

# The Malaria-Protective Human Glycophorin Structural Variant DUP4 Shows Somatic Mosaicism and Association with Hemoglobin Levels

Walid Algady,<sup>1,9</sup> Sandra Louzada,<sup>2,9</sup> Danielle Carpenter,<sup>1</sup> Paulina Brajer,<sup>1</sup> Anna Färnert,<sup>3,4</sup> Ingegerd Rooth,<sup>5</sup> Billy Ngasala,<sup>6,7</sup> Fengtang Yang,<sup>2</sup> Marie-Anne Shaw,<sup>8</sup> and Edward J. Hollox<sup>1,\*</sup>

Glycophorin A and glycophorin B are red blood cell surface proteins and are both receptors for the parasite *Plasmodium falciparum*, which is the principal cause of malaria in sub-Saharan Africa. DUP4 is a complex structural genomic variant that carries extra copies of a glycophorin A-glycophorin B fusion gene and has a dramatic effect on malaria risk by reducing the risk of severe malaria by up to 40%. Using fiber-FISH and Illumina sequencing, we validate the structural arrangement of the glycophorin locus in the DUP4 variant and reveal somatic variation in copy number of the glycophorin B-glycophorin A fusion gene. By developing a simple, specific, PCR-based assay for DUP4, we show that the DUP4 variant reaches a frequency of 13% in the population of a malaria-endemic village in south-eastern Tanzania. We genotype a substantial proportion of that village and demonstrate an association of DUP4 genotype with hemoglobin levels, a phenotype related to malaria, using a family-based association test. Taken together, we show that DUP4 is a complex structural variant that may be susceptible to somatic variation and show that DUP4 is associated with a malarial-related phenotype in a longitudinally followed population.

## Introduction

Structural variation (SV) of genomes, including inversions, deletions, duplications, and more complex rearrangements is seen at polymorphic frequencies across all species. Like single-nucleotide variation, much SV is likely to be evolving neutrally but in some cases there is evidence for balancing or adaptive evolution.<sup>1–5</sup> SV has also been shown to generate novel genes with functional consequences, for example generation of human-specific *SRGAP2* genes that increase the density of dendritic spines in the brain.<sup>6,7</sup> Regions that show extensive SV are thought to have a high mutation rate, due to recurrent non-allelic homologous recombination.<sup>8–11</sup> In addition to SV in the germline, large somatic SVs have been observed in human brain,<sup>12</sup> in skin fibroblasts,<sup>13</sup> and in blood of identical twins.<sup>14</sup>

Although SV is a source of variation between individuals<sup>15</sup> and between cells within an individual,<sup>16</sup> its contribution to disease resistance remains unclear, and the best-characterized examples of structural variants associated with a common disease involve identifying an inherited germline structural variant in linkage disequilibrium with a sentinel single-nucleotide variant (SNV) that has been previously highlighted in a large genome-wide association study.<sup>17</sup>

The human genome assembly carries three glycophorin genes, *GYPE* (MIM: 138590), *GYPB* (MIM: 617923), and

*GYP A* (MIM: 617922), tandemly arranged on three ~120 kb repeats sharing ~97% identity. Glycophorin A (encoded by *GYP A*) and glycophorin B (encoded by *GYP B*) are readily detectable on the surface of erythrocytes and carry the MNS blood group system.<sup>18</sup> Mature glycophorin E (encoded by *GYPE*) is predicted to be 59 amino acids long but has not been unambiguously detected on the erythrocyte surface. The genomic region carrying these genes is known to undergo extensive copy number variation and gene conversion, and rearrangements that shuffle the coding regions of the three genes can generate rare blood group antigens in the Miltenberger series (MIM: 111300).<sup>19</sup> This genomic region has also been highlighted as a region of balancing selection<sup>20–23</sup> and positive selection,<sup>24</sup> though the effect of extensive gene conversion between the 120 kb repeats on statistical measures of positive selection is not clear.

Recent work studying the host genetic contribution to severe malaria identified a SNV allele at a nearby non-repeated region in linkage disequilibrium with a complex structural variant at the glycophorin locus protective against severe malaria.<sup>5,20</sup> This structural variant, called DUP4, is restricted to East African populations and is responsible for a glycophorin B-glycophorin A fusion gene product that is detected using serology as the blood group antigen Dantu NE+. This DUP4 variant confers a clinically important protective effect, with carriers ~40% less likely to develop cerebral malaria.<sup>5</sup> Given that

<sup>1</sup>Department of Genetics and Genome Biology, University of Leicester, Leicester LE1 7RH, UK; <sup>2</sup>Wellcome Sanger Institute, Hinxton, Cambridge CB10 1SA, UK; <sup>3</sup>Division of Infectious Diseases, Department of Medicine Solna, Karolinska Institutet, 17176 Stockholm, Sweden; <sup>4</sup>Department of Infectious Diseases, Karolinska University Hospital, Stockholm 17176, Sweden; <sup>5</sup>Nyamiasati Malaria Research, Rufiji, National Institute for Medical Research, Dar-es-Salaam, Tanzania; <sup>6</sup>Department of Parasitology and Medical Entomology, Muhimbili University of Health and Allied Sciences, Dar es Salaam, Tanzania; <sup>7</sup>Department of Women's and Children's Health, International Maternal and Child Health (IMCH), Uppsala Universitet, 75185 Uppsala, Sweden; <sup>8</sup>Leeds Institute of Medical Research at St James's, University of Leeds, Leeds LS9 7TF, UK

<sup>9</sup>These authors contributed equally to this work

\*Correspondence: [ejh33@le.ac.uk](mailto:ejh33@le.ac.uk)

<https://doi.org/10.1016/j.ajhg.2018.10.008>

© 2018 The Author(s). This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).



glycophorin A and glycophorin B are erythrocyte receptors for the *Plasmodium falciparum* cell surface receptors MSP-1, EBA-175, and EBL-1, respectively, it is likely that this protective effect is mediated by an altered ability of the parasite to invade the host erythrocyte.<sup>25–28</sup>

A model for the DUP4 variant, based on analysis of mis-mapping positions of short-read sequences, was put forward that involved a duplication of the *GYPE* gene, deletion of the *GYPB* gene, and generation of two copies of a *GYPB-A* fusion gene.<sup>5</sup> It is not clear exactly how the *GYPB-A* fusion protein confers protection against malaria, but it has been suggested that it could affect interactions with *Plasmodium falciparum* receptors and host band 3 protein at the erythrocyte surface.<sup>5</sup> A complete characterization of the DUP4 allele is therefore important to understand the mechanistic basis for this protective effect and to facilitate design of treatments mimicking the mechanism of protection.

