <sup>81</sup>Universidad Nacional de la Amazonia Peruana, Iquitos, Peru
- <sup>82</sup>Nelson Mandela Institute of Science and Technology, Arusha, Tanzania
- <sup>83</sup>Papua New Guinea Institute of Medical Research, Goroka, Papua New Guinea
- <sup>84</sup>Center for Vaccine Development and Global Health, University of Maryland, School of Medicine, Baltimore, MD, USA
- <sup>85</sup>Department of Microbiology and Immunology, Columbia University Irving Medical Center, New York, New York, USA
- <sup>86</sup>University of Kinshasa, Kinshasa, Congo, Democratic Republic
- <sup>87</sup>Sapienza University of Rome, Rome, Italy
- <sup>88</sup>Yale School of Medicine, New Haven, CT, USA
- <sup>89</sup>Malaria Research and Control Center of the National Institute of Public Health, Abidjan, Cote d'Ivoire
- 90

Department of Medical Research, Yangon, Myanmar

**v2** First published: 24 Feb 2021, 6:42  
<https://doi.org/10.12688/wellcomeopenres.16168.1>  
 Latest published: 13 Jul 2021, 6:42  
<https://doi.org/10.12688/wellcomeopenres.16168.2>

## Abstract

MalariaGEN is a data-sharing network that enables groups around the world to work together on the genomic epidemiology of malaria. Here we describe a new release of curated genome variation data on 7,000 *Plasmodium falciparum* samples from MalariaGEN partner studies in 28 malaria-endemic countries. High-quality genotype calls on 3 million single nucleotide polymorphisms (SNPs) and short indels were produced using a standardised analysis pipeline. Copy number variants associated with drug resistance and structural variants that cause failure of rapid diagnostic tests were also analysed. Almost all samples showed genetic evidence of resistance to at least one antimalarial drug, and some samples from Southeast Asia carried markers of resistance to six commonly-used drugs. Genes expressed during the mosquito stage of the parasite life-cycle are prominent among loci that show strong geographic differentiation. By continuing to enlarge this open data resource we aim to facilitate research into the evolutionary processes affecting malaria control and to accelerate development of the surveillance toolkit required for malaria elimination.

## Keywords

malaria, plasmodium falciparum, genomics, genomic epidemiology, evolution, data resource, population genetics, drug resistance, rapid diagnostic test failure

## Open Peer Review

Reviewer Status  

Invited Reviewers

1

2

### version 2

(revision)

13 Jul 2021

### version 1




24 Feb 2021



report



report

1. **Maria Isabel Veiga** , University of Minho, Braga, Portugal  
**Nuno S. Osório** , University of Minho, Braga, Portugal
2. **Didier Menard** , Institut Pasteur, Paris, France

Any reports and responses or comments on the article can be found at the end of the article.

**Corresponding author:** MalariaGEN ([support@malariagen.net](mailto:support@malariagen.net))

**Author roles:** **Ahoudi A:** Investigation, Resources, Writing – Review & Editing; **Ali M:** Investigation, Writing – Review & Editing; **Almagro-Garcia J:** Formal Analysis, Investigation, Writing – Review & Editing; **Amambua-Ngwa A:** Investigation, Resources, Writing – Review & Editing; **Amaratunga C:** Investigation, Resources, Writing – Review & Editing; **Amato R:** Conceptualization, Data Curation, Formal Analysis, Investigation, Project Administration, Software, Supervision, Writing – Original Draft Preparation, Writing – Review & Editing; **Amenga-Etego L:** Investigation, Resources, Writing – Review & Editing; **Andagalu B:** Investigation, Resources, Writing – Review & Editing; **Anderson TJC:** Investigation, Resources, Writing – Review & Editing; **Andrianarajaka V:** Investigation, Resources, Writing – Review & Editing; **Apinjoh T:** Investigation, Resources, Writing – Review & Editing; **Ariani C:** Investigation, Writing – Review & Editing; **Ashley EA:** Investigation, Resources, Writing – Review & Editing; **Auburn S:** Investigation, Resources, Writing – Review & Editing; **Awandare GA:** Investigation, Resources, Writing – Review & Editing; **Ba H:** Investigation, Resources, Writing – Review & Editing; **Baraka V:** Investigation, Resources, Writing – Review & Editing; **Barry AE:** Investigation, Resources, Writing – Review & Editing; **Bejon P:** Investigation, Resources, Writing – Review & Editing; **Bertin GI:** Investigation, Resources, Writing – Review & Editing; **Boni MF:** Investigation, Resources, Writing – Review & Editing; **Borrmann S:** Investigation, Resources, Writing – Review & Editing; **Bousema T:** Investigation, Resources, Writing – Review & Editing; **Branch O:** Investigation, Resources, Writing – Review & Editing; **Bull PC:** Investigation, Resources, Writing – Review & Editing; **Busby GBJ:** Investigation, Software, Writing – Review & Editing; **Chookajorn T:** Formal Analysis, Investigation, Writing – Review & Editing; **Chotivanich K:** Investigation, Resources, Writing – Review & Editing; **Claessens A:** Investigation, Resources, Writing – Review & Editing; **Conway D:** Investigation, Resources, Writing – Review & Editing; **Craig A:** Investigation, Resources, Writing – Review & Editing; **D'Alessandro U:** Investigation, Resources, Writing – Review & Editing; **Dama S:** Investigation, Resources, Writing – Review & Editing; **Day NP:** Investigation, Resources, Writing – Review & Editing; **Denis B:** Investigation, Resources, Writing – Review & Editing; **Diakite M:** Investigation, Resources, Writing – Review & Editing; **Djimdé A:** Investigation, Resources, Writing – Review & Editing; **Dolecek C:** Investigation, Resources, Writing – Review & Editing; **Dondorp AM:** Investigation, Resources, Writing – Review & Editing; **Drakeley C:** Investigation, Resources, Writing – Review & Editing; **Drury E:**

Investigation, Writing – Review & Editing; **Duffy P**: Investigation, Resources, Writing – Review & Editing; **Echeverry DF**: Investigation, Resources, Writing – Review & Editing; **Egwang TG**: Investigation, Resources, Writing – Review & Editing; **Erko B**: Investigation, Resources, Writing – Review & Editing; **Fairhurst RM**: Investigation, Resources, Writing – Review & Editing; **Faiz A**: Investigation, Resources, Writing – Review & Editing; **Fanello CA**: Investigation, Resources, Writing – Review & Editing; **Fukuda MM**: Investigation, Resources, Writing – Review & Editing; **Gamboia D**: Investigation, Resources, Writing – Review & Editing; **Ghansah A**: Investigation, Resources, Writing – Review & Editing; **Golassa L**: Investigation, Resources, Writing – Review & Editing; **Goncalves S**: Investigation, Project Administration, Writing – Review & Editing; **Hamilton WL**: Formal Analysis, Investigation, Writing – Original Draft Preparation, Writing – Review & Editing; **Harrison GLA**: Investigation, Resources, Writing – Review & Editing; **Hart L**: Investigation, Software, Writing – Review & Editing; **Henrichs C**: Investigation, Project Administration, Writing – Review & Editing; **Hien TT**: Investigation, Resources, Writing – Review & Editing; **Hill CA**: Investigation, Resources, Writing – Review & Editing; **Hodgson A**: Investigation, Resources, Writing – Review & Editing; **Hubbart C**: Investigation, Writing – Review & Editing; **Imwong M**: Investigation, Resources, Writing – Review & Editing; **Ishengoma DS**: Investigation, Resources, Writing – Review & Editing; **Jackson SA**: Investigation, Resources, Writing – Review & Editing; **Jacob CG**: Investigation, Writing – Review & Editing; **Jeffery B**: Investigation, Software, Writing – Review & Editing; **Jeffreys AE**: Investigation, Writing – Review & Editing; **Johnson KJ**: Investigation, Project Administration, Writing – Review & Editing; **Jyothi D**: Data Curation, Investigation, Software, Writing – Review & Editing; **Kamaliddin C**: Investigation, Resources, Writing – Review & Editing; **Kamau E**: Investigation, Resources, Writing – Review & Editing; **Kekre M**: Investigation, Writing – Review & Editing; **Kluczynski K**: Investigation, Software, Writing – Review & Editing; **Kochakarn T**: Formal Analysis, Investigation, Writing – Review & Editing; **Konaté A**: Investigation, Resources, Writing – Review & Editing; **Kwiatkowski DP**: Conceptualization, Formal Analysis, Funding Acquisition, Investigation, Project Administration, Supervision, Writing – Original Draft Preparation, Writing – Review & Editing; **Kyaw MP**: Investigation, Resources, Writing – Review & Editing; **Lim P**: Investigation, Resources, Writing – Review & Editing; **Lon C**: Investigation, Resources, Writing – Review & Editing; **Loua KM**: Investigation, Resources, Writing – Review & Editing; **Maiga-Ascofaré O**: Investigation, Resources, Writing – Review & Editing; **Malangone C**: Data Curation, Investigation, Software, Writing – Review & Editing; **Manske M**: Investigation, Software, Writing – Review & Editing; **Marfurt J**: Investigation, Resources, Writing – Review & Editing; **Marsh K**: Investigation, Resources, Writing – Review & Editing; **Mayxay M**: Investigation, Resources, Writing – Review & Editing; **Miles A**: Investigation, Software, Writing – Review & Editing; **Miotto O**: Data Curation, Formal Analysis, Investigation, Project Administration, Software, Writing – Review & Editing; **Mobegi V**: Investigation, Resources, Writing – Review & Editing; **Mokuolu OA**: Investigation, Resources, Writing – Review & Editing; **Montgomery J**: Investigation, Resources, Writing – Review & Editing; **Mueller I**: Investigation, Resources, Writing – Review & Editing; **Newton PN**: Investigation, Resources, Writing – Review & Editing; **Nguyen T**: Data Curation, Investigation, Software, Writing – Review & Editing; **Nguyen TN**: Investigation, Resources, Writing – Review & Editing; **Noedi H**: Investigation, Resources, Writing – Review & Editing; **Nosten F**: Investigation, Resources, Writing – Review & Editing; **Noviyanti R**: Investigation, Resources, Writing – Review & Editing; **Nzila A**: Investigation, Resources, Writing – Review & Editing; **Ochola-Oyier LI**: Investigation, Resources, Writing – Review & Editing; **Ocholla H**: Investigation, Resources, Writing – Review & Editing; **Oduro A**: Investigation, Resources, Writing – Review & Editing; **Omedo I**: Investigation, Resources, Writing – Review & Editing; **Onyamboko MA**: Investigation, Resources, Writing – Review & Editing; **Ouedraogo JB**: Investigation, Resources, Writing – Review & Editing; **Oyebola K**: Investigation, Resources, Writing – Review & Editing; **Pearson RD**: Conceptualization, Data Curation, Formal Analysis, Investigation, Project Administration, Software, Supervision, Writing – Original Draft Preparation, Writing – Review & Editing; **Peshu N**: Investigation, Resources, Writing – Review & Editing; **Phyo AP**: Investigation, Resources, Writing – Review & Editing; **Plowe CV**: Investigation, Resources, Writing – Review & Editing; **Price RN**: Investigation, Resources, Writing – Review & Editing; **Pukrittayakamee S**: Investigation, Resources, Writing – Review & Editing; **Randrianariveolosia M**: Investigation, Resources, Writing – Review & Editing; **Rayner JC**: Investigation, Resources, Writing – Review & Editing; **Ringwald P**: Investigation, Resources, Writing – Review & Editing; **Rockett KA**: Investigation, Project Administration, Writing – Review & Editing; **Rowlands K**: Investigation, Writing – Review & Editing; **Ruiz L**: Investigation, Resources, Writing – Review & Editing; **Saunders D**: Investigation, Resources, Writing – Review & Editing; **Shayo A**: Investigation, Resources, Writing – Review & Editing; **Siba P**: Investigation, Resources, Writing – Review & Editing; **Simpson VJ**: Investigation, Project Administration, Writing – Review & Editing; **Stalker J**: Data Curation, Investigation, Software, Writing – Review & Editing; **Su Xz**: Investigation, Resources, Writing – Review & Editing; **Sutherland C**: Investigation, Resources, Writing – Review & Editing; **Takala-Harrison S**: Investigation, Resources, Writing – Review & Editing; **Tavul L**: Investigation, Resources, Writing – Review & Editing; **Thathy V**: Investigation, Resources, Writing – Review & Editing; **Tshefu A**: Investigation, Resources, Writing – Review & Editing; **Verra F**: Investigation, Resources, Writing – Review & Editing; **Vinetz J**: Investigation, Resources, Writing – Review & Editing; **Wellems TE**: Investigation, Resources, Writing – Review & Editing; **Wendler J**: Investigation, Resources, Writing – Review & Editing; **White NJ**: Investigation, Resources, Writing – Review & Editing; **Wright I**: Investigation, Software, Writing – Review & Editing; **Yavo W**: Investigation, Resources, Writing – Review & Editing; **Ye H**: Investigation, Resources, Writing – Review & Editing

**Competing interests:** No competing interests were disclosed.

**Grant information:** The sequencing, analysis, informatics and management of the Community Project are supported by Wellcome through Sanger Institute core funding (098051), a Strategic Award (090770/Z/09/Z) and the Wellcome Centre for Human Genetics core funding (203141/Z/16/Z), by the MRC Centre for Genomics and Global Health which is jointly funded by the Medical Research Council and the Department for International Development (DFID) (G0600718; M006212), and by the Bill & Melinda Gates Foundation (OPP1204628). *The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*

**Copyright:** © 2021 MalariaGEN *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. The author(s) is/are employees of the US Government and therefore domestic copyright protection in USA does not apply to this work. The work may be protected under the copyright laws of other jurisdictions when used in those jurisdictions.

**How to cite this article:** MalariaGEN, Ahouidi A, Ali M *et al.* **An open dataset of *Plasmodium falciparum* genome variation in 7,000 worldwide samples [version 2; peer review: 2 approved]** Wellcome Open Research 2021, 6:42  
<https://doi.org/10.12688/wellcomeopenres.16168.2>

**First published:** 24 Feb 2021, 6:42 <https://doi.org/10.12688/wellcomeopenres.16168.1>

**REVISED Amendments from Version 1**

We are grateful to the reviewers for their suggestions and have updated the manuscript in response. We now include gene IDs every time a gene is mentioned for the first time in the manuscript. We have replaced “complex rearrangements” in the results section with an explicit description of the event. We have added a paragraph to detail that sample collection is heterogeneous and due care is needed when interpreting the results. No changes have been made to the data or figures.

