1  **Structural variation of the malaria-associated human glycophorin A-B-E region**

2

3  Sandra Louzada (1,6,7), Walid Algady (2), Eleanor Weyell (2), Luciana W. Zuccherato (3),

4  Paulina Brajer (2), Faisal Almalki (2), Marilia O Scliar (4), Michel S Naslavsky (4), Guilherme L

5  Yamamoto (4), Yeda A O  Duarte (5), Maria Rita Passos-Bueno (4), Mayana Zatz (4), Fengtang

6  Yang (1), Edward J Hollox (2) *

7

8  1.  Wellcome Sanger Institute, Hinxton, Cambridge, UK

9  2.  Department of Genetics and Genome Biology, University of Leicester, UK

10  3.  Department of Parasitology, Universidade Federal de Minas Gerais, Belo Horizonte,

11  Brazil

12  4.  Human Genome and Stem Cell Research Center, Department of Genetics and

13  Evolutionary Biology, Instituto de Biociências, Universidade de São Paulo, São Paulo,

14  Brazil.

15  5.  School of Nursing, Universidade de São Paulo, São Paulo, Brazil

16  6.  Present address: Laboratory of Cytogenomics and Animal Genomics (CAG),

17  Department of Genetics and Biotechnology, University of Trás-os-Montes and Alto

18  Douro (UTAD), Vila Real, Portugal

19  7.  Present address: BioISI – Biosystems & Integrative Sciences Institute, Faculty of Sciences,

20  University of Lisboa, Lisbon, Portugal

21  * Corresponding author

22  **Abstract**

23  Approximately 5% of the human genome consists of structural variants, which are enriched for

24  genes involved in the immune response and cell-cell interactions. A well-established region of

25  extensive structural variation is the glycophorin gene cluster, comprising three tandemly-

26  repeated regions about 120kb in length, carrying the highly homologous genes *GYPA*, *GYPB*

27  and *GYPE*. Glycophorin A and glycophorin B are glycoproteins present at high levels on the

28  surface of erythrocytes, and they have been suggested to act as decoy receptors for viral

29  pathogens. They act as receptors for invasion of a causative agent of malaria, *Plasmodium*

30  *falciparum*. A particular complex structural variant (DUP4) that creates a *GYPB*/*GYPA* fusion

31  gene is known to confer resistance to malaria. Many other structural variants exist, and remain

32  poorly characterised. Here, we analyse sequences from 6466 genomes from across the world for

33  structural variation at the glycophorin locus, confirming 15 variants in the 1000 Genomes

34  project cohort, discovering 9 new variants, and characterising a selection using fibre-FISH and

35  breakpoint mapping. We identify variants predicted to create novel fusion genes and a common

36  inversion duplication variant at appreciable frequencies in West Africans. We show that almost

37  all variants can be explained by unequal cross over events (non-allelic homologous

38  recombination, NAHR) and. by comparing the structural variant breakpoints with

39    recombination hotspot maps, show the importance of a particular meiotic recombination

40    hotspot on structural variant formation in this region.

41

42    **Keywords**

43    Structural variation, copy number variation, inversion, immune response, glycophorin, *GYPA*,

44    *GYPB*, *GYPE*, erythrocytes, malaria

## Introduction

Human genetic variation encompasses single nucleotide variation, short insertion-deletions and structural variation. Structural variation includes copy number variation, tandem repeat variation, inversion and polymorphic retrotransposons. Structural variation is responsible for much of the differences in DNA sequence between individual human genomes (Sudmant et al. 2015; Zarrei et al. 2015; Hehir-Kwa et al. 2016), yet analysis of the phenotypic importance of structural variation has lagged behind the rapid progress made in studies of single nucleotide variation (Hollox and Hoh 2014; Usher and McCarroll 2015; Huddleston and Eichler 2016). This is mainly because of technical limitations in detecting, characterising, and genotyping structural variants both directly (Cantsilieris et al. 2014) and indirectly by imputation (Handsaker et al. 2015). However, a combination of new technical approaches using genome sequencing data to detect structural variation and larger datasets allowing more robust imputation of structural variation have begun to show that some structural variants at an appreciable frequency in populations do indeed contribute to clinically-important phenotypes (Sekar et al. 2016; Raffield et al. 2018).