## Subjects and Methods

### Study Population, Samples, and Phenotypic Data

The study population is from the coastal village of Nyamisati in the Rufiji delta in Tanzania.<sup>29</sup> Human genomic DNA was extracted from whole peripheral blood, as previously described,<sup>30</sup> with informed consent from participants or their guardians and approval of the local ethics committees of the Muhimbili University of Health and Allied Sciences and National Institute of Medical Research in Tanzania and the Regional Ethical Committee of Stockholm in Sweden.

Phenotypic data spanning 7 years, from 1993 to 1999, was collected using annual total population surveys and annual records of malarial episodes, as previously described.<sup>30</sup> The total population survey provided information on asymptomatic parasite load (parasites per  $\mu\text{L}$ ) and hemoglobin levels. A single hemoglobin value was generated for each individual from annual total population surveys carried out over 7 years, corrected for age, sex, and parasite load prior to genetic analysis. A single parasite load value was generated for each individual from asymptomatic parasitemia recorded during annual total population surveys carried out over 7 years. This single value was corrected for age and sex prior to genetic analysis.

All clinical malarial episodes were recorded and confirmed by microscopy. Multiple clinical episodes were recorded if the recurrence was greater than 4 weeks apart. A small proportion (1%) of individuals presented with a clinical episode at the time of the total population survey, and these samples were included in the analysis. A single clinical episode value was generated for each individual from records of all malarial episodes occurring in the village during a period of 7 years. The phenotypic data derived for clinical episodes, parasite load, and hemoglobin have previously been described in detail.<sup>30,31</sup>

### Fluorescence *In Situ* Hybridization using Single-Molecule DNA Fibers (Fiber-FISH)

The probes used in this study included four fosmid clones selected from the UCSC Genome Browser GRCh37/hg19 assembly and a 3,632-bp PCR product that is specific for the glycophorin E repeat (see below). Probes were made by whole-genome amplification

with GenomePlex Whole Genome Amplification Kits (Sigma-Aldrich) as described previously.<sup>32</sup> Briefly, the purified fosmid DNA and the PCR product were amplified and then labeled as follows: G248P86579F1 and glycophorin E repeat-specific PCR product were labeled with digoxigenin-11-dUTP, G248P8211G10 was labeled with biotin-16-dUTP, G248P85804F12 was labeled with DNP-11-dUTP, and G248P80757F7 was labeled with Cy5-dUTP. All labeled dUTPs were purchased from Jena Bioscience.

The preparation of single-molecule DNA fibers by molecular combing and fiber-FISH was as previously published,<sup>3,33</sup> with the exception of post-hybridization washes, which consisted of three 5-min washes in  $2\times$  SSC at  $42^\circ\text{C}$ , instead of two 20-min washes in 50% formamide/50%  $2\times$  SSC at room temperature.

### Interphase-, Metaphase-FISH, and Karyotyping by Multiplex-FISH

Metaphase chromosomes were prepared from a human lymphoblastoid cell line (HG02554) purchased from Coriell Cell Repositories. Briefly, colcemid (Thermo Fisher Scientific) was added to a final concentration of  $0.1\ \mu\text{g}/\text{mL}$  for 1 hr, followed by treatment with hypotonic buffer (0.4% KCl in 10 mM HEPES [pH 7.4]) for 10 min and then fixed using 3:1 (v/v) methanol:acetic acid.

For interphase- and metaphase-FISH, G248P8211G10 labeled with Texas Red-dUTP, G248P85804F12 labeled with Atto488-XX-dUTP (Jena Bioscience), and RP11-325A24 labeled with Atto425-dUTP (Jena Bioscience) were used as probes. Slides pre-treatments included a 10-min fixation in acetone (Sigma-Aldrich), followed by baking at  $65^\circ\text{C}$  for 1 hr. Metaphase spreads on slides were denatured by immersion in an alkaline denaturation solution (0.5 M NaOH, 1.0 M NaCl) for 10 min, followed by rinsing in 1 M Tris-HCl (pH 7.4) solution for 3 min,  $1\times$  PBS for 3 min, and dehydration through a 70%, 90%, and 100% ethanol series. The probe mix was denatured at  $65^\circ\text{C}$  for 10 min before being applied onto the denatured slides. Hybridization was performed at  $37^\circ\text{C}$  overnight. The post-hybridization washes included a 5-min stringent-wash in  $1\times$  SSC at  $73^\circ\text{C}$ , followed by a 5-min rinse in  $2\times$  SSC containing 0.05% Tween20 (VWR) and a 2-min rinse in  $1\times$  PBS, both at room temperature. Finally, slides were mounted with SlowFade Gold mounting solution containing 4'6-diamidino-2-phenylindole (Thermo Fisher Scientific). Multiplex-FISH (M-FISH) with human 24-color painting probe, as previously described.<sup>34</sup>

Slides were examined using AxioImager D1 microscope equipped with appropriate narrow-band pass filters for DAPI, Aqua, FITC, Cy3, and Cy5 fluorescence. Digital image capture and processing was carried out using the SmartCapture software (Digital Scientific UK). Ten randomly selected metaphase cells were karyotyped based on the M-FISH and DAPI-banding patterns using the SmartType Karyotyper software (Digital Scientific UK).