**Any further responses from the reviewers can be found at the end of the article**

**Introduction**

A major obstacle to malaria elimination is the great capacity of the parasite and vector populations to evolve in response to malaria control interventions. The widespread use of chloroquine and DDT in the 1950’s led to high levels of drug and insecticide resistance, and the same pattern has been repeated for other first-line antimalarial drugs and insecticides. Over the past 15 years, mass distribution of pyrethroid-treated bednets in Africa and worldwide use of artemisinin combination therapy (ACT) has led to substantial reductions in malaria prevalence and mortality, but there are rapidly increasing levels of resistance to ACT in Southeast Asian parasites and of pyrethroid resistance in African mosquitoes. A deep understanding of local patterns of resistance and the continually changing nature of the local parasite and vector populations is necessary to manage the use of drugs and insecticides and to deploy public health resources for maximum sustainability and impact.

Current methods for genetic surveillance of the parasite population are largely based on targeted genotyping of specific loci, e.g. known markers of drug resistance. Whole genome sequencing of malaria parasites is currently more expensive and complex, particularly at the stage of data analysis, but it is an important adjunct to targeted genotyping, as it provides a more comprehensive picture of parasite genetic variation. It is particularly important for discovery of new drug resistance markers and for monitoring patterns of gene flow and evolutionary adaptation in the parasite population.

The *Plasmodium falciparum* Community Project (*Pf* Community Project) was established with the aim of integrating parasite genome sequencing into clinical and epidemiological studies of malaria ([www.malariagen.net/projects](http://www.malariagen.net/projects)). It forms part of the Malaria Genomic Epidemiology Network (MalariaGEN), a global data-sharing network comprising multiple partner studies, each with its own research objectives and led by a local investigator<sup>1</sup>. Genome sequencing was performed centrally, and partner studies were free to analyse and publish the genetic data produced on their own samples, in line with MalariaGEN’s guiding principles on equitable data sharing<sup>1-3</sup>. A programme of capacity building for research into parasite genetics was developed at multiple sites in Africa alongside the *Pf* Community Project<sup>4</sup>.

The first phase of the project focused on developing simple methods to obtain purified parasite genome DNA from small blood samples collected in the field<sup>5,6</sup> and on establishing reliable computational methods for variant discovery and genotype calling from short-read sequencing data<sup>7</sup>. This presented a number of analytical challenges due to long tracts of highly repetitive sequence and hypervariable regions within the *P. falciparum* genome, and also because a single infection can contain a complex mixture of genotypes. Once a reliable analysis pipeline was in place, a process was established for periodic data releases to partners, with continual improvements in data quality as new analytical methods were developed.

Data from the *Pf* Community Project were initially released through a companion project called **Pf3k**, whose goal was to bring together leading analysts from multiple institutions to benchmark and standardise methods of variant discovery and genotyping calling. A **visual analytics web application** was developed<sup>8</sup> for researchers to explore the data. The open dataset was enlarged in 2016 when multiple partner studies contributed to a consortial publication on 3,488 samples from 23 countries<sup>9</sup>.

Data produced by the *Pf* Community Project have been used to address a broad range of research questions, both by the groups that generated samples and data and by the wider research community, and have generated over 50 previous publications (refs 5–55). These data have become a key resource for the epidemiology and population genetics of antimalarial drug resistance<sup>9-22</sup> and an important platform for the discovery of new genetic markers and mechanisms of resistance through genome-wide association studies<sup>23-27</sup> and combined genome-transcriptome analysis<sup>28</sup>. The data have also been used to study gene deletions that cause failure of rapid diagnostic tests<sup>29</sup>; to characterise genetic variation in malaria vaccine antigens<sup>30,31</sup>; to screen for new vaccine candidates<sup>32</sup>; to investigate specific host-parasite interactions<sup>33,34</sup>; and to describe the evolutionary adaptation and diversification of local parasite populations<sup>7,9,12,35-40</sup>.

The *Pf* Community Project data also provide an important resource for developing and testing new analytical and computational methods. A key area of methods development is quantification of within-host diversity<sup>7,41-46</sup>, estimation of inbreeding<sup>7,47</sup>, and deconvolution of mixed infections into individual strains<sup>48,49</sup>. The data have also been used to develop and test methods for estimating identity by descent<sup>50,51</sup>, imputation<sup>52</sup>, typing structural variants<sup>53</sup>, designing other SNP genotyping platforms<sup>54</sup> and data visualisation<sup>8,55</sup>. In a companion study we performed whole genome sequencing of experimental genetic crosses of *P. falciparum*, and this provided a benchmark to test the accuracy of our genotyping methods, and to conduct an in-depth analysis of indels, structural variants and recombination events which are complicated to ascertain in these population genetic samples<sup>56</sup>.

Here we describe a new release of curated genome variation data on 7,113 samples of *P. falciparum* collected by 49 partner

studies from 73 locations in Africa, Asia, South America and Oceania between 2002 and 2015 (Table 1, Supplementary Data; Supplementary Table 1 and 2).

## Results

### Variant discovery and genotyping

We used the Illumina platform to produce genome sequencing data on all samples and we mapped the sequence reads against the *P. falciparum* 3D7 v3 reference genome. The median depth of coverage was 73 sequence reads averaged across the whole genome and across all samples. We constructed an analysis pipeline for variant discovery and genotyping, including stringent quality control filters that took into account the unusual features of the *P. falciparum* genome, incorporating lessons learnt from our previous work<sup>7,56</sup> and the Pf3k project, as outlined in the *Methods* section.

In the first stage of analysis we discovered variation at over six million positions, corresponding to about a quarter of the 23 Mb *P. falciparum* genome (Supplementary Data; Supplementary Table 3). These included 3,168,721 single nucleotide polymorphisms (SNPs): these were slightly more common in coding than non-coding regions and were mostly biallelic. The remaining 2,882,975 variants were predominantly short indels but also included more complex combinations of SNPs and indels: these were much more abundant in non-coding than coding regions, and mostly had at least three alleles. The predominance of indels in non-coding regions has been previously observed and is most likely a consequence of the extreme AT bias which leads to many short repetitive sequences<sup>56,57</sup>.

For the purpose of this analysis, we excluded all variants in subtelomeric and internal hypervariable regions, mitochondrial and apicoplast genomes, and some other regions of the genome where the mapping of short sequence reads is prone to a high error rate due to extremely high rates of variation<sup>56</sup>. A total of 1,838,733 SNPs (of which 1,626,886 were biallelic) and 1,276,027 indels (or SNP/indel combinations) passed all these filters. The pass rate for SNPs in coding regions (66%) was considerably higher than that for SNPs in non-coding regions (47%), indels in coding regions (37%) and indels in non-coding regions (47%). Finally, we removed samples with a low genotyping success rate or other quality control issues. We also removed replicates and 41 samples with genetic markers of infection by multiple *Plasmodium* species, leaving 5,970 high-quality samples from 28 countries (Table 1).

We used coverage and read pair analysis to determine duplication genotypes around *mdr1* (PF3D7\_0523000), *plasmepsin2/3* (PF3D7\_1408000 and PF3D7\_1408100) and *gch1* (PF3D7\_1224000), each of which are associated with drug resistance. For each of these three genes we discovered many different sets of breakpoints (29, 10 and 3 pairs of breakpoints for *mdr1*, *gch1*, and *plasmepsin 2/3*, respectively), including a large and complex structural rearrangement involving a triplicated segment embedded within a duplication, in which the triplicated segment is inverted (“dup-trpinv-dup”)<sup>58</sup> that to the best of our knowledge has not been observed before in *Plasmodium* species (Supplementary Data; Supplementary Note,

Supplementary Tables 4–6). We also used sequence reads coverage to identify large structural variants that appear to delete or disrupt *hrp2* (PF3D7\_0831800) and *hrp3* (PF3D7\_1372200), an event that can cause rapid diagnostic tests to malfunction.

The population genetic analyses in this paper are based on the filtered dataset of high-quality SNP genotypes in 5,970 samples. These data are openly available, together with annotated genotyping data on 6 million putative variants in all 7,113 samples, plus details of partner studies and sampling locations, at [www.malariagen.net/resource/26](http://www.malariagen.net/resource/26).

### Global population structure

The genetic structure of the global parasite population reflects its geographic regional structure<sup>7,9,10</sup> as illustrated by a neighbour-joining tree and a principal component analysis of all samples based on their SNP genotypes (Figure 1). Based on these observations we grouped the samples into eight geographic regions: West Africa, Central Africa, East Africa, South Asia, the western part of Southeast Asia, the eastern part of Southeast Asia, Oceania and South America. Each of these can be viewed as a regional sub-population of parasites, which is more or less differentiated from other regional sub-populations depending on rates of gene flow and other factors. The different regions encompass a range of epidemiological and environmental settings, varying in transmission intensity, vector species and history of antimalarial drug usage. Note these regional classifications are intentionally broad, and therefore overlook many interesting aspects of local population structure, e.g. a distinctive Ethiopian sub-population can be identified by more detailed analysis of African samples<sup>12</sup>.

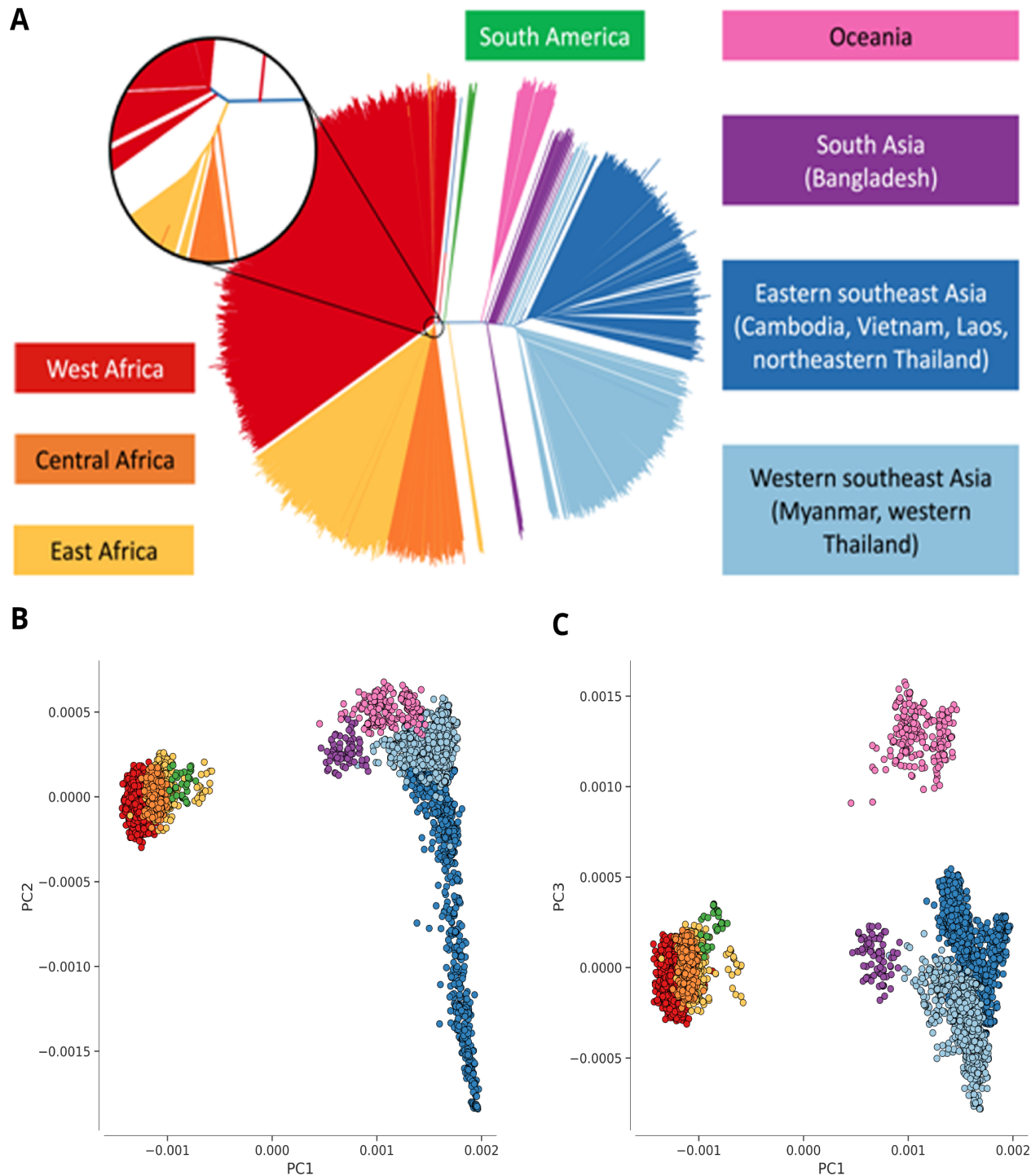
Genetically mixed infections were considerably more common in Africa than other regions, consistent with the high intensity of malaria transmission in Africa (Figure 2a). Analysis of  $F_{ws}$ , a measure of within-host diversity<sup>7</sup>, shows that most samples from Southeast Asia (1763/2341), South America (37/37) and Oceania (158/201) have  $F_{ws} > 0.95$ , which to a first approximation indicates that the infection is dominated by a clonal population of parasite<sup>41</sup>. In contrast, nearly half of samples from Africa (1625/3314) have  $F_{ws} < 0.95$ , indicating the presence of more complex infections. Genetically mixed infections were also common in Bangladesh (41/77 samples have  $F_{ws} < 0.95$ ), another area of high malaria transmission and the only South Asian country represented in this dataset, but did not reach the extremely high levels of within-host diversity ( $F_{ws} < 0.2$ ) observed in some samples from Africa.