One such example is the identification of a structural variant called DUP4 at the human glycophorin gene locus, which confers a reduced risk of severe malaria and protection against malarial anemia (Leffler et al. 2017; Algady et al. 2018; Ndila et al. 2018). The glycophorin gene locus consists of three ~120 kb tandem repeats sharing ~97% identity, each repeat carrying a closely-related glycophorin gene, starting from the centromeric end: glycophorin E (*GYPE*), glycophorin B (*GYPB*) and glycophorin A (*GYPA*) (Vignal et al. 1990; Onda et al. 1994). Large tandem repeats, like the glycophorin locus, are prone to genomic rearrangements, and indeed the DUP4 variant is a complex variant that generates a *GYPB-GYPA* fusion gene, with potential somatic variation in fusion gene copy number (Leffler et al. 2017; Algady et al. 2018). The mechanism of resistance to malaria of this gene is not fully understood, but although both glycophorin A and glycophorin B interact with receptors on *Plasmodium falciparum*, recent data suggest that alteration of receptor-ligand interactions is not important. Instead, it seems likely that DUP4 is associated with more complex alterations in the protein levels at the red blood cell surface resulting in increased red blood cell tension, mediating its protective effect against *P. falciparum* invasion (Kariuki et al. 2018). Given the size of effect of DUP4 in protection against malaria (odds ratio ~0.6) and the frequency of the allele (up to 13% in Tanzania), it is clinically very significant, although it appears to be geographically restricted to East Africa (Leffler et al. 2017; Algady et al. 2018).

Other structural variants in the glycophorin region have been identified in the 1000 Genomes project samples by using sequence read depth analysis of 1.6kb bins combined with a Hidden

83    Markov Model approach to identify regions of copy number gain and loss (Leffler et al. 2017).
84    This builds upon identification of extensive CNV is this area by array CGH (Conrad et al. 2009)
85    and indeed by previous analysis of rare MNS (Miltenberger) blood groups, such as $M^K$, caused
86    by homozygous deletion of both *GYPA* and *GYPB* (Vignal et al. 1990). The variants were
87    classified as DUP and DEL representing gain and loss of sequence read depth respectively.
88    Although only DUP4 has been found to be robustly associated with clinical malaria phenotypes,
89    it is possible that some of the other structural variants are also protective, but are either rare,
90    recurrent, or both rare and recurrent, making imputation from flanking SNP haplotypes and
91    genetic association with clinical phenotypes challenging.
92
93    It is important, therefore, to extend the catalogue of structural variants at this locus and robustly
94    characterise their nature and likely effect on the number of full-length and fusion glycophorin
95    genes. In this study we use sequence read depth analysis of 6466 genomes from across the
96    world, followed by direct analysis of structural variants using fibre-FISH and breakpoint
97    mapping using paralogue-specific PCR and Sanger sequencing. This will allow future
98    development of robust yet simple PCR-based assays for each structural variant and detailed
99    analysis of the phenotypic consequences of particular structural variants on malaria infection
100   and other traits. We also examine the pattern of structural variation breakpoints in relation to
101   their mechanism of generation and known meiotic recombination hotspots within the region,
102   and the relative allele frequencies across the world. Together, this allows us to gain some
103   insight into the evolutionary context of the extensive structural variation at the glycophorin
104   locus.
105
106
107   **Methods**
108
109   *Sequencing data*
110   Sequence alignment files (.bam format7) from four cohorts (1000 Genomes Project ENA
111   accession number PRJNA262923) with a mean coverage of 7.4x (Auton et al. 2015), Simons
112   Diversity Project ENA accession number PRJEB9586 with a mean coverage of 43x (Mallick et al.
113   2016), and the Gambian Genome Diversity project mean coverage 4x, ENA study IDs
114   ERP001420, ERP001781, ERP002150, ERP002385) (Band et al. 2019) were downloaded from the
115   European Nucleotide Archive or from the International Genome Sample Resource site
116   http://www.internationalgenome.org/data-portal/ (Clarke et al. 2017). Brazilian sequence
117   alignment files from the SABE (Health, Wellbeing and Aging) study (Barbosa et al. 2005) and a
118   sample of cognitively healthy octogenarians enrolled at the Human Genome and Stem Cell
119   Research Center (80+), with a mean coverage of 30x for 1324 individuals generated at Human
120   Longevity Inc. (HLI, San Diego, California) (Telenti et al. 2016).