### PCR for Fiber-FISH Probe Generation

The 3,632 bp glycophorin E repeat-specific PCR product for use as a fiber-FISH probe was generated by long PCR. Long PCRs were performed in a total volume of  $25\ \mu\text{L}$  using a *Taq/Pfu* DNA polymerase blend (0.6U *Taq* DNA polymerase/0.08U *Pfu* DNA polymerase), a final concentration of  $0.2\ \mu\text{M}$  primers Specific\_glycophorinE\_F and Specific\_glycophorinE\_R (Table S2), in 45 mM Tris-HCl (pH8.8), 11 mM  $(\text{NH}_4)_2\text{SO}_4$ , 4.5 mM  $\text{MgCl}_2$ , 6.7 mM 2-mercaptoethanol, 4.4 mM EDTA, 1 mM of each dNTP (sodium salt), 113  $\mu\text{g}/\text{mL}$  bovine serum albumin. Cycling

conditions were an initial denaturation of 94°C for 1 min, a first stage consisting of 20 cycles each of 94°C for 15 s and 65°C for 10 min, and a second stage consisting of 12 cycles each of 94°C for 15 s and 65°C for 10 min (plus 15 s/cycle); these were followed by a single incubation phase of 72°C for 10 min.

### Illumina Sequencing of DUP4 Samples

1 µg genomic DNA was randomly fragmented to a size of 350 bp by shearing, DNA fragments were end polished, A-tailed, and ligated with the NEBNext adaptor for Illumina sequencing, and further PCR enriched by P5 and indexed P7 oligos. The PCR products were purified (AMPure XP system) and the resulting libraries were analyzed for size distribution by an Agilent 2100 Bioanalyzer and quantified using real-time PCR.

Following sequencing on an Illumina platform, the resulting 150 bp paired-end sequences were examined for sequencing quality, aligned using BWA to the human reference genome (hg19 plus decoy sequences), sorted using samtools<sup>35</sup> and duplicate reads marked using Picard, generating the final bam file. Sequencing and initial bioinformatics was done by Novogene Ltd.

Sequence read depth was calculated using samtools to count mapped reads in non-overlapping 5 kb windows across the glycoporphin region. Read counts were normalized for coverage to a non-CNV region (chr4:145516270–145842585), then to the first 50 kb of the glycoporphin region which has diploid copy number of 2.

### DUP4 Junction Fragment PCR Genotyping

Primer sequences are shown in Table S2. PCR was conducted in a final volume of 10 µL in 1× Kapa A PCR buffer (a standard ammonium sulfate PCR buffer) with a final concentration of 1.5 mM MgCl<sub>2</sub>, ~10 ng genomic DNA, 0.2 mM of each of dATP, dCTP, dGTP, and dTTP, 1U *Taq* DNA Polymerase, 0.1 µM each of rs186873296F and rs186873296R, and 0.5 µM each of DUP4F2 and DUP4R2. Thermal cycling used an ABI Veriti Thermal cycler with an initial denaturation of 95°C for 2 min, followed by 35 cycles of 95°C 30 s, 58°C 30 s, 70°C 30 s, then followed by a final extension of 70°C for 5 min. 5 µL of each the PCR products were analyzed using standard horizontal electrophoresis on an ethidium-bromide-stained 2% agarose gel.

Routine genotyping included a DUP4 positive control in every experiment (sample HG02554). We distinguished homozygotes by quantification of DUP4 and control band intensity on agarose gels using ImageJ, and calculating the ratio of DUP4:control band intensity for each individual. At low allele frequencies, homozygotes are expected to be rare. After log<sub>2</sub> transformation, a cluster of four outliers of high ratio (>2 SD, log<sub>2</sub>ratio > 1.43) were clearly separated from the 278 other DUP4-positive samples, and these four were classified as homozygotes. The remaining 278 DUP4-positive samples were classified as heterozygotes.

### Family-Based Association Analysis

Associations between the three clinical phenotypes and DUP4 genotype were tested using QTD T v.2.6.1<sup>36</sup> on the full dataset of 167 pedigrees, using an orthogonal model. The heritability for all the clinical phenotypes was initially estimated using a model for polygenic variance. A test for total evidence of association was performed which included all individuals within the samples, retaining as much information as possible. This total test of association included environmental, polygenic heritability, and additive major locus variance components within the model. To control for population stratification within family, association was tested in

an orthogonal model including environmental, polygenic heritability, and additive major locus variance components. Direction of effect was estimated by comparing the hemoglobin level values, expressed as residuals from the regression model used to correct for age and sex, between the 262 unrelated individuals with hemoglobin level values from the Nyamisati cohort carrying ( $n = 70$ , mean = 1.67 g/L, standard deviation = 16.0 g/L) and not carrying ( $n = 192$ , mean = -0.12 g/L, standard deviation = 14.1 g/L) the DUP4 variant.

## Results

### The Physical Structure of the DUP4 Structural Variant

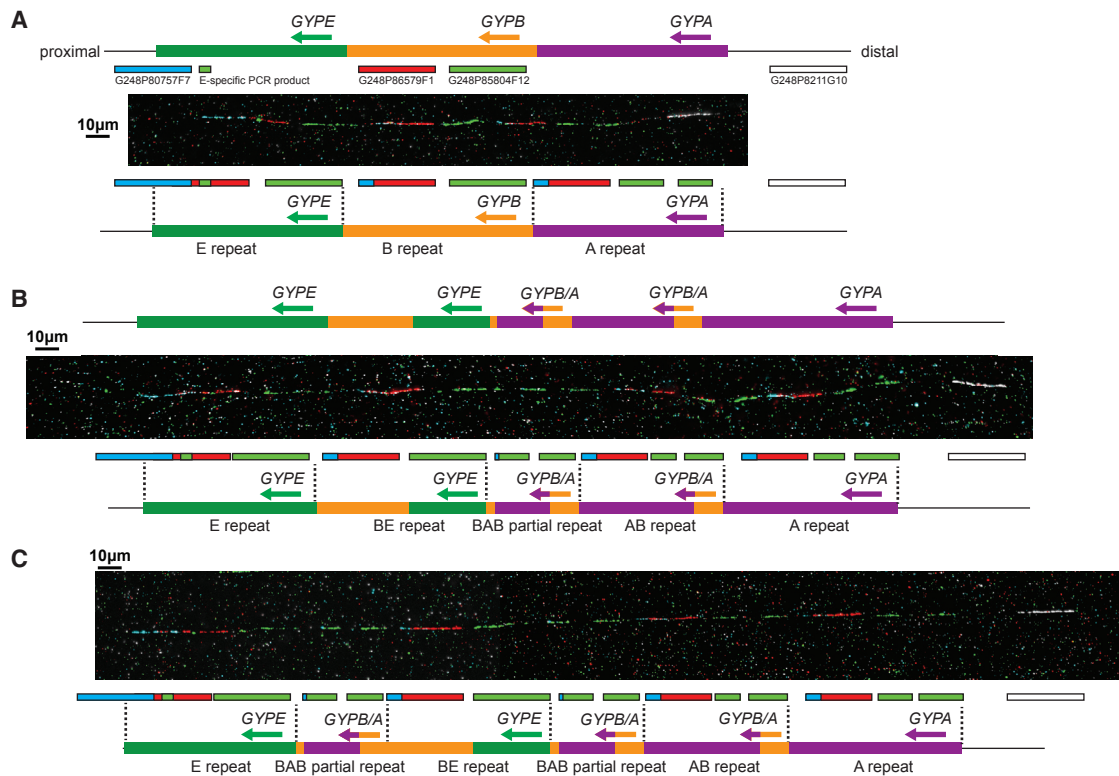
We initially decided to confirm this structural model by physical mapping of the DUP4 variant using fiber-FISH. We grew lymphoblastoid cells from the known DUP4 heterozygote from the 1000 Genomes project sample collection, HG02554, derived from a man with African ancestry from Barbados, and selected clones to act as fiber-FISH probes from the WIBR-2 human fosmid library spanning the region, based on fosmid end sequences previously mapped to the GRChr37 reference genome.