The average nucleotide diversity across the global sample collection was 0.040% (median=0.028%), i.e. two randomly-selected samples differ by an average of 4 nucleotide positions per 10kb. Levels of nucleotide diversity vary greatly across the genome<sup>56</sup> and also geographically (Figure 2b). Distributions of values were highest in Africa, followed by Bangladesh, but the scale of regional differences was relatively modest, ranging from an average of 0.030% in Eastern Southeast Asia to 0.040% in West Africa (median=0.019% and 0.028% respectively; Figure 2b). In other words, the nucleotide diversity of each regional parasite population was not much less than that of

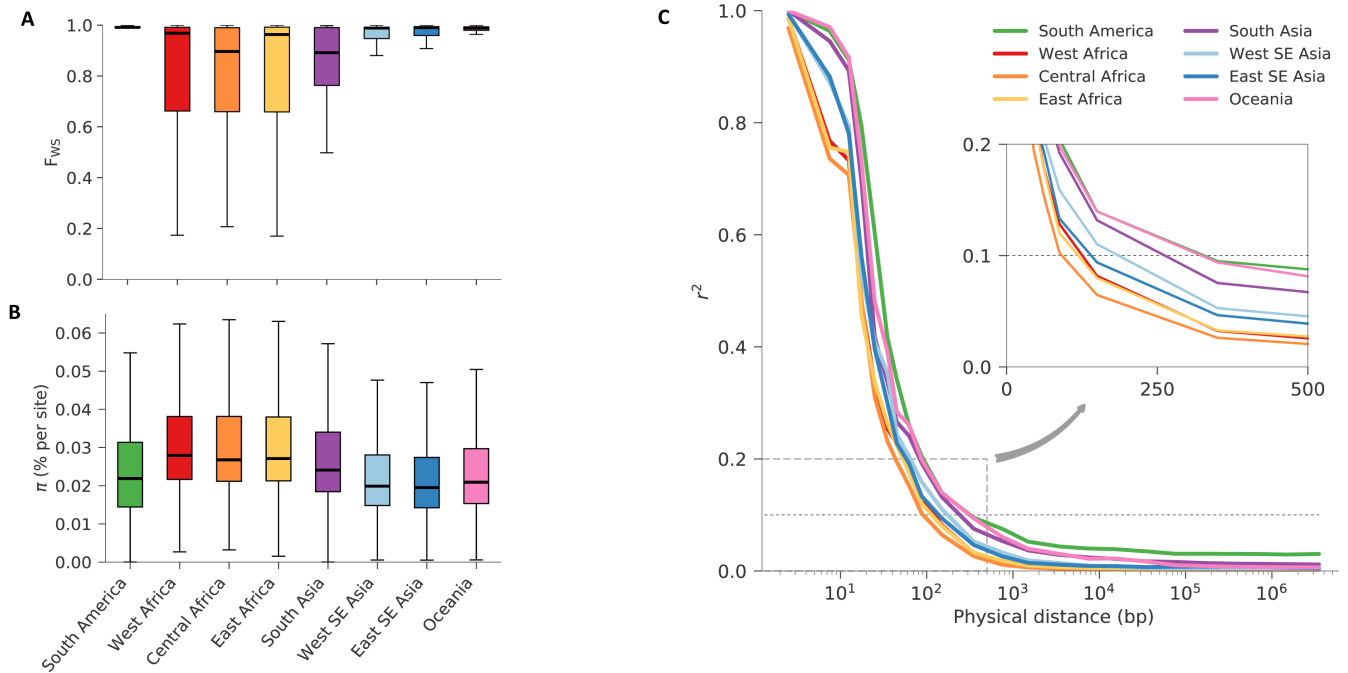


**Table 1. Count of samples in the dataset.** Countries are grouped into eight geographic regions based on their geographic and genetic characteristics. For each country, the table reports: the number of distinct sampling locations; the total number of samples sequenced; the number of high-quality samples included in the analysis; and the percentage of samples collected between 2012–2015, the most recent sampling period in the dataset. Eight samples were obtained from travellers returning from an endemic country, but where the precise site of the infection could not be determined. These were reported from Ghana (3 sequenced samples/2 analysis set samples), Kenya (2/1), Uganda (2/1) and Mozambique (1/1). “Lab samples” contains all sequences obtained from long-term *in vitro* cultured and adapted isolates, e.g. laboratory strains. The breakdown by site is reported in Supplementary table 1 and the list of contributing studies in Supplementary table 2.

Region	Country	Sampling locations	Sequenced samples	Analysis set samples	% analysis samples 2012–2015
<b>South America (SAM)</b>	<b>Colombia</b>	4	16	16	0%
	<b>Peru</b>	2	23	21	0%
<b>West Africa (WAF)</b>	<b>Benin</b>	1	102	36	100%
	<b>Burkina Faso</b>	1	57	56	0%
	<b>Cameroon</b>	1	239	235	100%
	<b>Gambia</b>	4	277	219	67%
	<b>Ghana</b>	3	1,003	849	56%
	<b>Guinea</b>	2	197	149	0%
	<b>Ivory Coast</b>	3	70	70	100%
	<b>Mali</b>	5	449	426	80%
	<b>Mauritania</b>	4	86	76	100%
	<b>Nigeria</b>	2	42	29	97%
<b>Senegal</b>	1	86	84	100%	
<b>Central Africa (CAF)</b>	<b>Congo DR</b>	1	366	344	100%
<b>East Africa (EAF)</b>	<b>Ethiopia</b>	2	34	21	100%
	<b>Kenya</b>	3	129	109	55%
	<b>Madagascar</b>	3	25	24	100%
	<b>Malawi</b>	2	351	254	0%
	<b>Tanzania</b>	5	350	316	85%
	<b>Uganda</b>	1	14	12	0%
<b>South Asia (SAS)</b>	<b>Bangladesh</b>	2	93	77	64%
<b>Western Southeast Asia (WSEA)</b>	<b>Myanmar</b>	5	250	211	71%
	<b>Western Thailand</b>	2	962	868	24%
<b>Eastern Southeast Asia (ESEA)</b>	<b>Cambodia</b>	5	1,214	896	32%
	<b>Northeastern Thailand</b>	1	28	20	75%
	<b>Laos</b>	2	131	120	21%
	<b>Viet Nam</b>	2	264	226	11%
<b>Oceania (OCE)</b>	<b>Indonesia</b>	1	92	80	73%
	<b>Papua New Guinea</b>	3	139	121	63%
<b>Returning travellers</b>	<b>Various locations</b>	0	8	5	0%
<b>Lab samples</b>	<b>Various locations</b>	0	16	0	0%
<b>Total</b>		<b>73</b>	<b>7,113</b>	<b>5,970</b>	<b>52%</b>



**Figure 1. Population structure.** (A) Genome-wide unrooted neighbour-joining tree showing population structure across all sites, with sample branches coloured according to country groupings (Table 1): South America (green, n=37); West Africa (red, n=2231); Central Africa (orange, n=344); East Africa (yellow, n=739); South Asia (purple, n=77); West Southeast Asia (light blue; n=1079); East Southeast Asia (dark blue; n=1262); Oceania (magenta; n=201). The circular inset shows a magnified view of the part of the tree where the majority of samples from Africa coalesce, showing that the three African sub-regions are genetically close but distinct. (B, C) First three component of a genome-wide principal coordinate analysis. The first axis (PC1) captures the separation of African and South American from Asian samples. The following two axes (PC2 and PC3) capture finer levels of population structure due to geographical separation and selective forces. Each point represents a sample and the colour legend is the same as above.



**Figure 2. Characteristics of the eight regional parasite populations.** (A) Distribution of within-host diversity, as measured by  $F_{ws}$ , showing that genetically mixed infections were considerably more common in Africa than other regions, consistent with the high intensity of malaria transmission in Africa. (B) Distribution of per site nucleotide diversity calculated in non-overlapping 25kbp genomic windows. We only considered coding biallelic SNPs to reduce the ascertainment bias caused by poor accessibility of non-coding regions. In both previous panels, thick lines represent median values, boxes show the interquartile range, and whiskers represent the bulk of the distribution, discounting outliers. (C) Genome-wide median LD (y-axis, measured by  $r^2$ ) between pairs of SNPs as function of their physical distance (x-axis, in bp), showing a rapid decay in all regional parasite populations. The inset panel shows a magnified view of the decay, showing that in all populations  $r^2$  decayed below 0.1 (dashed horizontal line) within 500 bp. All panels utilise the same palette, with colours denoting each geographic region.

the global parasite population. This is consistent with the idea that the global *P. falciparum* population has a common African origin and that historically there must have been significant levels of migration.

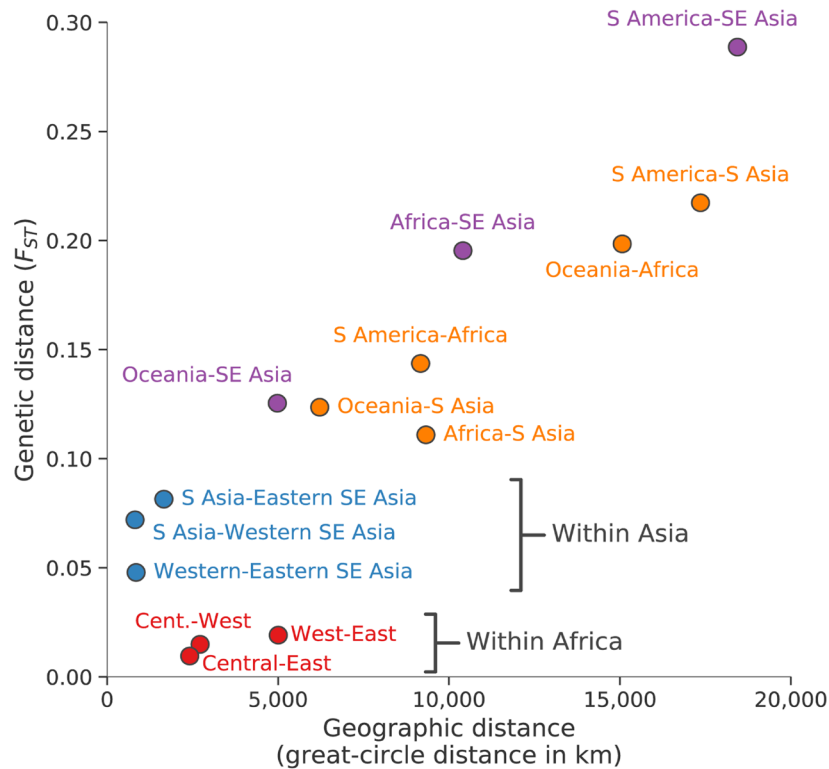
All regional sub-populations showed very low levels of linkage disequilibrium relative to human populations, e.g.  $r^2$  decayed to  $<0.1$  within 500 bp (Figure 2c). As expected, African populations had the highest rates of LD decay, implying the highest levels of haplotype diversity.

### Geographic patterns of population differentiation and gene flow

Parasite sub-populations in different locations naturally tend to differentiate over time unless there is sufficient gene flow to counterbalance genetic drift. Genome-wide estimates of  $F_{ST}$  provide an indicator of this process of genetic differentiation, which is partly determined by geographic distance (Figure 3). For example, we observe much greater genetic differentiation between South America and South Asia (genome-wide average  $F_{ST}$  0.22) or between Africa and Oceania (0.20) than between sub-regions within Asia ( $<0.1$ ) or within Africa ( $<0.02$ ).

These data reveal some interesting exceptions to the general rule that genome-wide  $F_{ST}$  is correlated with geographic distance. For example, African parasites are more strongly differentiated from Southeast Asian parasites (genome-wide average  $F_{ST}$  0.20) than they are from parasites in neighbouring Bangladesh (0.11). If this is examined in more detail, there is an unexpectedly steep gradient of genetic differentiation at the geographical boundary between South Asia and Southeast Asia, i.e. parasites sampled in Myanmar and Western Thailand are much more strongly differentiated from parasites sampled in Bangladesh (genome-wide  $F_{ST}$  0.07) than would be expected given that these are neighbouring countries. As discussed later, Southeast Asia is the global epicentre of antimalarial drug resistance, and these observations add to a growing body of evidence that Southeast Asian parasites have acquired a wide range of genomic features that are likely due to natural selection rather than genetic drift<sup>23,40</sup>.

It is noteworthy that the level of genetic differentiation between western and eastern parts of Southeast Asia (genome-wide  $F_{ST}$  0.05) is greater than between West Africa and East Africa (0.02) although the geographic distances are much greater in Africa. This is likely due to the lower intensity of malaria



**Figure 3. Geographic patterns of population differentiation and gene flow.** Each point represents one pairwise comparison between two regional parasite populations. The x-axis reports the geographic separation between the two populations, measured as great-circle distance between the centre of mass of each population and without taking into account natural barriers. The y-axis reports the genetic differentiation between the two populations, measured as average genome-wide  $F_{ST}$ . Points are coloured based on the regional populations they represent: between African populations (red); between Asian populations (blue); between Southeast Asia (as a whole) and Oceania, Africa or South America (purple); all the rest (orange).

transmission in Southeast Asia, and in particular the presence of a malaria-free corridor running through Thailand, which act as barriers to gene flow across the region<sup>23,40</sup>.

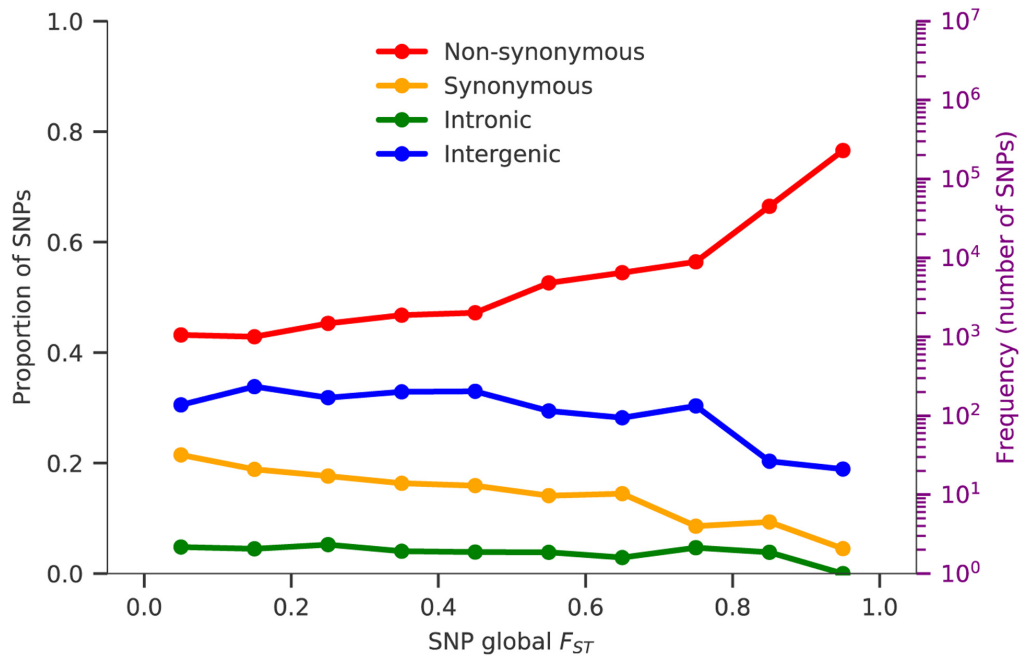
#### Genes with high levels of geographic differentiation

The  $F_{ST}$  metric can also be calculated for individual variants to identify specific genes that have acquired high levels of geographic differentiation relative to the genome as a whole. This can be done either at the global level (to identify variants that are highly differentiated between different regions of the world) or at the local level (to identify variants that are highly differentiated between different sampling locations within a region).

To identify variants that are strongly differentiated at the global level, we began by estimating  $F_{ST}$  for each SNP across all of the eight regional sub-populations. The group of SNPs with the highest global  $F_{ST}$  levels were found to be strongly enriched for non-synonymous mutations, suggesting that the process of differentiation is at least in part due to natural selection (Figure 4). After ranking all SNPs according to their global  $F_{ST}$  value, we calculated a *global differentiation score* for each gene based on the highest-ranking non-synonymous SNP

within the gene (see *Methods*). All genes are ranked according to their global differentiation score in the accompanying data release, and those with the highest score are listed in Supplementary Table 7 (Supplementary Data). The most highly differentiated gene, *p47* (PF3D7\_1346800), is known to interact with the mosquito immune system<sup>39</sup> and has two variants (S242L and V247A) that are at fixation in South America but absent in other geographic regions. Also among the five most highly differentiated genes are *gig* (PF3D7\_0935600, implicated in gametocytogenesis<sup>60</sup>), *pfs16*, (PF3D7\_0406200, expressed on the surface of gametes<sup>61</sup>) and *ctrp* (PF3D7\_0315200, expressed on the ookinete cell surface and essential for mosquito infection<sup>62</sup>). Thus, four of the five most highly differentiated parasite genes are involved in the process of transmission by the mosquito vector, raising the possibility that this reflects evolutionary adaptation of the *P. falciparum* population to the different *Anopheles* species that transmit malaria in different geographical regions.