121

122 DNA sequences from the 1000 Genomes project and the Simons diversity project had been
123 previously aligned to reference GRCh37 (hg19) to generate the alignment bam files. The
124 exception is sample NA18605, which was previously sequenced at high coverage (Lan et al.
125 2017) downloaded as paired-end Illumina sequences in fastq format (ENA sample accession
126 number SAMN00001619), and aligned to GRCh37 using standard approaches: FastQC v0.11.5
127 and Cutadapt v01.11 to trim reads and adapters, mapping using  BWA-MEM v0.7.15,
128 processing of the BAM files using SAMtools v1.8, local realignment was done using GATK v3.6
129 and duplicate reads marked using Picard v.1 and removed using SAMtools. Samples from the
130 Brazilian genomes and the Gambian genome diversity project had been aligned to GRCh38.
131

132 Throughout this paper, all loci are given using GRCh37 reference genome coordinates. For
133 analyses on GRCh38 alignments, genome coordinates were translated from the GRCh37
134 coordinates using the Liftover tool within the UCSC Genome Browser (Kent et al. 2002).
135

136 *Structural variant detection*
137 For each sample, we used SAMtools (SAMtools view –c –F 4) (Li et al. 2009) on indexed
138 bam files to count mapped reads to the glycophorin region (chr4:144745739-145069133) and a
139 reference region chr4:145516270-145842585. The reference region has no segmental duplications,
140 and is absent from copy number variation according to the gold standard track of the database
141 of Genomic Variants (DGV) (MacDonald et al. 2014). A ratio of the number or reads mapping to
142 the glycophorin region to the number of reads mapping to the reference region allows an
143 estimate of the total increase or decrease of sequence depth spanning the glycophorin region
144 (reflecting copy number gain or copy number loss, respectively). Following plotting these data
145 for each cohort on a histogram and observation of distinct clusters (supplementary figure 1),
146 samples with a ratio below 0.9 were classified as potential deletions and those above 1.1
147 potential duplications.
148 For the samples with potential deletions and duplications, number of mapped reads was
149 calculated across the glycophorin region in 5kb non-overlapping windows, and values,
150 normalised to average read count and diploid copy number, were plotted. Presence and nature
151 of structural variants were assessed by examination of the plots, and particular variants called
152 by plotting together with a reference sample for that variant. For the Simons Diversity Project
153 samples, 114 potential deletions were identified, much more than in other cohorts
154 (Supplementary figure 1). Inspection of these plots showed that 101 of these samples showed a
155 small apparent ~15kb deletion at the *GYPE* gene. This deletion was not found previously by
156 others (Leffler et al. 2017) or by us in any other cohort, and coincides with a region of low
157 mappability, suggesting that this may be an artefact caused either by particular filtering

158 conditions or the particular genome assembly (GRCh37d5) that includes decoy sequences.