Fiber-FISH showed that the reference haplotype generates signals consistent with the genome reference sequence (Figure 1A). Because of the high sequence identity between the tandemly repeated glycoporphin regions, there is extensive cross-hybridization of probes that map to the *GYPB* repeat with the *GYPE* and *GYP A* repeats. The *GYPE* repeat can be distinguished by hybridization of a small *GYPE*-repeat-specific PCR product, and the *GYP A* repeat can be identified by a gap in the green fosmid probe signal, caused by 16 kb of unique sequence in the *GYP A* repeat. Also, the overlap of the distal end of the blue fosmid probe with the proximal end of the *GYPE* repeat means that a small amount of blue signal at the distal end of the *GYP B* and *GYPE* repeats is detected, confirming repeat length and orientation. We identified DNA fibers showing an arrangement completely consistent with the DUP4 model proposed previously (Figure 1B).

### Identification of a Somatic DUP4 Variant

However, we also visualized fibers with an extra partial repeat unit, which we called DUP4b (Figure 1C). This novel variant carries an extra copy of the partial A-B repeat, which harbors the *GYP B*/*GYP A* fusion gene. We selected a fosmid probe that spanned the 16 kb insertion specific to the *GYP A* repeat and showed that the extra copy was at least partly derived from the A repeat, consistent with the extra copy being an extra copy of the partial A-B repeat (Figure S1).

To rule out large-scale karyotype changes being responsible for our observations of the additional novel variant (DUP4b), we analyzed metaphase spreads of HG02554 lymphoblastoid cell line using metaphase-FISH, interphase FISH, and multiplex-FISH karyotyping (Figure S2). DUP4 and reference chromosomes could be distinguished by interphase-FISH on the basis of hybridization intensity of



**Figure 1. Fiber FISH Analysis of the DUP4 Heterozygote Sample HG02554**

(A) An example DNA fiber from the reference haplotype. The position and label color of the fosmid probes is indicated above the fiber on a representation of the human reference genome, and the interpretation of the FISH signals shown below the fiber. (B) An example DNA fiber from the DUP4a haplotype. The Leffler model of the DUP4 haplotype is indicated above the fiber. The interpretation of the FISH signals shown below the fiber. (C) An example DNA fiber from the DUP4b haplotype. The interpretation of the FISH signals is shown below the fiber.

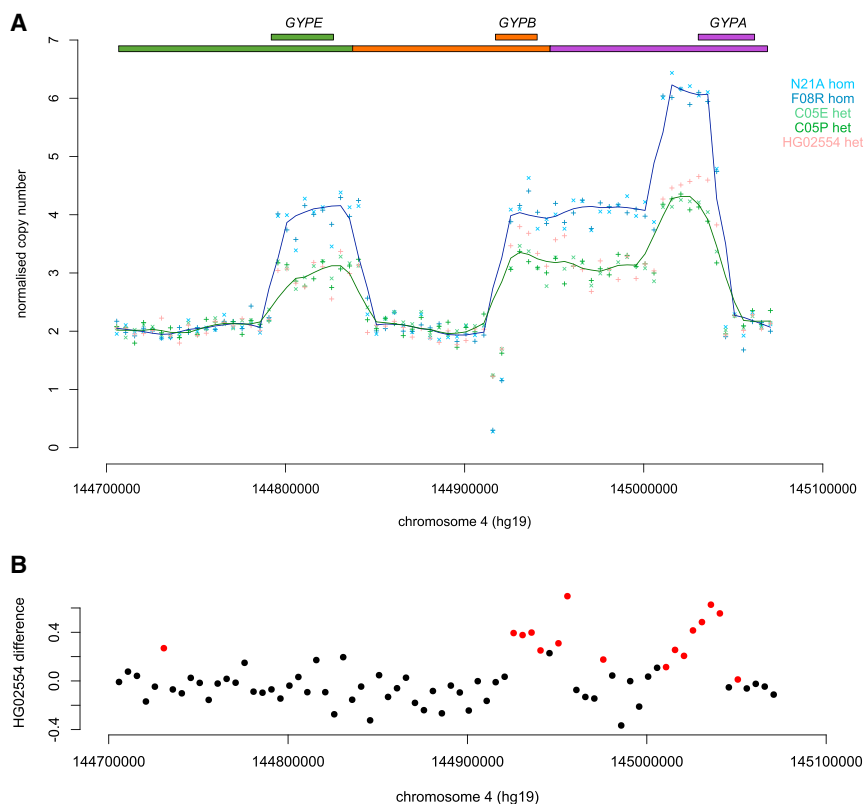
a fosmid probe mapping to *GYPB* (Figure S2B). No evidence of large-scale inter- and intrachromosomal rearrangements or aneuploidy was found in any of our experiments.