It is more difficult to characterise variants that are strongly differentiated at the local level, due to smaller sample sizes and various sources of sampling bias, but a crude estimate can be obtained by analysis of each of the six geographical regions



**Figure 4. SNPs geographic differentiation.** Coloured lines show the proportions of SNPs in ten  $F_{ST}$  bins, stratified by genomic regions: non-synonymous (red), synonymous (yellow), intronic (green) and intergenic (blue).  $F_{ST}$  is calculated between all eight regional parasite populations and the number of SNPs in each bin is indicated in the background histogram. The y-axis on the right-hand side refers to the histogram and is on a log scale.

with samples from multiple countries.  $F_{ST}$  was estimated for each SNP across different sampling locations within each geographical region, and the results for different regions were combined by a heuristic approach to obtain a *local differentiation score* for each gene (see *Methods*). A range of genes associated with drug resistance (*crt* (PF3D7\_0709000), *dhfr* (PF3D7\_0417200), *dhps* (PF3D7\_0810800), *kelch13* (PF3D7\_1343700), *mdr1* (PF3D7\_0523000), *mdr2* (PF3D7\_1447900) and *fd* (PF3D7\_1318100)) were in the top centile of local differentiation scores (Supplementary Data; Supplementary Figure 1, Supplementary Table 8, Supplementary Note).

### Geographic patterns of drug resistance

#### *Classification of samples based on markers of drug resistance.*

Antimalarial drug resistance represents a major focus of research for many partner studies within the *Pf* Community Project, and this dataset therefore contains a significant body of data that have appeared in previous reports on drug resistance. Readers are referred to these publications for more detailed analyses of local patterns of resistance<sup>9-14,16-22</sup> and of resistance to specific drugs including chloroquine<sup>16,21</sup>, sulfadoxine-pyrimethamine<sup>16,19,21</sup> and artemisinin combination therapy<sup>9-11,13-15,17,18,21,22</sup>.

Here we have classified all samples into different types of drug resistance based on published genetic markers and current knowledge of the molecular mechanisms (see [www.malariagen.net/resource/26](http://www.malariagen.net/resource/26) for details of the heuristic used). Table 2 summarises the frequency of different types of drug resistance in samples from different geographical regions. Overall, we

observed higher prevalence of samples classified as resistant in Southeast Asia than anywhere else, with multiple samples resistant to all drugs considered. Note that samples were collected over a relatively long time period (2002–15) during which there were major changes in global patterns of drug resistance, and that the sampling locations represented in a given year depended on which partner studies were operative at the time. To alleviate this problem, we have also divided the data into samples collected before and after 2011 (Supplementary Data; Supplementary table 10), but temporal trends in aggregated data should be interpreted with due caution.

Below we summarise the overall profile of drug resistance types in the regional sub-populations: this is intended simply to provide context for users of this dataset, and should not be regarded as a statement of the current epidemiological situation. The Supplementary Notes (Supplementary Data) contain a more detailed description of the geographical distribution of haplotypes, CNV breakpoints, interactions between genes, and variants associated with less commonly used antimalarial drugs. In the accompanying data release, we also identify samples with *mdr1*, *plasmepsin2/3* and *gch1* gene amplifications that can affect drug resistance.

**Chloroquine resistance.** Samples were classified as chloroquine resistant if they carried the *crt* 76T allele. As shown in Table 2, this was found in almost all samples from Southeast Asia, South America and Oceania. It was also found across Africa but at lower frequencies, particularly in East Africa where chloroquine resistance is known to have declined since

**Table 2. Cumulative frequency of different types of drug resistance in samples from different geographical regions.** All samples were classified into different types of drug resistance based on published genetic markers, and represent best attempt based on the available data. Each type of resistance was considered to be either present, absent or unknown for a given sample. For each resistance type, the table reports: the genetic markers considered; the drug they are associated with; the proportion of samples in each region classified as resistant out of the samples where the type was not unknown. The number of samples classified as either resistant or not resistant varies for each type of resistance considered (e.g. due to different levels of genomic accessibility); numbers in brackets reports the minimum and maximum number analysed while the exact numbers considered are reported in Supplementary table 9. SP: sulfadoxine-pyrimethamine; treatment: SP used for the clinical treatment of uncomplicated malaria; IPTp: SP used for intermittent preventive treatment in pregnancy; AS-MQ: artesunate + mefloquine combination therapy; DHA-PPQ: dihydroartemisinin + piperazine combination therapy. Details of the rules used to infer resistance status from genetic markers can be found on the resource page at [www.malariaigen.net/resource/26](http://www.malariaigen.net/resource/26).

Marker	Associated with resistance to	South America (n=33-37)	West Africa (n=1851-2231)	Central Africa (n=262-344)	East Africa (n=678-739)	South Asia (n=62-77)	Western Southeast Asia (n=906-1079)	Eastern Southeast Asia (n=867-1256)	Oceania (n=185-201)
<i>crt</i> 76T	Chloroquine	100%	41%	66%	14%	93%	100%	97%	99%
<i>dhfr</i> 108N	Pyrimethamine	97%	84%	100%	98%	100%	100%	100%	100%
<i>dhps</i> 437G	Sulfadoxine	30%	75%	97%	93%	97%	100%	87%	61%
<i>mdr1</i> 2+ copies	Mefloquine	0%	0%	0%	0%	0%	44%	12%	1%
<i>kelch13</i> WHO list	Artemisinin	0%	0%	0%	0%	0%	28%	46%	0%
<i>plasmeprin 2-3</i> 2+ copies	Piperaquine	0%	0%	0%	0%	0%	0%	17%	0%
<i>dhfr</i> triple mutant	SP (treatment)	0%	75%	82%	91%	43%	90%	92%	0%
<i>dhfr</i> and <i>dhps</i> sextuple mutant	SP (IPTp)	0%	0%	1%	10%	19%	82%	19%	0%
<i>kelch13</i> and <i>mdr1</i>	AS-MQ	0%	0%	0%	0%	0%	13%	9%	0%
<i>kelch13</i> and <i>plasmeprin 2-3</i>	DHA-PPQ	0%	0%	0%	0%	0%	0%	15%	0%

chloroquine was discontinued<sup>63–65</sup>. Supplementary Table 11 (Supplementary Data) shows the geographical distribution of different *crt* haplotypes (based on amino acid positions 72–76) which is consistent with the theory that chloroquine resistance spread from Southeast Asia to Africa with multiple independent origins in South America and Oceania<sup>66,67</sup>. The *crt* locus is also relevant to other types of drug resistance, e.g. *crt* variants that are relatively specific to Southeast Asia form the genetic background of artemisinin resistance, and newly emerging *crt* alleles have been associated with the spread of ACT failure due to piperazine resistance<sup>13,14,22,68</sup>.

**Sulfadoxine-pyrimethamine resistance.** Clinical resistance to sulfadoxine-pyrimethamine (SP) is determined by multiple mutations and their interactions, so following current practice<sup>69</sup> we classified SP resistant samples into four overlapping types: (i) carrying the *dhfr* 108N allele, associated with pyrimethamine resistance; (ii) the *dhps* 437G allele, associated with sulfadoxine resistance; (iii) carrying the *dhfr* triple mutant, which is strongly associated with SP failure; (iv) carrying the *dhfr/dhps* sextuple mutant, which confers a higher level of SP resistance. As shown in Table 2, *dhfr* 108N was found in almost all samples in all regions apart from West Africa, while *dhps* 437G was at very high frequency throughout most of Africa and Asia, and at lower frequencies in South America and Oceania (see also Supplementary Data; Supplementary Table 12). Triple mutant *dhfr* parasites were common throughout Africa and Asia, whereas sextuple mutant *dhfr/dhps* parasites were at much lower frequency except in Western Southeast Asia. In the accompanying data release, we also identify samples with *gchl* gene amplifications (Supplementary Data; Supplementary Table 4) that can modulate SP resistance<sup>70</sup>, although their effect on the clinical outcome and interaction with mutations in *dhfr* and *dhps* is not fully established.

**Resistance to artemisinin combination therapy.** We classified samples as artemisinin resistant based on the World Health Organization classification of non-synonymous mutations in the propeller region of the *kelch13* gene that have been associated with delayed parasite clearance<sup>71</sup>. By this definition, artemisinin resistance was confined to Southeast Asia but, as previously reported, this dataset contains a substantial number of non-synonymous *kelch13* propeller SNPs occurring at <5% frequency in Africa and elsewhere<sup>9</sup>. The most common ACT formulations in Southeast Asia are artesunate-mefloquine (AS-MQ) and dihydroartemisinin-piperazine (DHA-PPQ). We classified samples as mefloquine resistant if they had *mdr1* amplification<sup>72</sup> or as piperazine resistant if they had *plasmepsin 2/3* amplification<sup>25</sup>. Mefloquine resistance was observed throughout Southeast Asia and was most common in the western part. Piperazine resistance was confined to eastern Southeast Asia with a notable concentration in western Cambodia. Elsewhere<sup>11,13</sup> we describe the KEL1/PLA1 lineage of artemisinin- and piperazine-resistant parasites that expanded in western Cambodia during 2008–13, and then spread to other countries during 2013–18, causing high rates of DHA-PPQ treatment failure across eastern Southeast Asia: since the current dataset extends only to 2015 it captures only the first phase of the KEL1/PLA1 lineage expansion.

## HRP2/3 deletions that affect rapid diagnostic tests

Rapid diagnostic tests (RDTs) provide a simple and inexpensive way to test for parasites in the blood of patients who are suspected to have malaria, and have become a vital tool for malaria control<sup>73,74</sup>. The most widely used RDTs are designed to detect *P. falciparum* histidine-rich protein 2 and cross-react with histidine-rich protein 3, encoded by the *hrp2* and *hrp3* genes respectively. Parasites with gene deletions of *hrp2* and/or *hrp3* have emerged as an important cause of RDT failure in a number of locations<sup>75–79</sup>. It is difficult to devise a simple genetic assay to monitor for risk of RDT failure because *hrp2* and *hrp3* deletions comprise a diverse mixture of large structural variations with multiple independent origins, and both genes are located in subtelomeric regions of the genome with very high levels of natural variation<sup>29,80–83</sup>. In the absence of a well-validated algorithmic method, we visually inspected sequence read coverage and identified samples with clear evidence of large structural variants that disrupted or deleted the *hrp2* and *hrp3* genes. We took a conservative approach: samples that appeared to have a mixture of deleted and non-deleted genotypes were classified as non-deleted.

Deletions were found at relatively high frequency in Peru (8 of 21 samples had *hrp2* deletions, 14 had *hrp3* deletions and 6 had both) but were not seen in samples from Colombia and were relatively rare outside South America. Oceania was the only other region where we observed *hrp2* deletions, but at very low frequency (4%, n=3/80), and also had *hrp3* deletions (25%) though no combined deletions were seen. Deletions of *hrp3* only were more geographically widespread than *hrp2* deletions, being common in Ethiopia (43%, n=9/21) and in Senegal (7%, n=6/84), and at relatively low frequency (<5%) in Kenya, Cambodia, Laos, and Vietnam (Supplementary Data; Supplementary Table 13). Note that these findings might underestimate the true prevalence of *hrp2/hrp3* deletions, due to sampling bias (our samples were primarily collected from RDT-positive cases) and also because we focused on large structural variants and did not consider polymorphisms that might also cause RDT failure but would require more sophisticated analytical approaches. There is a need for more reliable diagnostics of *hrp2* and *hrp3* deletions, and we hope that these open data will accelerate this important area of applied methodological research.

## Discussion

This open dataset comprises sequence reads and genotype calls on over 7,000 *P. falciparum* samples from MalariaGEN partner studies in 28 countries. After excluding variants and samples that failed to meet stringent quality control criteria, the dataset contains high-quality genotype calls for 3 million polymorphisms including SNPs, indels, CNVs and large structural variations, in almost 6,000 samples. The data can be analysed in their entirety or can be filtered to select for specific genes, or geographical locations, or samples with particular genotypes. This is twice the sample size of our previous consortial publication<sup>9</sup> and is the largest available data resource for analysis of *P. falciparum* population structure, gene flow and evolutionary adaptation. Each sample has been annotated to show its profile of resistance to six major antimalarial drugs

and whether it carries structural variations that can cause RDT failure. The classification scheme is heuristic and based on a subset of known genetic markers, so it should not be treated as a failsafe predictor of the phenotype of a particular sample. Our purpose in providing these annotations is to make it easy for users without specialist training in genetics to explore the global dataset and to analyse any subset of samples for key features that are relevant to malaria control. Samples were collected by independent groups that were operative at a given time and in a given place with distinct objectives; while care needs to be taken when interpreting results spanning multiple years and geographical settings (e.g. aggregated trends of drug resistance prevalence), this heterogeneity also allows for the exploration of a wide range of epidemiological and transmission settings.

An important function of this curated dataset is to provide information on the provenance and key features of samples associated with each partner study, thus allowing the findings reported in different publications to be linked and compared. Data produced by the *Pf* Community Project have been analysed in more than 50 publications (refs 5–55) and a few examples will serve to illustrate the diverse ways in which the data are being used. An analysis of samples collected across Africa by Amambua-Ngwa, Djimde and colleagues found evidence that parasite population structure overlaps with historical patterns of human migration and that the *P. falciparum* population in Ethiopia is significantly diverged from other parts of the continent<sup>12</sup>. A series of studies by Amato, Miotto and colleagues have documented the evolution of a multidrug-resistant lineage of *P. falciparum* that originated in Western Cambodia over ten years ago and is now expanding rapidly across Southeast Asia, acquiring additional resistance mutations as it spreads<sup>11,13,14</sup>. McVean and colleagues have developed a computational method for deconvolution of the haplotypic structure of mixed infections, allowing analysis of the pedigree structure of parasites that are cotransmitted by the same mosquito<sup>49</sup>. Bahlo and colleagues have developed a different haplotype-based method to describe the relatedness structure of the parasite population and to identify new genomic loci with evidence of recent positive selection<sup>50</sup>.

A recent report from the World Health Organization highlights the need for improved surveillance systems in sustaining malaria control and achieving the long-term goal of malaria eradication<sup>84</sup>. To be of practical value for national malaria control programmes, genetic data must address well-defined use cases and be readily accessible<sup>85</sup>. Amplicon sequencing technologies provide a powerful new tool for targeted genotyping that could feasibly be implemented locally in malaria-endemic countries<sup>86,87</sup>, but there remains a need for the international malaria control community to generate and share whole genome sequencing data, e.g. to monitor for newly emerging forms of drug resistance and to understand regional patterns of parasite migration. The next generation of long-read sequencing technologies will improve the precision of population genomic inference, e.g. by enabling analysis of hypervariable regions of the genome, and of pedigree structures within mixed infections. The accuracy with which the resistance phenotype of a

sample can be predicted from genome sequencing data will also improve as we gain better functional understanding of the polygenic determinants of drug resistance.

Thus, the next few years are likely to see major advances in both the scale and information content of parasite genomic data. The practical value for malaria control will be greatly enhanced by the progressive acquisition of longitudinal time-series data, particularly if this is linked to other sources of epidemiological data and translated into reliable, actionable information with sufficient rapidity to allow control programmes to monitor the impact of their interventions on the parasite population in near real time. The *Pf* Community Project provides proof of concept that systems can be developed for groups in different countries to share data, to analyse it using standardised methods, and to make it readily accessible to other researchers and the malaria control community.

## Methods

Here we summarise the bioinformatics methods used to produce and analyse the data; further details are available at [www.malariagen.net/resource/26](http://www.malariagen.net/resource/26).

## Ethical approval

All samples in this study were derived from blood samples obtained from patients with *P. falciparum* malaria, collected with informed consent from the patient or a parent or guardian. At each location, sample collection was approved by the appropriate local and institutional ethics committees. The following local and institutional committees gave ethical approval for the partner studies: Human Research Ethics Committee of the Northern Territory Department of Health & Families and Menzies School of Health Research, Darwin, Australia; National Research Ethics Committee of Bangladesh Medical Research Council, Bangladesh; Comité d'Éthique de la Recherche - Institut des Sciences Biomedicales Appliquées, Benin; Ministère de la Santé – République du Bénin, Benin; Comité d'Éthique, Ministère de la Santé, Bobo-Dioulasso, Burkina Faso; Institutional Review Board Centre Muraz, Burkina Faso; Ministry of Health National Ethics Committee for Health Research, Cambodia; Institutional Review Board University of Buea, Cameroon; Comité Institucional de Ética de investigaciones en humanos de CIDEIM, Colombia; Comité National d'Éthique de la Recherche, Cote d'Ivoire; Comité d'Éthique Université de Kinshasa, Democratic Republic of Congo; Armauer Hansen Research Institute Institutional Review Board, Ethiopia; Addis Ababa University, Aklilu Lemma Institute of Pathobiology Institutional Review Board, Ethiopia; Kintampo Health Research Centre Institutional Ethics Committee, Ghana; Ghana Health Service Ethical Review Committee, Ghana; University of Ghana Noguchi Medical Research Institute, Ghana; Navrongo Health Research Centre Institutional Review Board, Ghana; Comité d'Éthique National Pour la Recherche en Santé, République de Guinée; Indian Council of Medical Research, India; Eijkman Institute Research Ethics Commission, Eijkman Institute for Molecular Biology, Jakarta, Indonesia; KEMRI Scientific and Ethics Review Unit, Kenya; Ministry of Health National Ethics Committee For Health Research, Laos; Ethical Review Committee of University of Ilorin Teaching Hospital, Nigeria; Comité National



d’Ethique auprès du Ministère de la Santé Publique, Madagascar; College of Medicine Regional Ethics Committee University of Malawi, Malawi; Faculté de Médecine, de Pharmacie et d’Odonto-Stomatologie, University of Bamako, Bamako, Mali; Ethics Committee of the Ministry of Health, Mali; Ethics committee of the Ministry of Health, Mauritania; Department of Medical Research (Lower Myanmar); Ministry of Health, Government of The Republic of the Union of Myanmar; : Institutional Review Board, Papua New Guinea Institute of Medical Research, Goroka, Papua New Guinea; PNG Medical Research Advisory Council (MRAC), Papua New Guinea; Institutional Review Board, Universidad Nacional de la Amazonia Peruana, Iquitos, Peru; Ethics Committee of the Ministry of Health, Senegal; National Institute for Medical Research and Ministry of Health and Social Welfare, Tanzania; Medical Research Coordinating Committee of the National Institute for Medical Research, Tanzania; Ethics Committee, Faculty of Tropical Medicine, Mahidol University, Bangkok, Thailand; Ethics Committee at Institute for the Development of Human Research Protections, Thailand; Gambia Government/MRC Joint Ethics Committee, Banjul, The Gambia; London School of Hygiene and Tropical Medicine Ethics Committee, London, UK; Oxford Tropical Research Ethics Committee, Oxford, UK; Walter Reed Army Institute of Research, USA; National Institute of Allergy and Infectious Diseases, Bethesda, MD, USA; Ethical Committee, Hospital for Tropical Diseases, Ho Chi Minh City, Vietnam; Ministry of Health Institute of Malariology-Parasitology-Entomology, Vietnam.

Standard laboratory protocols were used to determine DNA quantity and proportion of human DNA in each sample as previously described<sup>7,56</sup>.

### Data generation and curation

Reads mapping to the human reference genome were discarded before all analyses, and the remaining reads were mapped to the *P. falciparum* 3D7 v3 reference genome using [bwa mem](#)<sup>88</sup> version 0.7.15. “Improved” BAMs were created using the [Picard](#) tools [CleanSam](#), [FixMateInformation](#) and [MarkDuplicates](#) version 2.6.0 and [GATK](#) v3 base quality score recalibration. All lanes for each sample were merged to create sample-level BAM files.

We discovered potential SNPs and indels by running [GATK’s HaplotypeCaller](#)<sup>89</sup> independently across each of the 7,182 sample-level BAM files and genotyped these for each of the 16 reference sequences (14 chromosomes, 1 apicoplast and 1 mitochondria) using [GATK’s CombineGVCFs](#) and [GenotypeGVCFs](#).

SNPs and indels were filtered using [GATK’s Variant Quality Score Recalibration \(VQSR\)](#). Variants with a VQSLOD score  $\leq 0$  were filtered out. Functional annotations were applied using [snpEff](#)<sup>90</sup> version 4.1. Genome regions were annotated using [vcftools](#) version 0.1.10 and masked if they were outside the core genome. Unless otherwise specified, we used biallelic SNPs that pass all quality filters for all the analysis.

We removed 69 samples from lab studies to create the release VCF files which contain 7,113 samples. VCF files were

converted to [ZARR](#) format and subsequent analyses were mainly performed using [scikit-allele](#) version 1.1.18 and the ZARR files.

We identified species using nucleotide sequence from reads mapping to six different loci in the mitochondrial genome, using custom java code (available at <https://github.com/malariagen/GeneticReportCard>). The loci were located within the *cox3* gene (PF3D7\_MIT01400), as described in a previously published species detection method<sup>91</sup>. Alleles at various mitochondrial positions within the six loci were genotyped and used for classification as shown in Supplementary Table 14 (Supplementary Data).

We created a final analysis set of 5,970 samples after removing replicate, low coverage, suspected contaminations or mislabelling and mixed-species samples.

### Genotyping of drug resistance markers and samples classification

We used two complementary methods to determine tandem duplication genotypes around *mdr1*, *plasmepsin2/3* and *gch1*, namely a coverage-based method and a method based on position and orientation of reads near discovered duplication breakpoints. In brief, the outline algorithm is: (1) Determine copy number at each locus using a coverage based hidden Markov model (HMM); (2) Determine breakpoints of identified duplications by manual inspection of reads and face-away read pairs around all sets of breakpoints; (3) for each locus in each sample, initially set copy number to that determined by the HMM if  $\leq 10$  CNVs discovered in total, else consider undetermined; (4) if face-away pairs provide self-sufficient evidence for the presence or absence of the amplification, override the HMM call; (5) for each locus in each sample, set the breakpoint to be that with the highest proportion of face-away reads.

We genotyped deletions in *hrp2* and *hrp3* by manual inspection of sequence read coverage plots.

The procedure used to map genetic markers to inferred resistance status classification is described in detail for each drug in the accompanying data release (<https://www.malariagen.net/resource/26>).

In brief, we called amino acids at selected loci by first determining the reference amino acids and then, for each sample, applying all variations using the GT field of the VCF file. The amino acid and copy number calls generated were used to classify all samples into different types of drug resistance. Our methods of classification were heuristic and based on the available data and current knowledge of the molecular mechanisms. Each type of resistance was considered to be either present, absent or unknown for a given sample.

### Population-level analysis and characterisation

We calculate genetic distance between samples using biallelic SNPs that pass filters using a method previously described<sup>9</sup>. In addition to calculating genetic distance between all pairs of samples from the current data set, we also calculated the genetic

distance between each sample and the lab strains 3D7, 7G8, GB4, HB3 and Dd2 from the [Pf3k project](#).

The matrix of genetic distances was used to generate neighbour-joining trees and principal coordinates. Based on these observations we grouped the samples into eight geographic regions: South America, West Africa, Central Africa, East Africa, South Asia, the western part of Southeast Asia, the eastern part of Southeast Asia and Oceania, with samples assigned to region based on the geographic location of the sampling site. Five samples from returning travellers were assigned to region based on the reported country of travel.

$F_{ws}$  was calculated using custom python scripts using the method previously described<sup>7</sup>. Nucleotide diversity ( $\pi$ ) was calculated in non-overlapping 25 kbp genomic windows, only considering coding biallelic SNPs to reduce the ascertainment bias caused by poor accessibility of non-coding regions. LD decay ( $r^2$ ) was calculated using the method of Rogers and Huff and biallelic SNPs with low missingness and regional allele frequency >10%. Mean  $F_{ST}$  between populations was calculated using Hudson's method.

Allele frequencies stratified by geographic regions and sampling sites were calculated using the genotype calls produced by GATK.  $F_{ST}$  was calculated between all 8 regions, and also between all sites with at least 25 QC pass samples.  $F_{ST}$  between different locations for individual SNPs was calculated using Weir and Cockerham's method.

We defined the global differentiation score for a gene as  $1 - \frac{N}{\max(N)}$ , where  $N$  is the rank of the non-synonymous SNP with the highest global  $F_{ST}$  value within that gene. To define the local differentiation score, we first calculated for each region containing multiple sites (WAF, EAF, SAS, WSEA, ESEA and OCE)  $F_{ST}$  for each SNP between sites within that region. For each gene, we then calculated the rank of the highest  $F_{ST}$  non-synonymous SNP within that gene for each of the six regions. We defined the local differentiation score for each gene using the second highest of these six ranks ( $N$ ), to ensure that the gene was highly ranked in at least two populations, i.e. to minimise the chance of artefactually ranked a gene highly due to a single variant in a single population. The final local differentiation score was normalised to ensure that the range of possible scores was between 0 and 1, local differentiation score was defined as  $1 - \frac{N}{\max(N)}$ .

An earlier version of this article can be found on bioRxiv (DOI: <https://doi.org/10.1101/824730>).

## Data availability

### Underlying data

Data are available under the MalariaGEN terms of use for the *Pf* Community Project: <https://www.malariagen.net/data/terms-use/p-falciparum-community-project-terms-use>. Depending on the

nature, format and content of the data, appropriate mechanisms have been utilised for data access, as detailed below.

This project contains the following underlying data that are available as an online resource: [www.malariagen.net/resource/26](http://www.malariagen.net/resource/26). Data are also available from Figshare.

Figshare: Supplementary data to: An open dataset of Plasmodium falciparum genome variation in 7,000 worldwide samples. <https://doi.org/10.6084/m9.figshare.13388603><sup>92</sup>.

- Study information: Details of the 49 contributing partner studies, including description, contact information and key people.
- Sample provenance and sequencing metadata: sample information including partner study information, location and year of collection, ENA accession numbers, and QC information for 7,113 samples from 28 countries.
- Measure of complexity of infections: characterisation of within-host diversity (FWS) for 5,970 QC pass samples.
- Drug resistance marker genotypes: genotypes at known markers of drug resistance for 7,113 samples, containing amino acid and copy number genotypes at six loci: *crt*, *dhfr*, *dhps*, *mdr1*, *kelch13*, *plasmepsin 2–3*.
- Inferred resistance status classification: classification of 5,970 QC pass samples into different types of resistance to 10 drugs or combinations of drugs and to RDT detection: chloroquine, pyrimethamine, sulfadoxine, mefloquine, artemisinin, piperaquine, sulfadoxine-pyrimethamine for treatment of uncomplicated malaria, sulfadoxine-pyrimethamine for intermittent preventive treatment in pregnancy, artesunate-mefloquine, dihydroartemisinin-piperaquine, *hrp2* and *hrp3* genes deletions.
- Drug resistance markers to inferred resistance status: details of the heuristics utilised to map genetic markers to resistance status classification.
- Gene differentiation: estimates of global and local differentiation for 5,561 genes.
- Short variants genotypes: Genotype calls on 6,051,696 SNPs and short indels in 7,113 samples from 29 countries, available both as VCF and zarr files.

### Extended data

This project contains the following underlying supplementary data available as a single document download: [www.malariagen.net/resource/26](http://www.malariagen.net/resource/26). Extended data are also available from Figshare.

Figshare: Supplementary data to: An open dataset of Plasmodium falciparum genome variation in 7,000 worldwide samples. <https://doi.org/10.6084/m9.figshare.13388603><sup>92</sup>.

'File9\_Pf\_6\_supplementary' contains the Supplementary Note, Supplementary Tables and Supplementary Figure:

- Supplementary Note
  - Analysis of local differentiation score

- The classic 76T chloroquine resistance mutation in *crt* is found on multiple haplotypes
- Sulphadoxine-pyrimethamine resistance is widespread and associated with many haplotypes
- *mdr1* duplications have many different breakpoints
- Artemisinin, piperazine, and mefloquine resistance
- No evidence of resistance to less commonly used antimalarials
- Supplementary Table 1. Breakdown of analysis set samples by geography.
- Supplementary Table 2. Studies contributing samples.
- Supplementary Table 3. Summary of discovered variant positions.
- Supplementary Table 4. Breakpoints of duplications of *gch1*.
- Supplementary Table 5. Breakpoints of duplications of *mdr1*.
- Supplementary Table 6. Breakpoints of duplications of *plasmepsin 2–3*.
- Supplementary Table 7. Genes ranked by global differentiation score.
- Supplementary Table 8. Genes ranked by local differentiation score.
- Supplementary Table 9. Number of samples used to determine proportions in Table 2.
- Supplementary Table 10. Frequencies of mutations associated with mono- and multi-drug resistance pre- and post-2011.
- Supplementary Table 11. Frequency of *crt* amino acid 72–76 haplotypes.
- Supplementary Table 12. Frequencies of *dhfr* (51, 59, 108, 164) and *dhps* (437, 540, 581, 613) multi-locus haplotypes.
- Supplementary Table 13. Frequency of *HRP2* and *HRP3* deletions by country.
- Supplementary Table 14. Alleles at six mitochondrial positions used for the species identification.
- Supplementary Figure 1. Histogram of local differentiation score for all genes.

Data hosted with Figshare are available under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/) (CC-BY 4.0).