We hypothesized that DUP4b is a somatic variant that occurred through rearrangement of the original DUP4 variant (which we call DUP4a), but not the reference variant. If this is true, we would expect to observe an equal number of reference and DUP4 fibers from each of the parental chromosomes confirming the heterozygous DUP4 genotype of the source cells, but for the DUP4 fibers to be subdivided into DUP4a and DUP4b variants. Of 24 fibers examined from HG02554, 12 were reference and 12 were DUP4, and, of the 12 DUP4 variants, 7 were DUP4a and 5 were DUP4b, strongly supporting the model where DUP4b is a somatic rearrangement of DUP4a and the presence of two sub-clones (populations) of cells, one with reference and DUP4a haplotypes the other with reference and DUP4b haplotypes. We also analyzed the HG02554 cell line from the Oxford laboratory used in their study<sup>5</sup> and confirmed the existence of DUP4b by fiber-FISH. The high frequency of DUP4b variant chromosomes within the cell lines together with the observation of DUP4b in two cell line cultures suggests that DUP4b is a somatic variant of DUP4a that has arisen prior to the passage received by the Oxford laboratory or the Wellcome

Sanger Institute, either in the donor individual or early in the cell-culturing process, perhaps increasing in frequency due to the associated transformation cell bottleneck.<sup>37</sup>

To further characterize the somatic variation observed in HG02554, we Illumina sequenced at high depth (50×) HG02554 DNA purchased directly from Coriell Cell Repositories and extracted from their HG02554 lymphoblastoid cell line rather than extracted from our cell lines, together with peripheral-blood derived genomic DNA from two Tanzanian DUP4 homozygotes and two Tanzanian DUP4 heterozygotes. Analysis of sequence read depth across the glycoprotein repeat region showed the same pattern as that observed previously,<sup>5</sup> leading to a model that is confirmed by our fiber-FISH data (Figure 2A). DUP4 homozygotes show the expected increase to four copies and six copies in duplicated and triplicated regions, respectively.

We then compared the sequence read depth of HG02554 to the other two DUP4 heterozygotes to search for evidence of an increased copy number of the BAB partial repeat carrying the *GYPA/GYPB* fusion gene (Figure 1) suggested by our fiber-FISH data, which would reflect somatic mosaicism. HG02554 indeed shows a significant increase in DNA dosage in regions matching the BAB partial repeat, of around about 0.5, reflecting an extra copy of the region in ~50% of cells (Figure 2B).



**Figure 2. Sequence Read Depth Analysis of DUP4 Homozygotes and Heterozygotes** (A) Normalized sequence read depth of 5 kb windows spanning the reference sequence glycoprotein region for five samples. The lines show the Loess regression line ( $f = 0.1$ ) for homozygotes (blue) and heterozygotes (green). Gene positions and repeats, with respect to the reference sequence, are shown above the plot.

(B) The difference in HG02554 sequence read depth compared to the average sequence read depth of the two other heterozygotes C05E and C05P is shown in 5 kb windows across the glycoprotein region. Points highlighted in red are significantly different ( $p < 0.01$ ).

to malaria: hemoglobin levels in peripheral blood, parasite load, and mean number of clinical episodes of malaria, with hemoglobin levels showing the highest heritability of the three phenotypes in this cohort (Table 1). The DUP4 structural variant has previously been associated with both severe cerebral malaria and severe malarial anemia,<sup>5,39</sup> and both

are diagnoses related to our quantitative traits. For example, although the causes of hemoglobin level variation between individuals from a malaria-endemic region will be multifactorial, they will be strongly affected by malaria infection status of the individual, with infected individuals showing lower levels of hemoglobin.<sup>40</sup> At the extreme low end of the distribution of hemoglobin levels is anemia, a sign of malaria that is one important feature in the pathology of the disease.<sup>41</sup>

We analyzed data from a longitudinal study of a population from the village of Nyamisati, in the Rufiji river delta, 150 km south of Dar-es-Salaam, Tanzania, described previously.<sup>29–31</sup> This region was holoendemic for malaria, predominantly *P. falciparum*, which causes 99.5% of all recorded clinical episodes of malaria. Parasite prevalence was recorded as 75% at the start of the study in 1993, falling to 48% in 1998, as measured by microscopy in the 2- to 9-year-old children. A total of 962 individuals with pedigree information were genotyped; of these 278 were DUP4 heterozygotes and 4 were DUP4 homozygotes. Previous work has suggested that the DUP4 variant is at a frequency of about 3% in the Wasambaa of north-eastern Tanzania,<sup>5</sup> and our analysis found an allele frequency of 13.4% (95% confidence intervals 11.0%–16.1%) in the 348 unrelated individuals from Nyamisati village. For these unrelated individuals, 87 were DUP4 heterozygotes and 3 were DUP4 homozygotes, with genotype frequencies in Hardy-Weinberg equilibrium. We used the pedigree and genotype information from our full cohort to test for association of three malaria-related phenotypes with the

### Development of a Simple Robust DUP4 Genotyping Assay

Having characterized the structure of DUP4 variants, we designed a simple robust junction fragment PCR assay that would allow detection of the DUP4 variant (both DUP4a and DUP4b) in nanograms of genomic DNA, at a large scale. This involved designing allele-specific and paralog-specific PCR primers across a known breakpoint, a process made more challenging by the high sequence identity between paralogs. DUP4-specific primers had a modified locked nucleic acid base incorporated in the terminal 3' nucleotide to enhance specificity for the correct paralog.<sup>38</sup> We initially targeted the GYPA-GYPB breakpoint that created the fusion gene but found that a similar breakpoint was present in a frequent gene conversion allele. We therefore designed primers to target the breakpoint between the GYPE repeat and the GYPB repeat, which was predicted to be unique to DUP4.

The DUP4 genotype was determined using a duplex PCR approach, with one pair of primers specific for the DUP4 variant and a second pair amplifying across the SNP rs186873296, outside the structurally variable region, acting as a control for PCR amplification. The assay was validated against control samples showing different structural variants,<sup>5</sup> and samples showing no structural variation, to ensure DUP4 specificity (Table S1, Figure S3).

### Association of DUP4 with Malaria-Related Phenotypes

The DUP4 genotyping assay allowed us to investigate the association of DUP4 with three quantitative traits related

**Table 1. Association of DUP4 Allele with Malarial Phenotypes in the Nyamisiti Cohort**

Phenotype	Heritability of Phenotype (95% CI) <sup>a</sup>	Individuals	Association p Value
Hemoglobin	0.302 (0.136–0.469)	800	0.0054
Parasite load	0.104 (0.002–0.206)	864	0.39
Clinical episodes	0.221 (0.131–0.311)	939	0.72

<sup>a</sup>Calculated on this cohort using SOLAR<sup>30,31,46</sup>

DUP4 variant. Using a family-based association method modeled in QTDI,<sup>36</sup> we found a statistically significant association of the DUP4 variant with hemoglobin levels ( $p = 0.0054$ , Table 1). We estimated the direction of effect by comparing the mean corrected hemoglobin levels of unrelated individuals with and without the DUP4 variant. Individuals with the DUP4 allele showed a higher hemoglobin level compared to those without a DUP4 variant, showing that DUP4 variant is associated with higher hemoglobin levels.

## Discussion

In summary, we directly demonstrate that the DUP4 variant has a complex structure involving a duplication of *GYPE*, deletion of *GYPB*, and generation of two *GYPB/GYPA* fusion genes. The evolution of this particular rearrangement remains unclear. A model involving three intermediates has been suggested but none of these putative intermediates have yet been found.<sup>5</sup> Given the relatively limited numbers of individuals analyzed for glycoprotein structural variation so far, it is possible that these intermediate variants are rare or have been lost from the population. Indeed, given the extensive structural variation seen already at this locus, it seems likely that a high rate of genomic rearrangement generates complex variants that are mostly lost by genetic drift, with a few, such as DUP4, increasing in frequency due to positive selection. Further studies on the extensive variation in Africa are needed to fully characterize the variation at this locus.

We show that the DUP4 variant is associated with hemoglobin levels in a community setting indicating protection from malaria. Low levels of hemoglobin indicate anemia, which can reflect sub-clinical levels of malaria infection, and the village studied has a very high prevalence of *P. falciparum* infection, so our study supports the importance of the DUP4 variant in malaria protection. However, the absence of an association with either the number of clinical episodes of malaria or the parasite load is perhaps more puzzling. This may reflect the lower heritability of these traits compared to hemoglobin levels, and therefore the increased effect of non-genetic variation (Table 1). A recent case-control study of severe malaria in Kenyan children found an association of DUP4 with higher hemoglobin levels but not with parasite load, repeating the results we present here.<sup>39</sup> How DUP4 protects against malaria is unknown and alternatively these results may point to a

role in directly affecting erythrocyte invasion by the parasite, which is detectable in our cohort, rather than the more general phenotypes such as number of clinical malaria episodes or parasite load.

We also show that a novel somatic variant exists (DUP4b) with an extra *GYPB/GYPA* fusion gene, suggesting that this region may be prone to somatic rearrangements. We cannot rule out a somatic rearrangement in the transformation and culturing of the lymphoblastoid cell line, although it has been shown previously that such genetic changes introduced by EBV transformation are either rare<sup>15,42</sup> or overlap with regions known to undergo extensive programmed somatic rearrangement, such as the immunoglobulin loci.<sup>43</sup> It is possible, therefore, that the somatic variant originated in the donor patient given that the HG02554 B-lymphoblastoid cells from Oxford and the Wellcome Sanger Institute were both from the same batch of cells (passage #4, according to Coriell Cell Repositories); recent evidence suggests that such structural variant mosaics are likely to occur at a significant frequency, at least at certain loci.<sup>44,45</sup> We demonstrate that this somatic variant is able to be detected from high coverage short read sequence data, which will allow further analysis of somatic variation at this locus without cell material. Our data raise the intriguing possibility of heightened somatic instability and somatic mosaicism at this locus in DUP4 carriers, which might confer added protection against malaria.

## Accession Numbers

Sequence data are available from the European Nucleotide Archive for HG02554 (accession number ERP110671) and European Genome-Phenome Archive for the four Tanzanian samples (study accession number EGAS00001003239). Access to the Tanzanian sample sequences is restricted to projects related to malarial research.

## Supplemental Data

Supplemental Data include three figures and two tables and can be found with this article online at <https://doi.org/10.1016/j.ajhg.2018.10.008>.

## Acknowledgments

This work was funded by a SACB PhD studentship to W.A. and Wellcome Trust grant WT098051 (F.Y. and S.L.). This research used the SPECTRE High Performance Computing Facility at the University of Leicester. We wish to thank the villagers of Nyamisiti and the research team for continuous engagement

and contributions. We thank Ellen Leffler and Gavin Band for helpful comments on a previous version of this manuscript, Kirk Rockett for providing the HG02554 cells used by the Oxford laboratory, and Chris Tyler-Smith for support.