#### Data analysis group

Pearson, RD\*, Amato, R\*, Hamilton, WL, Almagro-Garcia, J, Chookajorn, T, Kochakarn, T, Miotto, O, Kwiatkowski, DP

\*Joint analysis lead

#### Local study design, implementation and sample collection

Ahouidi, A, Amambua-Ngwa, A, Amaratunga, C, Amenga-Etego, L, Andagalu, B, Anderson, TJC, Apinjoh, T, Ashley, EA, Auburn, S, Awandare, G, Ba, H, Baraka, V, Barry, AE, Bejon, P, Bertin, GI, Boni, MF, Borrmann, S, Bousema, T, Branch, O, Bull, PC, Chotivanich, K, Claessens, A, Conway, D, Craig, A, D'Alessandro, U, Dama, S, Day, N, Denis, B, Diakite, M, Djimdé, A, Dolecek, C, Dondorp, A, Drakeley, C, Duffy, P, Echeverry, DF, Egwang, TG, Erko, B, Fairhurst, RM, Faiz, A, Fanello, CA, Fukuda, MM, Gamboa, D, Ghansah, A, Golassa, L, Harrison, GLA, Hien, TT, Hill, CA, Hodgson, A, Imwong, M, Ishengoma, DS, Jackson, SA, Kamaliddin, C, Kamau, E, Konaté, A, Kyaw, MP, Lim, P, Lon, C, Loua, KM, Maïga-Ascofaré, O, Marfurt, J, Marsh, K, Mayxay, M, Mobegi, V, Mokuolu, OA, Montgomery, J, Mueller, I, Newton, PN, Nguyen, TN, Noedl, H, Nosten, F, Noviyanti, R, Nzila, A, Ochola-Oyier, LI, Ocholla, H, Oduro, A, Omedo, I, Onyamboko, MA, Ouedraogo, J, Oyebola, K, Peshu, N, Phyto, AP, Plowe, CV, Price, RN, Pukrittayakamee, S, Randrianarivelojosia, M, Rayner, JC, Ringwald, P, Ruiz, L, Saunders, D, Shayo, A, Siba, P, Su, X, Sutherland, C, Takala-Harrison, S, Tavul, L, Thathy, V, Tshefu, A, Verra, F, Vinetz, J, Wellem, TE, Wendler, J, White, NJ, Yavo, W, Ye, H

#### Sequencing, data production and informatics

Pearson, RD, Stalker, J, Ali, M, Amato, R, Ariani, C, Busby, G, Drury, E, Hart, L, Hubbart, C, Jacob, CG, Jeffery, B, Jeffreys, AE, Jyothi, D, Kekre, M, Kluczynski, K, Malangone, C, Manske, M, Miles, A, Nguyen, T, Rowlands, K, Wright, I, Goncalves, S, Rockett, KA

#### Partner study support and coordination

Simpson, VJ, Miotto, O, Amato, R, Goncalves, S, Henrichs, C, Johnson, KJ, Pearson, RD, Rockett, KA, Kwiatkowski, DP

#### Acknowledgements

This study was conducted by the MalariaGEN *Plasmodium falciparum* Community Project, and was made possible by clinical parasite samples contributed by partner studies, whose investigators are represented in the author list and in the associated data release (<https://www.malariagen.net/resource/26>). This research was supported in part by the Intramural Research Programme of the NIH, NIAID. In addition, the authors would like to thank the following individuals who contributed to partner studies, making this study possible: Dr Eugene Laman for work in sample collection in the Republic of Guinea; Dr Abderahmane Tandia and Dr Yacine Deh and Dr Samuel Assefa for work in sample collection in Mauritania; Dr Ibrahim Sanogo for work in sample collection in Mali; Dr James Abugri and Dr Nicholas Amoako for work coordinating sample collection in Ghana. Genome sequencing was undertaken by the Wellcome Sanger Institute and we thank the staff of the Wellcome Sanger Institute Sample Logistics, Sequencing, and Informatics facilities for their contribution. The authors would like to thank Erin Courtier for her assistance with the journal submission. The views expressed here are solely those of the authors and do not reflect the views, policies or positions of the U.S.

Government or Department of Defense. Material has been reviewed by the Walter Reed Army Institute of Research. There is no objection to its presentation and/or publication. The opinions or assertions contained herein are the private views of the author, and are not to be construed as official, or as reflecting true views of the Department of the Army or the Department

of Defense. The investigators have adhered to the policies for protection of human subjects as prescribed in AR 70–25. PR is a staff member of the World Health Organization. PR alone is responsible for the views expressed in this publication and they do not necessarily represent the decisions, policy or views of the World Health Organization.

## References

- Malaria Genomic Epidemiology Network: **A global network for investigating the genomic epidemiology of malaria.** *Nature.* 2008; **456**(7223): 732–7. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Chokshi DA, Parker M, Kwiatkowski DP: **Data sharing and intellectual property in a genomic epidemiology network: policies for large-scale research collaboration.** *Bull World Health Organ.* 2006; **84**(5): 382–7. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Parker M, Bull SJ, de Vries J, et al.: **Ethical data release in genome-wide association studies in developing countries.** *PLoS Med.* 2009; **6**(11): e1000143. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Ghansah A, Amenga-Etego L, Amambua-Ngwa A, et al.: **Monitoring parasite diversity for malaria elimination in sub-Saharan Africa.** *Science.* 2014; **345**(6202): 1297–8. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Auburn S, Campino S, Clark TG, et al.: **An effective method to purify *Plasmodium falciparum* DNA directly from clinical blood samples for whole genome high-throughput sequencing.** *PLoS One.* 2011; **6**(7): e22213. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Venkatesan M, Amaratunga C, Campino S, et al.: **Using CF11 cellulose columns to inexpensively and effectively remove human DNA from *Plasmodium falciparum*-infected whole blood samples.** *Malar J.* 2012; **11**: 41. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Manske M, Miotto O, Campino S, et al.: **Analysis of *Plasmodium falciparum* diversity in natural infections by deep sequencing.** *Nature.* 2012; **487**(7407): 375–9. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Vauterin P, Jeffery B, Miles A, et al.: **Panoptes: Web-based exploration of large scale genome variation data.** *Bioinformatics.* 2017; **33**(20): 3243–3249. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- MalariaGEN *Plasmodium falciparum* Community Project: **Genomic epidemiology of artemisinin resistant malaria.** *eLife.* 2016; **5**: e08714. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Miotto O, Almagro-Garcia J, Manske M, et al.: **Multiple populations of artemisinin-resistant *Plasmodium falciparum* in Cambodia.** *Nat Genet.* 2013; **45**(6): 648–55. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Amato R, Pearson RD, Almagro-Garcia J, et al.: **Origins of the current outbreak of multidrug-resistant malaria in southeast Asia: a retrospective genetic study.** *Lancet Infect Dis.* 2018; **18**(3): 337–45. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Amambua-Ngwa A, Amenga-Etego L, Kamau E, et al.: **Major subpopulations of *Plasmodium falciparum* in sub-Saharan Africa.** *Science.* 2019; **365**(6455): 813–6. [PubMed Abstract](#) | [Publisher Full Text](#)
- Hamilton WL, Amato R, van der Pluijm RW, et al.: **Evolution and expansion of multidrug-resistant malaria in southeast Asia: a genomic epidemiology study.** *Lancet Infect Dis.* 2019; **19**(9): 943–51. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- van der Pluijm RW, Imwong M, Chau NH, et al.: **Determinants of dihydroartemisinin-piperazine treatment failure in *Plasmodium falciparum* malaria in Cambodia, Thailand, and Vietnam: a prospective clinical, pharmacological, and genetic study.** *Lancet Infect Dis.* 2019; **19**(9): 952–61. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Ariey F, Witkowski B, Amaratunga C, et al.: **A molecular marker of artemisinin-resistant *Plasmodium falciparum* malaria.** *Nature.* 2014; **505**(7481): 50–5. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Nwakanma DC, Duffy CW, Amambua-Ngwa A, et al.: **Changes in malaria parasite drug resistance in an endemic population over a 25-year period with resulting genomic evidence of selection.** *J Infect Dis.* 2014; **209**(7): 1126–35. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Ashley EA, Dhorda M, Fairhurst RM, et al.: **Spread of Artemisinin Resistance in *Plasmodium falciparum* Malaria.** *N Engl J Med.* 2014; **371**(5): 411–23. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Kamau E, Campino S, Amenga-Etego L, et al.: **K13-propeller polymorphisms in *Plasmodium falciparum* parasites from sub-Saharan Africa.** *J Infect Dis.* 2015; **211**(8): 1352–5. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Ravenhall M, Benavente ED, Mipando M, et al.: **Characterizing the impact of sustained sulfadoxine/pyrimethamine use upon the *Plasmodium falciparum* population in Malawi.** *Malar J.* 2016; **15**(1): 575. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Gomes AR, Ravenhall M, Benavente ED, et al.: **Genetic diversity of next generation antimalarial targets: A baseline for drug resistance surveillance programmes.** *Int J Parasitol Drugs Drug Resist.* 2017; **7**(2): 174–180. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Apinjoh TO, Mugri RN, Miotto O, et al.: **Molecular markers for artemisinin and partner drug resistance in natural *Plasmodium falciparum* populations following increased insecticide treated net coverage along the slope of mount Cameroon: Cross-sectional study.** *Infect Dis Poverty.* 2017; **6**(1): 136. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Ross LS, Dhingra SK, Mok S, et al.: **Emerging Southeast Asian PfCRT mutations confer *Plasmodium falciparum* resistance to the first-line antimalarial piperazine.** *Nat Commun.* 2018; **9**(1): 3314. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Miotto O, Amato R, Ashley EA, et al.: **Genetic architecture of artemisinin-resistant *Plasmodium falciparum*.** *Nat Genet.* 2015; **47**(3): 226–34. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Takala-Harrison S, Jacob CG, Arze C, et al.: **Independent Emergence of Artemisinin Resistance Mutations Among *Plasmodium falciparum* in Southeast Asia.** *J Infect Dis.* 2015; **211**(5): 670–9. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Amato R, Lim P, Miotto O, et al.: **Genetic markers associated with dihydroartemisinin-piperazine failure in *Plasmodium falciparum* malaria in Cambodia: a genotype-phenotype association study.** *Lancet Infect Dis.* 2017; **17**(2): 164–73. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Borrmann S, Straimer J, Mwai L, et al.: **Genome-wide screen identifies new candidate genes associated with artemisinin susceptibility in *Plasmodium falciparum* in Kenya.** *Sci Rep.* 2013; **3**: 3318. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Wendler JP, Okombo J, Amato R, et al.: **A Genome Wide Association Study of *Plasmodium falciparum* Susceptibility to 22 Antimalarial Drugs in Kenya.** *PLoS One.* 2014; **9**(5): e96486. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Zhu L, Tripathi J, Rocamora FM, et al.: **The origins of malaria artemisinin resistance defined by a genetic and transcriptomic background.** *Nat Commun.* 2018; **9**(1): 5158. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Sepúlveda N, Phelan J, Diez-Benavente E, et al.: **Global analysis of *Plasmodium falciparum* histidine-rich protein-2 (*pfrp2*) and *pfrp3* gene deletions using whole-genome sequencing data and meta-analysis.** *Infect Genet Evol.* 2018; **62**: 211–9. [PubMed Abstract](#) | [Publisher Full Text](#)
- Williams AR, Douglas AD, Miura K, et al.: **Enhancing blockade of *Plasmodium falciparum* erythrocyte invasion: assessing combinations of antibodies against PFRH5 and other merozoite antigens.** *PLoS Pathog.* 2012; **8**(11): e1002991. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Benavente ED, Oresegun DR, de Sessions PF, et al.: **Global genetic diversity of *var2csa* in *Plasmodium falciparum* with implications for malaria in pregnancy and vaccine development.** *Sci Rep.* 2018; **8**(1): 15429. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

32. Amambua-Ngwa A, Tetteh KKA, Manske M, *et al.*: Population genomic scan for candidate signatures of balancing selection to guide antigen characterization in malaria parasites. *PLoS Genet.* 2012; **8**(11): e1002992. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
33. Campino S, Marin-Menendez A, Kemp A, *et al.*: A forward genetic screen reveals a primary role for *Plasmodium falciparum* Reticulocyte Binding Protein Homologue 2a and 2b in determining alternative erythrocyte invasion pathways. *PLoS Pathog.* 2018; **14**(11): e1007436. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
34. Crosnier C, Iqbal Z, Knuepfer E, *et al.*: Binding of *Plasmodium falciparum* merozoite surface proteins DBLMSP and DBLMSP2 to human immunoglobulin M is conserved among broadly diverged sequence variants. *J Biol Chem.* 2016; **291**(27): 14285–99. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
35. Amambua-Ngwa A, Jeffries D, Amato R, *et al.*: Consistent signatures of selection from genomic analysis of pairs of temporal and spatial *Plasmodium falciparum* populations from the Gambia. *Sci Rep.* 2018; **8**(1): 9687. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
36. Duffy CW, Amambua-Ngwa A, Ahouidi AD, *et al.*: Multi-population genomic analysis of malaria parasites indicates local selection and differentiation at the *gdf1* locus regulating sexual development. *Sci Rep.* 2018; **8**(1): 15763. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
37. Duffy CW, Ba H, Assefa S, *et al.*: Population genetic structure and adaptation of malaria parasites on the edge of endemic distribution. *Mol Ecol.* 2017; **26**(11): 2880–2894. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
38. Duffy CW, Assefa SA, Abugri J, *et al.*: Comparison of genomic signatures of selection on *Plasmodium falciparum* between different regions of a country with high malaria endemicity. *BMC Genomics.* 2015; **16**(1): 527. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
39. Mobegi VA, Duffy CW, Amambua-Ngwa A, *et al.*: Genome-wide analysis of selection on the malaria parasite *Plasmodium falciparum* in West African populations of differing infection endemicity. *Mol Biol Evol.* 2014; **31**(6): 1490–9. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
40. Shetty AC, Jacob CG, Huang F, *et al.*: Genomic structure and diversity of *Plasmodium falciparum* in Southeast Asia reveal recent parasite migration patterns. *Nat Commun.* 2019; **10**(1): 2665. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
41. Auburn S, Campino S, Miotto O, *et al.*: Characterization of within-host *Plasmodium falciparum* diversity using next-generation sequence data. *PLoS One.* 2012; **7**(2): e32891. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
42. Assefa SA, Preston MD, Campino S, *et al.*: estMOI: estimating multiplicity of infection using parasite deep sequencing data. *Bioinformatics.* 2014; **30**(9): 1292–4. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
43. Murray L, Mobegi VA, Duffy CW, *et al.*: Microsatellite genotyping and genome-wide single nucleotide polymorphism-based indices of *Plasmodium falciparum* diversity within clinical infections. *Malar J.* 2016; **15**(1): 275. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
44. Chang HH, Worby CJ, Yeka A, *et al.*: THE REAL McCOIL: A method for the concurrent estimation of the complexity of infection and SNP allele frequency for malaria parasites. *PLoS Comput Biol.* 2017; **13**(1): e1005348. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
45. O'Brien JD, Iqbal Z, Wendler J, *et al.*: Inferring Strain Mixture within Clinical *Plasmodium falciparum* Isolates from Genomic Sequence Data. *PLoS Comput Biol.* 2016; **12**(6): e1004824. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
46. Robinson T, Campino SG, Auburn S, *et al.*: Drug-resistant genotypes and multi-clonality in *Plasmodium falciparum* analysed by direct genome sequencing from peripheral blood of malaria patients. *in press. PLoS One.* 2011; **6**(8): e23204. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
47. O'Brien JD, Amenga-Etego L, Li R: Approaches to estimating inbreeding coefficients in clinical isolates of *Plasmodium falciparum* from genomic sequence data. *Malar J.* 2019; **18**: 473. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
48. Zhu SJ, Almagro-Garcia J, McVean G: Deconvolution of multiple infections in *Plasmodium falciparum* from high throughput sequencing data. *Bioinformatics.* 2018; **34**(1): 9–15. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
49. Zhu SJ, Hendry JA, Almagro-Garcia J, *et al.*: The origins and relatedness structure of mixed infections vary with local prevalence of *P. falciparum* malaria. *eLife.* 2019; **8**: e40845. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
50. Henden L, Lee S, Mueller I, *et al.*: Identity-by-descent analyses for measuring population dynamics and selection in recombining pathogens. *PLoS Genet.* 2018; **14**(5): e1007279. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
51. Schaffner SF, Taylor AR, Wong W, *et al.*: hmmIBD: software to infer pairwise identity by descent between haploid genotypes. *Malar J.* 2018; **17**(1): 196. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
52. Samad H, Coll F, Preston MD, *et al.*: Imputation-Based Population Genetics Analysis of *Plasmodium falciparum* Malaria Parasites. *PLoS Genet.* 2015; **11**(4): e1005131. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
53. Ravenhall M, Campino S, Clark TG: SV-Pop: population-based structural variant analysis and visualization. *BMC Bioinformatics.* 2019; **20**(1): 136. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
54. Jacob CG, Tan JC, Miller BA, *et al.*: A microarray platform and novel SNP calling algorithm to evaluate *Plasmodium falciparum* field samples of low DNA quantity. *BMC Genomics.* 2014; **15**(1): 719. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
55. Preston MD, Assefa SA, Ocholla H, *et al.*: PlasmoView: A Web-based Resource to Visualise Global *Plasmodium falciparum* Genomic Variation. *J Infect Dis.* 2014; **209**(11): 1808–15. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
56. Miles A, Iqbal Z, Vauterin P, *et al.*: Indels, structural variation, and recombination drive genomic diversity in *Plasmodium falciparum*. *Genome Res.* 2016; **26**(9): 1288–99. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
57. Hamilton WL, Claessens A, Otto TD, *et al.*: Extreme mutation bias and high AT content in *Plasmodium falciparum*. *Nucleic Acids Res.* 2017; **45**(4): 1889–901. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
58. Carvalho CMB, Ramocki MB, Pehlivan D, *et al.*: Inverted genomic segments and complex triplication rearrangements are mediated by inverted repeats in the human genome. *Nat Genet.* 2011; **43**(11): 1074–81. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
59. Molina-Cruz A, Garver LS, Alabaster A, *et al.*: The human malaria parasite *Pfs47* gene mediates evasion of the mosquito immune system. *Science.* 2013; **340**(6135): 984–7. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
60. Gardiner DL, Dixon MWA, Spielmann T, *et al.*: Implication of a *Plasmodium falciparum* gene in the switch between asexual reproduction and gametocytogenesis. *Mol Biochem Parasitol.* 2005; **140**(2): 153–60. [PubMed Abstract](#) | [Publisher Full Text](#)
61. Moelans II, Meis JF, Kocken C, *et al.*: A novel protein antigen of the malaria parasite *Plasmodium falciparum*, located on the surface of gametes and sporozoites. *Mol Biochem Parasitol.* 1991; **45**(2): 193–204. [PubMed Abstract](#) | [Publisher Full Text](#)
62. Dessens JT, Beetsma AL, Dimopoulos G, *et al.*: CTRP is essential for mosquito infection by malaria ookinetes. *EMBO J.* 1999; **18**(22): 6221–7. [PubMed Abstract](#) | [Free Full Text](#)
63. Laufer MK, Thesing PC, Eddington ND, *et al.*: Return of Chloroquine Antimalarial Efficacy in Malawi. *N Engl J Med.* 2006; **355**(19): 1959–66. [PubMed Abstract](#) | [Publisher Full Text](#)
64. Laufer MK, Takala-Harrison S, Dzinjalimala FK, *et al.*: Return of Chloroquine-Susceptible *Plasmodium falciparum* in Malawi Was a Reexpansion of Diverse Susceptible Parasites. *J Infect Dis.* 2010; **202**(5): 801–8. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
65. Frosch AEP, Laufer MK, Mathanga DP, *et al.*: Return of Widespread Chloroquine-Sensitive *Plasmodium falciparum* to Malawi. *J Infect Dis.* 2014; **210**(7): 1110–4. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
66. Wootton JC, Feng X, Ferdig MT, *et al.*: Genetic diversity and chloroquine selective sweeps in *Plasmodium falciparum*. *Nature.* 2002; **418**(6895): 320–3. [PubMed Abstract](#) | [Publisher Full Text](#)
67. Mita T, Tanabe K, Kita K: Spread and evolution of *Plasmodium falciparum* drug resistance. *Elsevier, Parasitol Int.* 2009; **58**(3): 201–9. [PubMed Abstract](#) | [Publisher Full Text](#)
68. Agrawal S, Moser KA, Morton L, *et al.*: Association of a Novel Mutation in the *Plasmodium falciparum* Chloroquine Resistance Transporter With Decreased Piperazine Sensitivity. *J Infect Dis.* 2017; **216**(4): 468–76. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
69. Naidoo I, Roper C: Mapping 'partially resistant', 'fully resistant', and 'super resistant' malaria. *Trends Parasitol.* 2013; **29**(10): 505–15. [PubMed Abstract](#) | [Publisher Full Text](#)
70. Heinberg A, Kirkman L: The molecular basis of antifolate resistance in *Plasmodium falciparum*: looking beyond point mutations. *Ann N Y Acad Sci.* 2015; **1342**(1): 10–8. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
71. World Health Organization: Artemisinin and artemisinin-based combination therapy resistance: status report. 2018. [Reference Source](#)
72. Price RN, Uhlemann AC, Brockman A, *et al.*: Mefloquine resistance in *Plasmodium falciparum* and increased *pfmdr1* gene copy number. *Lancet.* 2004; **364**(9432): 438–47. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
73. Cheng Q, Gatton ML, Barnwell J, *et al.*: *Plasmodium falciparum* parasites lacking histidine-rich protein 2 and 3: a review and recommendations for accurate reporting. *Malar J.* 2014; **13**: 283. [PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

74. WHO: Malaria rapid diagnostic test performance. Results of WHO product testing of malaria RDTs: round 8 (2016-2018). WHO, 2018; (accessed Aug 22, 2019).  
[Reference Source](#)
75. Gamboa D, Ho MF, Bendezu J, *et al.*: A Large Proportion of *P. falciparum* Isolates in the Amazon Region of Peru Lack *pfhrp2* and *pfhrp3*: Implications for Malaria Rapid Diagnostic Tests. *PLoS One*. 2010; 5(1): e8091.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
76. Rachid Viana GM, Akinyi Okoth S, Silva-Flannery L, *et al.*: Histidine-rich protein 2 (*pfhrp2*) and *pfhrp3* gene deletions in *Plasmodium falciparum* isolates from select sites in Brazil and Bolivia. *PLoS One*. 2017; 12(3): e0171150.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
77. Parr JB, Verity R, Doctor SM, *et al.*: *Pfhrp2*-deleted *Plasmodium falciparum* parasites in the democratic republic of the congo: a national cross-sectional survey. *J Infect Dis*. 2017; 216(1): 36–44.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
78. Menegon M, L'Episcopia M, Nurahmed AM, *et al.*: Identification of *Plasmodium falciparum* isolates lacking histidine-rich protein 2 and 3 in Eritrea. *Infect Genet Evol*. 2017; 55: 131–4.  
[PubMed Abstract](#) | [Publisher Full Text](#)
79. Bharti PK, Chandel HS, Ahmad A, *et al.*: Prevalence of *pfhrp2* and/or *pfhrp3* Gene Deletion in *Plasmodium falciparum* Population in Eight Highly Endemic States in India. *PLoS One*. 2016; 11(8): e0157949.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
80. Baker J, Ho MF, Pelecanos A, *et al.*: Global sequence variation in the histidine-rich proteins 2 and 3 of *Plasmodium falciparum*: implications for the performance of malaria rapid diagnostic tests. *Malar J*. 2010; 9: 129.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
81. Akinyi S, Hayden T, Gamboa D, *et al.*: Multiple genetic origins of histidine-rich protein 2 gene deletion in *Plasmodium falciparum* parasites from Peru. *Sci Rep*. 2013; 3: 2797.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
82. Akinyi Okoth S, Abdallah JF, Ceron N, *et al.*: Variation in *Plasmodium falciparum* Histidine-Rich Protein 2 (*Pfhrp2*) and *Plasmodium falciparum* Histidine-Rich Protein 3 (*Pfhrp3*) Gene Deletions in Guyana and Suriname. *PLoS One*. 2015; 10(5): e0126805.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
83. Parr JB, Anderson O, Juliano JJ, *et al.*: Streamlined, PCR-based testing for *pfhrp2*- and *pfhrp3*-negative *Plasmodium falciparum*. *Malar J*. 2018; 17(1): 137.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
84. World Health Organisation: WHO Strategic Advisory Group on Malaria Eradication. Malaria eradication: benefits, future scenarios and feasibility. Executive Summary. WHO Strategic Advisory Group on Malaria Eradication. Executive Summary. Geneva: World Health Organisation, 2019.  
[Reference Source](#)
85. Dalmat R, Naughton B, Kwan-Gett TS, *et al.*: Use cases for genetic epidemiology in malaria elimination. *Malar J*. 2019; 18(1): 163.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
86. Early AM, Daniels RF, Farrell TM, *et al.*: Detection of low-density *Plasmodium falciparum* infections using amplicon deep sequencing. *Malar J*. 2019; 18(1): 219.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
87. Boyce RM, Hathaway N, Fulton T, *et al.*: Reuse of malaria rapid diagnostic tests for amplicon deep sequencing to estimate *Plasmodium falciparum* transmission intensity in western Uganda. *Sci Rep*. 2018; 8(1): 10159.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
88. Li H, Durbin R: Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*. 2009; 25(14): 1754–60.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
89. DePristo MA, Banks E, Poplin R, *et al.*: A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet*. 2011; 43(5): 491–8.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
90. Cingolani P, Platts A, Wang LL, *et al.*: A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)*. 2012; 6(2): 80–92.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
91. Echeverry DF, Deason NA, Davidson J, *et al.*: Human malaria diagnosis using a single-step direct-PCR based on the *Plasmodium cytochrome oxidase III* gene. *Malar J*. 2016; 15: 128.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
92. MalariaGEN: Supplementary data to: An open dataset of *Plasmodium falciparum* genome variation in 7,000 worldwide samples. *figshare*. Dataset. 2021.  
<http://www.doi.org/10.6084/m9.figshare.13388603.v1>

# Open Peer Review

Current Peer Review Status:  

---

## Version 1

Reviewer Report 25 March 2021

<https://doi.org/10.21956/wellcomeopenres.17752.r42796>

© 2021 Menard D. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



**Didier Menard** 

Malaria Genetics and Resistance Unit, Parasites and Insect Vectors Department, Institut Pasteur, Paris, France

This manuscript from the MalariaGEN consortium, a data-sharing community of teams working on *Plasmodium falciparum* genomic epidemiology, presents the new release of curated *P. falciparum* genomes from isolates collected in 73 locations in Africa, Asia, South America and Oceania.

Based on robust and perfectly detailed methods (ranging from the treatment of the blood samples, the DNA extraction, the Illumina and computational platforms developed to produce genome sequencing for variant discovery and genotype calling), they analyzed 7000 *P. falciparum* genome sequences and provided numerous exciting data. For instance, they found that variations (SNPs and indels) in *P. falciparum* genome affected about a quarter of the 23 Mb genome (and mostly coding regions), or that duplication genotypes are frequent around *mdr1*, *plasmepsin2/3* and *gch1*, which are known to be associated with antimalarial drug resistance (including mefloquine, piperaquine and sulfadoxine/pyrimethamine).

Moreover, population genetic analyses conducted on this largest available data resource, depict a comprehensive picture of *P. falciparum* parasite populations globally and sub populations at continental level. In the results, a large section is devoted to the description of the geographic patterns of validated molecular markers (SNPs and CNVs) associated with antimalarial drug resistance. By compiling data on all samples collected from 2002–2015, they present clear profiles of drug resistance by regional sub-populations for the most used antimalarial drugs. Finally, they reveal a global landscape regarding a major challenge for malaria elimination, that are deletions in *hrp2* and *3* genes linked with false negative results of HRP2-based malaria RDT.

Written in a very clear way, it must be point out that the authors have made huge efforts so that these data are understandable for a general audience, especially for the non-experts in genomics or for policy makers in malaria endemic countries. Their data effectively depict the main challenges currently encountered in the fight against malaria: the monitoring of the strategies deployed by the assessment of the impact on *P. falciparum* parasite populations, the geographical evolution of antimalarial drug resistances and the effectiveness of diagnostic tools used in malaria

endemic areas (i.e. malaria RDT).

Of note, the authors fairly expose the main issues and drawbacks related to the methods used (i.e. the analytical challenges due to long tracts of highly repetitive sequence and hypervariable regions within the *P. falciparum* genome, and the challenges of studying a complex mixture of genotypes from polyclonal infections),

Although, I am impressed by the work done by the consortium, I have several minor comments that could improve the manuscript:

- Sample collection - *P. falciparum* samples investigated are not from systematic sampling collections dedicated to this study but rather from multiple studies conducted by groups with different objectives and from heterogeneous populations (patients living in malaria endemic areas, travelers, etc.) . I think this issue should be discussed in the manuscript.
- Likewise, the long time period covering the samples collection (2002–2015) is also a major bias which can alter the final results.
- I guess that all samples were collected from symptomatic patients seen at health facilities level? Unfortunately, this makes that data presented capture only *P. falciparum* populations infected this population. With the rise of new technologies, I am wondering whether the MalariaGEN consortium could investigate samples collected from asymptomatic individuals and explore the genomic profiles of this hidden reservoir but representing the major parasite biomass?
- I am aware that the authors have performed a difficult and complex exercise by providing high quality genomic data and comprehensive description of their data for a large audience. The major challenge that is not addressed in the manuscript is how these important data can be translated into concrete actions in the field by health providers.
- Last comment regarding the database. It will be helpful to provide for each sample/genome sequence, the location (country) and the date of collection.