## Declaration of Interests

The authors have no competing interests

Received: July 27, 2018

Accepted: October 4, 2018

Published: November 1, 2018

## Web Resources

European Genome-phenome Archive (EGA), <https://www.ebi.ac.uk/ega>

European Nucleotide Archive, <https://www.ebi.ac.uk/ena>

OMIM, <http://www.omim.org/>

## References

- Schrider, D.R., Hahn, M.W., and Begun, D.J. (2016). Parallel evolution of copy-number variation across continents in *Drosophila melanogaster*. *Mol. Biol. Evol.* 33, 1308–1316.
- Chain, F.J.J., Feulner, P.G.D., Panchal, M., Eizaguirre, C., Samonte, I.E., Kalbe, M., Lenz, T.L., Stoll, M., Bornberg-Bauer, E., Milinski, M., and Reusch, T.B. (2014). Extensive copy-number variation of young genes across stickleback populations. *PLoS Genet.* 10, e1004830.
- Polley, S., Louzada, S., Forni, D., Sironi, M., Balaskas, T., Hains, D.S., Yang, E., and Hollox, E.J. (2015). Evolution of the rapidly mutating human salivary agglutinin gene (DMBT1) and population subsistence strategy. *Proc. Natl. Acad. Sci. USA* 112, 5105–5110.
- Schrider, D.R., Navarro, F.C.P., Galante, P.A.F., Parmigiani, R.B., Camargo, A.A., Hahn, M.W., and de Souza, S.J. (2013). Gene copy-number polymorphism caused by retrotransposition in humans. *PLoS Genet.* 9, e1003242.
- Leffler, E.M., Band, G., Busby, G.B.J., Kivinen, K., Le, Q.S., Clarke, G.M., Bojang, K.A., Conway, D.J., Jallow, M., Sisay-Joof, F., et al.; Malaria Genomic Epidemiology Network (2017). Resistance to malaria through structural variation of red blood cell invasion receptors. *Science* 356, 356.
- Charrier, C., Joshi, K., Coutinho-Budd, J., Kim, J.-E., Lambert, N., de Marchena, J., Jin, W.-L., Vanderhaeghen, P., Ghosh, A., Sassa, T., and Polleux, F. (2012). Inhibition of SRGAP2 function by its human-specific paralogs induces neoteny during spine maturation. *Cell* 149, 923–935.
- Dennis, M.Y., Nuttle, X., Sudmant, P.H., Antonacci, F., Graves, T.A., Nefedov, M., Rosenfeld, J.A., Sajjadian, S., Malig, M., and Kotkiewicz, H. (2012). Human-specific evolution of novel SRGAP2 genes by incomplete segmental duplication. *Cell* 149, 912.
- Peng, Z., Zhou, W., Fu, W., Du, R., Jin, L., and Zhang, F. (2015). Correlation between frequency of non-allelic homologous recombination and homology properties: evidence from homology-mediated CNV mutations in the human genome. *Hum. Mol. Genet.* 24, 1225–1233.
- MacArthur, J.A.L., Spector, T.D., Lindsay, S.J., Mangino, M., Gill, R., Small, K.S., and Hurles, M.E. (2014). The rate of non-allelic homologous recombination in males is highly variable, correlated between monozygotic twins and independent of age. *PLoS Genet.* 10, e1004195.
- Shwan, N.A.A., Louzada, S., Yang, F., and Armour, J.A.L. (2017). Recurrent rearrangements of human amylase genes create multiple independent CNV series. *Hum. Mutat.* 38, 532–539.
- Xu, D., Pavlidis, P., Taskent, R.O., Alachiotis, N., Flanagan, C., DeGiorgio, M., Blekhman, R., Ruhl, S., and Gokcumen, O. (2017). Archaic hominin introgression in Africa contributes to functional salivary MUC7 genetic variation. *Mol. Biol. Evol.* 34, 2704–2715.
- McConnell, M.J., Moran, J.V., Abyzov, A., Akbarian, S., Bae, T., Cortes-Ciriano, I., Erwin, J.A., Fasching, L., Flasch, D.A., Freed, D., et al.; Brain Somatic Mosaicism Network (2017). Intersection of diverse neuronal genomes and neuropsychiatric disease: The Brain Somatic Mosaicism Network. *Science* 356, 356.
- Abyzov, A., Mariani, J., Palejev, D., Zhang, Y., Haney, M.S., Tomasini, L., Ferrandino, A.F., Rosenberg Belmaker, L.A., Szekely, A., Wilson, M., et al. (2012). Somatic copy number mosaicism in human skin revealed by induced pluripotent stem cells. *Nature* 492, 438–442.
- Bruder, C.E., Piotrowski, A., Gijsbers, A.A., Andersson, R., Erickson, S., Diaz de Ståhl, T., Menzel, U., Sandgren, J., von Tell, D., Poplawski, A., et al. (2008). Phenotypically concordant and discordant monozygotic twins display different DNA copy-number-variation profiles. *Am. J. Hum. Genet.* 82, 763–771.
- Sudmant, P.H., Rausch, T., Gardner, E.J., Handsaker, R.E., Abyzov, A., Huddleston, J., Zhang, Y., Ye, K., Jun, G., Fritz, M.H., et al.; 1000 Genomes Project Consortium (2015). An integrated map of structural variation in 2,504 human genomes. *Nature* 526, 75–81.
- O’Huallachain, M., Karczewski, K.J., Weissman, S.M., Urban, A.E., and Snyder, M.P. (2012). Extensive genetic variation in somatic human tissues. *Proc. Natl. Acad. Sci. USA* 109, 18018–18023.
- Usher, C.L., and McCarroll, S.A. (2015). Complex and multi-allelic copy number variation in human disease. *Brief. Funct. Genomics* 14, 329–338.
- Reid, M.E. (2009). MNS blood group system: a review. *Immunohematology* 25, 95–101.
- Daniels, G. (2008). *Human Blood Groups* (John Wiley & Sons).
- Band, G., Rockett, K.A., Spencer, C.C., Kwiatkowski, D.P.; and Malaria Genomic Epidemiology Network (2015). A novel locus of resistance to severe malaria in a region of ancient balancing selection. *Nature* 526, 253–257.
- Bigham, A.W., Magnaye, K., Dunn, D.M., Weiss, R.B., and Bamshad, M. (2018). Complex signatures of natural selection at GYPA. *Hum. Genet.* 137, 151–160.
- Ko, W.-Y., Kaercher, K.A., Giombini, E., Marcatili, P., Froment, A., Ibrahim, M., Lema, G., Nyambo, T.B., Omar, S.A., Wambebe, C., et al. (2011). Effects of natural selection and gene conversion on the evolution of human glycoporphins coding for MNS blood polymorphisms in malaria-endemic African populations. *Am. J. Hum. Genet.* 88, 741–754.
- Baum, J., Ward, R.H., and Conway, D.J. (2002). Natural selection on the erythrocyte surface. *Mol. Biol. Evol.* 19, 223–229.
- Johnson, K.E., and Voight, B.F. (2018). Patterns of shared signatures of recent positive selection across human populations. *Nat Ecol Evol* 2, 713–720.
- Orlandi, P.A., Klotz, F.W., and Haynes, J.D. (1992). A malaria invasion receptor, the 175-kilodalton erythrocyte binding