**Is the work clearly and accurately presented and does it cite the current literature?**

Yes

**Is the study design appropriate and is the work technically sound?**

Yes

**Are sufficient details of methods and analysis provided to allow replication by others?**

Yes

**If applicable, is the statistical analysis and its interpretation appropriate?**

I cannot comment. A qualified statistician is required.

**Are all the source data underlying the results available to ensure full reproducibility?**

Yes

**Are the conclusions drawn adequately supported by the results?**



Yes

**Competing Interests:** No competing interests were disclosed.

**Reviewer Expertise:** Expert in antimalarial drug resistance

**I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

Author Response 06 Jul 2021

**Richard Pearson**, Wellcome Sanger Institute, Hinxton, UK

We thank the reviewers for the extremely positive and supportive feedback. In their comments and suggestions both reviewers have well captured the spirit of this data resource and of the large collaborative network behind it. We are pleased to submit detailed responses and a revised version of the manuscript that addresses their comments.

**2.1) Sample collection - *P. falciparum* samples investigated are not from systematic sampling collections dedicated to this study but rather from multiple studies conducted by groups with different objectives and from heterogeneous populations (patients living in malaria endemic areas, travelers, etc.). I think this issue should be discussed in the manuscript.**

Thanks for raising this point. On one hand, the heterogeneity of sampling approaches offers a unique opportunity to investigate questions in a variety of epidemiological settings in a systematic way. Specifics of each study are provided in [ftp://ngs.sanger.ac.uk/production/malaria/pfcommunityproject/Pf6/Pf\\_6\\_partner\\_studies.pdf](ftp://ngs.sanger.ac.uk/production/malaria/pfcommunityproject/Pf6/Pf_6_partner_studies.pdf) and users of the resource can contact individual investigators for further details. At the same time, we agree that this can also act as a confounder in some analysis, which is why we've devoted significant time to the curation of the dataset to make it "analysis ready".

As suggested, we have amended the manuscript in version 2 to include the considerations above in the paragraph: "Samples were collected by independent groups that were operative at a given time and in a given place with distinct objectives; while care needs to be taken when interpreting results spanning multiple years and geographical settings (e.g. aggregated trends of drug resistance prevalence), this heterogeneity also allows for the exploration of a wide range of epidemiological and transmission settings."

**2.2) Likewise, the long time period covering the samples collection (2002–2015) is also a major bias which can alter the final results.**

This is an important point in particular for interpreting drug resistance results, and one we explicitly bring out in the paragraph: "Note that samples were collected over a relatively long time period (2002–15) during which there were major changes in global patterns of drug resistance, and that the sampling locations represented in a given year depended on which partner studies were operative at the time. To alleviate this problem, we have also divided the data into samples collected before and after 2011 (Supplementary Data);

Supplementary table 10), but temporal trends in aggregated data should be interpreted with due caution.”. Following the reviewer’s suggestion, we have now stressed this point further in our reply to point (2.1) above.

**2.3) I guess that all samples were collected from symptomatic patients seen at health facilities level? Unfortunately, this makes that data presented capture only *P. falciparum* populations infected this population. With the rise of new technologies, I am wondering whether the MalariaGEN consortium could investigate samples collected from asymptomatic individuals and explore the genomic profiles of this hidden reservoir but representing the major parasite biomass?**

Asymptomatic infections are indeed an incredibly significant reservoir that needs to be explicitly considered to achieve a complete and accurate picture of the transmission landscape. The development of new technologies has begun to dig deeper and deeper in this area and initial results seem to be very encouraging that good quality data can indeed be obtained from asymptomatic and/or low parasitemia subjects. MalariaGEN would certainly be supportive of this kind of effort and we have indeed active collaborations with partners exploring these questions. To the best of our knowledge, though, some of these methodologies are still of limited sensitivity and in part experimental and will require further work in order to be deployed on the large scale required by this scientific question, but that is certainly an area for future investigation.

**2.4) I am aware that the authors have performed a difficult and complex exercise by providing high quality genomic data and comprehensive description of their data for a large audience. The major challenge that is not addressed in the manuscript is how these important data can be translated into concrete actions in the field by health providers.**

This data resource represents a clear step towards the ultimate objective of translating genomic surveillance outputs into actionable actions, although it is fair to say that this is a long journey with many different components. The ability for multiple groups to share data, to analyse it using standardised methods, and to make it readily accessible is the foundation for translational impact to reach maturity.

In the discussion we highlighted a series of future translational directions which have been and will be facilitated by resources like this one (and future ones) but it is certainly true that these results require careful interpretation due to the caveats highlighted in the paper and by the reviewer, which inevitably limit their impact. At the same time this dataset does create a systematic framework to enact and contextualize future discoveries of that nature and, indirectly, contributes to them.

Ultimately, the practical value for malaria control will be greatly enhanced by the progressive acquisition of longitudinal time-series data and their integration with other sources of epidemiological data which will allow control programmes to monitor the impact of their interventions on the parasite population in near real time.

**2.5) Last comment regarding the database. It will be helpful to provide for each sample/genome sequence, the location (country) and the date of collection.**

This information is included in the "Sample provenance and sequencing metadata" file available at the resource page <https://www.malariagen.net/resource/26>

**Competing Interests:** No competing interests were disclosed.

Reviewer Report 22 March 2021

<https://doi.org/10.21956/wellcomeopenres.17752.r42794>

© 2021 Veiga M et al. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



**Maria Isabel Veiga** 

ICVS/3B's - PT Government Associate Laboratory, University of Minho, Braga, Portugal

**Nuno S. Osório** 

ICVS/3B's - PT Government Associate Laboratory, University of Minho, Braga, Portugal

The analysis of whole-genome sequences obtained from *Plasmodium falciparum* is particularly challenging due to the presence of hypervariable regions, highly repetitive sequences, and frequent mixture of parasites due to multiple infections of the host. The authors of this study describe a curated list of over three million high-confidence polymorphisms obtained from the genome sequence analysis of more than 7000 samples of *P. falciparum* collected by several studies in 73 locations in Africa, Asia, South America and Oceania.

This work, reporting a laudable effort to substantially enrich publicly available genome data of *P. falciparum* worldwide, is of paramount importance for the field. The contribution goes in line with authors' previous consortia publications, extending largely the number of available data that can be analysed via web with powerful data analysis pipelines. By providing open access to a curated list of polymorphisms based on reproducible and high-quality protocols for the sequencing and analysis of *P. falciparum* genomes this study is likely to decrease the difficulties that have delayed the research on genomic epidemiology and population genomics of *P. falciparum*. Among other advances, studies in this area are likely to have important implications for a better understanding of the evolution towards drug resistance of the different global parasite populations ultimately contributing for a better control of this devastating disease. The manuscript is very well written and clear. It presents eight genetically distinct populations of parasites each endemic to different world regions, including South America, West Africa, Central Africa, East Africa, South Asia, West Southeast Asia, East Southeast Asia and Oceania. An interesting genetic and geographic characterization of the eight parasite populations is also shown. Of note, the finding of higher within-host diversity in the parasite populations endemic to Africa, the identification of single nucleotide polymorphism with high levels of geographic differentiation, and further characterization of geographic patterns of drug resistance and polymorphisms with potential impact in rapid diagnostic tests. We do not have major criticisms of the study.

**Our minor suggestions for the improval of the manuscript focus on:**

- Increasing the accessibility of the table listing polymorphisms in supplementary data. The authors do provide the data in VCF and zarr files, which are not very user friendly nor allow a fast search of a specific polymorphism. We understand that developing a web interface for this purpose would be a challenge beyond this research article but possibly exporting the VCF file data into tables that could be available in online repositories.
- Add to the supplementary file 4, describing the drug resistance markers genotype, the PfMDR1 N86Y. This SNP is a well-known modulator of antimalarial response and considered a risk factor for the treatment of artemether-lumefantrine.
- Add the ID of the genes most mentioned in the main article. The gene ID (PF3D7\_xxxxxxx), is provided in supplementary file 7, but to clarify the reader, we recommend to add it also in the main article when first describing the genes.
- In the results section, when describing gene amplification and different sets of breakpoints, the authors describe complex rearrangements that have not been observed before in Plasmodium species. In regards to pfmdr1 duplication events has been described to vary in size while spanning different genes in different parasites<sup>1,2,3,4</sup>. In a genome walking like approach, it has been described different amplicon sizes containing the pfmdr1 in clinical isolates from Southeast Asia where they also investigated if the type (i.e., which genes are included) and size of the amplicon influence drug susceptibility phenotypes<sup>5</sup>.

**References**

1. Foote S, Thompson J, Cowman A, Kemp D: Amplification of the multidrug resistance gene in some chloroquine-resistant isolates of *P. falciparum*. *Cell*. 1989; **57** (6): 921-930 [Publisher Full Text](#)
2. Nair S, Nash D, Sudimack D, Jaidee A, et al.: Recurrent gene amplification and soft selective sweeps during evolution of multidrug resistance in malaria parasites. *Mol Biol Evol*. 2007; **24** (2): 562-73 [PubMed Abstract](#) | [Publisher Full Text](#)
3. Triglia T, Foote SJ, Kemp DJ, Cowman AF: Amplification of the multidrug resistance gene pfmdr1 in *Plasmodium falciparum* has arisen as multiple independent events. *Mol Cell Biol*. 1991; **11** (10): 5244-50 [PubMed Abstract](#) | [Publisher Full Text](#)
4. Ribacke U, Mok BW, Wirta V, Normark J, et al.: Genome wide gene amplifications and deletions in *Plasmodium falciparum*. *Mol Biochem Parasitol*. 2007; **155** (1): 33-44 [PubMed Abstract](#) | [Publisher Full Text](#)
5. Veiga MI, Ferreira PE, Malmberg M, Jörnhammar L, et al.: pfmdr1 amplification is related to increased *Plasmodium falciparum* in vitro sensitivity to the bisquinoline piperazine. *Antimicrob Agents Chemother*. 2012; **56** (7): 3615-9 [PubMed Abstract](#) | [Publisher Full Text](#)

**Is the work clearly and accurately presented and does it cite the current literature?**

Yes

**Is the study design appropriate and is the work technically sound?**

Yes

**Are sufficient details of methods and analysis provided to allow replication by others?**

Yes

**If applicable, is the statistical analysis and its interpretation appropriate?**

Yes

**Are all the source data underlying the results available to ensure full reproducibility?**

Yes

**Are the conclusions drawn adequately supported by the results?**

Yes

**Competing Interests:** No competing interests were disclosed.

**Reviewer Expertise:** Molecular epidemiology, antimalarial drug resistance

**We confirm that we have read this submission and believe that we have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

Author Response 06 Jul 2021

**Richard Pearson**, Wellcome Sanger Institute, Hinxton, UK

We thank the reviewers for the extremely positive and supportive feedback. In their comments and suggestions both reviewers have well captured the spirit of this data resource and of the large collaborative network behind it. We are pleased to submit detailed responses and a revised version of the manuscript that addresses their comments.

***1.1) Increasing the accessibility of the table listing polymorphisms in supplementary data. The authors do provide the data in VCF and zarr files, which are not very user friendly nor allow a fast search of a specific polymorphism. We understand that developing a web interface for this purpose would be a challenge beyond this research article but possibly exporting the VCF file data into tables that could be available in online repositories.***

We thank the reviewer for this important feedback on how to increase the reach of this resource. Since the publication of this article, we have been working on an initial web interface that allows users to navigate some aspects of the data: please see <https://www.malariagen.net/apps/pf6>. The current version mainly focuses on epidemiologically relevant data and emphasises the community behind the project and at the moment doesn't provide access to the genomic variation information, which will require further work.

Of course accessibility is a relative criteria and as such it requires balancing out different priorities. In the past we have provided tabular versions of the data ([www.malariagen.net/data](http://www.malariagen.net/data)) but the benefits have been very limited. For example, handling multiallelic and non-SNP variations requires somewhat arbitrary encoding decisions that significantly affect the simplicity and intuitiveness of the tabular format. Increasing the sample size has made these variations more common (e.g. in this release there are about 50% non-SNP variants and 50% multiallelic variants) to the point that there was no real

advantage in maintaining the format. The decision of primarily utilising the VCF format comes from the recognition that these files are the standard de facto in the genomic community, which in turn has developed a large ecosystem of tools to handle them: please see the README at [ftp://ngs.sanger.ac.uk/production/malaria/pfcommunityproject/Pf6/Pf\\_6\\_README\\_20191010.txt](ftp://ngs.sanger.ac.uk/production/malaria/pfcommunityproject/Pf6/Pf_6_README_20191010.txt) for some examples, e.g. to subset the data.

However we agree this might still be limiting for some use cases and we are working towards a more integrated solution. As an example of our direction of travel, please see <https://malariagen.github.io/vector-data/landing-page.html>, which presents some simplified data access workflows for the MalariaGEN *Anopheles gambiae* 1000 Genomes Project.

**1.2) Add to the supplementary file 4, describing the drug resistance markers genotype, the PfMDR1 N86Y. This SNP is a well-known modulator of antimalarial response and considered a risk factor for the treatment of artemether-lumefantrine.**

We recognise that there is growing evidence of the role of *PfMDR1* N86Y in artemether-lumefantrine resistance. In particular, multiple studies have shown that lumefantrine appears to select for N86. Despite that, WHO still reports markers of resistance to lumefantrine as “Yet to be validated” (p. 22 - <https://www.who.int/publications/i/item/9789240012813>). In this release, supplementary file 4 only contains validated markers so it would be inconsistent to add the markers. However, we will consider adding putative markers in future releases where appropriate.

**1.3) Add the ID of the genes most mentioned in the main article. The gene ID (PF3D7\_XXXXXX), is provided in supplementary file 7, but to clarify the reader, we recommend to add it also in the main article when first describing the genes.**

We have implemented the recommendation and added gene IDs every time a gene is mentioned for the first time in the manuscript version 2.

**1.4) In the results section, when describing gene amplification and different sets of breakpoints, the authors describe complex rearrangements that have not been observed before in Plasmodium species. In regards to pfmdr1 duplication events has been described to vary in size while spanning different genes in different parasites<sup>1,2,3,4</sup>. In a genome walking like approach, it has been described different amplicon sizes containing the pfmdr1 in clinical isolates from Southeast Asia where they also investigated if the type (i.e., which genes are included) and size of the amplicon influence drug susceptibility phenotypes<sup>5</sup>.**

The complex rearrangements that have not been observed before which we were referring to here are “dup-trpinv-dup” rearrangements that to the best of our knowledge have only previously been described in human data (see ref 58). This complex and large structural rearrangement involves a triplicated segment embedded within a duplication, in which the triplicated segment is inverted. We recognise that the original wording in the text was

ambiguous and we've replaced "complex rearrangements" with an explicit description of the event.

**Competing Interests:** No competing interests were disclosed.