- antigen of *Plasmodium falciparum* recognizes the terminal Neu5Ac(alpha 2-3)Gal- sequences of glycophorin A. *J. Cell Biol.* *116*, 901–909.
26. Sim, B.K., Chitnis, C.E., Wasniowska, K., Hadley, T.J., and Miller, L.H. (1994). Receptor and ligand domains for invasion of erythrocytes by *Plasmodium falciparum*. *Science* *264*, 1941–1944.
  27. Mayer, D.C., Cofie, J., Jiang, L., Hartl, D.L., Tracy, E., Kabat, J., Mendoza, L.H., and Miller, L.H. (2009). Glycophorin B is the erythrocyte receptor of *Plasmodium falciparum* erythrocyte-binding ligand, EBL-1. *Proc. Natl. Acad. Sci. USA* *106*, 5348–5352.
  28. Baldwin, M.R., Li, X., Hanada, T., Liu, S.C., and Chishti, A.H. (2015). Merozoite surface protein 1 recognition of host glycophorin A mediates malaria parasite invasion of red blood cells. *Blood* *125*, 2704–2711.
  29. Färnert, A., Yman, V., Homann, M.V., Wandell, G., Mhoja, L., Johansson, M., Jesaja, S., Sandlund, J., Tanabe, K., Hammar, U., et al. (2014). Epidemiology of malaria in a village in the Rufiji River Delta, Tanzania: declining transmission over 25 years revealed by different parasitological metrics. *Malar. J.* *13*, 459.
  30. Carpenter, D., Rooth, I., Färnert, A., Abushama, H., Quinnell, R.J., and Shaw, M.A. (2009). Genetics of susceptibility to malaria related phenotypes. *Infect. Genet. Evol.* *9*, 97–103.
  31. Carpenter, D., Färnert, A., Rooth, I., Armour, J.A., and Shaw, M.-A. (2012). CCL3L1 copy number and susceptibility to malaria. *Infect. Genet. Evol.* *12*, 1147–1154.
  32. Gribble, S.M., Wiseman, F.K., Clayton, S., Prigmore, E., Langley, E., Yang, F., Maguire, S., Fu, B., Rajan, D., Sheppard, O., et al. (2013). Massively parallel sequencing reveals the complex structure of an irradiated human chromosome on a mouse background in the Tc1 model of Down syndrome. *PLoS ONE* *8*, e60482.
  33. Louzada, S., Komatsu, J., and Yang, F. (2017). Fluorescence in situ hybridization onto DNA fibres generated using molecular combing. In *Fluorescence In Situ Hybridization (FISH) Application Guide*, T. Liehr, B. Heidelberg, ed. (Springer-Verlag), pp. 275–293.
  34. Agu, C.A., Soares, F.A., Alderton, A., Patel, M., Ansari, R., Patel, S., Forrest, S., Yang, F., Lineham, J., Vallier, L., and Kirton, C.M. (2015). Successful generation of human induced pluripotent stem cell lines from blood samples held at room temperature for up to 48 hr. *Stem Cell Reports* *5*, 660–671.
  35. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R.; and 1000 Genome Project Data Processing Subgroup (2009). The sequence alignment/map format and SAMtools. *Bioinformatics* *25*, 2078–2079.
  36. Abecasis, G.R., Cardon, L.R., and Cookson, W.O. (2000). A general test of association for quantitative traits in nuclear families. *Am. J. Hum. Genet.* *66*, 279–292.
  37. Plagnol, V., Uz, E., Wallace, C., Stevens, H., Clayton, D., Ozcelik, T., and Todd, J.A. (2008). Extreme clonality in lymphoblastoid cell lines with implications for allele specific expression analyses. *PLoS ONE* *3*, e2966.
  38. Latorra, D., Hopkins, D., Campbell, K., and Hurley, J.M. (2003). Multiplex allele-specific PCR with optimized locked nucleic acid primers. *Biotechniques* *34*, 1150–1152, 1154, 1158.
  39. Ndila, C.M., Uyoga, S., Macharia, A.W., Nyutu, G., Peshu, N., Ojal, J., Shebe, M., Awuondo, K.O., Mturi, N., Tsofa, B., et al.; MalariaGEN Consortium (2018). Human candidate gene polymorphisms and risk of severe malaria in children in Kilifi, Kenya: a case-control association study. *Lancet Haematol.* *5*, e333–e345.
  40. Ekvall, H., Premji, Z., Bennett, S., and Bjorkman, A. (2001). Hemoglobin concentration in children in a malaria holoendemic area is determined by cumulated *Plasmodium falciparum* parasite densities. *Am. J. Trop. Med. Hyg.* *64*, 58–66.
  41. Phillips, R.E., and Pasvol, G. (1992). Anaemia of *Plasmodium falciparum* malaria. *Baillieres Clin. Haematol.* *5*, 315–330.
  42. Nickles, D., Madireddy, L., Yang, S., Khankhanian, P., Lincoln, S., Hauser, S.L., Oksenberg, J.R., and Baranzini, S.E. (2012). In depth comparison of an individual's DNA and its lymphoblastoid cell line using whole genome sequencing. *BMC Genomics* *13*, 477.
  43. Shirley, M.D., Baugher, J.D., Stevens, E.L., Tang, Z., Gerry, N., Beiswanger, C.M., Berlin, D.S., and Pevsner, J. (2012). Chromosomal variation in lymphoblastoid cell lines. *Hum. Mutat.* *33*, 1075–1086.
  44. Campbell, I.M., Yuan, B., Robberecht, C., Pfundt, R., Szfranski, P., McEntagart, M.E., Nagamani, S.C.S., Erez, A., Bartnik, M., Wiśniowiecka-Kowalnik, B., et al. (2014). Parental somatic mosaicism is underrecognized and influences recurrence risk of genomic disorders. *Am. J. Hum. Genet.* *95*, 173–182.
  45. Gajicka, M. (2016). Unrevealed mosaicism in the next-generation sequencing era. *Mol. Genet. Genomics* *291*, 513–530.
  46. Almsy, L., and Blangero, J. (1998). Multipoint quantitative-trait linkage analysis in general pedigrees. *Am. J. Hum. Genet.* *62*, 1198–1211.