159 These 101 samples were treated as being homozygous for the reference structure.

160

161 *Fiber-FISH*

162 The probes used in this study included four WIBR-2 fosmid clones selected from the

163 UCSC Genome Browser GRCh37/hg19 assembly and a 3632-bp PCR product that is specific for

164 the glycophorin E repeat (Algady et al. 2018). Probes were made by amplification with

165 GenomePlex Whole Genome Amplification Kits (Sigma-Aldrich) as described previously

166 (Gribble et al. 2013). Briefly, the purified fosmid DNA and the PCR product were amplified and

167 then labeled as follow: G248P86579F1, G248P89366H1 and glycophorin E repeat-specific PCR

168 product were labeled with digoxigenin-11-dUTP, G248P8211G10 was labeled with biotin-16-

169 dUTP, G248P85804F12 was labeled with DNP-11-dUTP and G248P80757F7 was labeled with

170 Cy5-dUTP. All labeled dUTPs were purchased from Jena Bioscience.

171 The preparation of single-molecule DNA fibers by molecular combing and fiber-FISH

172 was as previously published (Louzada et al. 2017), with the exception of post-hybridization

173 washes, which consisted of three 5-min washes in 2× SSC at 42°C, instead of two 20-min washes

174 in 50% formamide/50% 2× SSC at room temperature.

175

176 *Breakpoint analysis using PCR and Sanger sequencing*

177 Using the 5kb window sequence read count data, PCR primers were designed so that a

178 PCR product spanned the predicted breakpoints for each deletion and duplication. The 3'

179 nucleotide for each PCR primer was designed to match uniquely to a particular glycophorin

180 repeat, and to mismatch the other two glycophorin repeats. Annealing specificity of the PCR

181 primer was enhanced by incorporating a locked nucleic acid at that particular 3' position of the

182 PCR primer (Latorra et al. 2003). Long-range PCR amplification used 10 ng genomic DNA in a

183 final volume of 25.5 μl, including 0.5 μl of each 10μM primer, 0.075U *Pfu* DNA polymerase,

184 0.625U *Taq* DNA polymerase, and 2.25 μl of PCR buffer (45 mM Tris-HCl (pH 8.8), 11 mM

185 ammonium sulphate, 4.5 mM magnesium chloride, 6.7 mM 2-mercaptoethanol, 4.4 mM EDTA

186 (pH 8.0), 113 μg/mL non-acetylated Bovine Serum Albumin (BSA) (Ambion®) and 1 mM of

187 each dNTP (Promega) (Jeffreys et al. 1990)). The reaction was thermal cycled as follows: 94°C 1

188 minute, followed by 20 cycles of 94°C for 15s, x°C for 10 minutes, followed by 15 cycles of 94°C

189 for 15s, x°C for 10 minutes+15s for each successive cycle, followed by a final extension at 72°C

190 for 10 minutes, where x is the annealing temperature for a particular primer pair shown in

191 supplementary table 1. PCR products were purified using agarose gel electrophoresis (Ma and

192 Difazio 2008) and Sanger sequenced using standard approaches. PCR primers are shown in

193 supplementary table 1. Multiple alignments with paralogous reference sequences used MAFFT

194 v7 (Katoh and Standley 2013) available at the EMBL-EBI Job Dispatcher framework (Li et al.

195 2015). A breakpoint was called in the transition region between three paralogous sequence

196    variants corresponding to one glycophorin repeat and three paralogous sequence variants

197    corresponding to the alternative glycophorin repeat.

198

199    *Breakpoint analysis using high depth sequences*

200        For particular variants, copy number breakpoints were refined by inspecting sequence

201    read depth in 1kb windows spanning the likely breakpoints identified by the 5kb window

202    analysis. Changes in read depth were then confirmed directly using the Integrative Genome

203    Viewer (Thorvaldsdóttir et al. 2012).

204

205    *Nomenclature of variants*

206        We used the same nomenclature as Leffler et al. 2017 when our variant could be

207    identified as the same variant in the same sample from the 1000 Genomes project. In some

208    instances, we could not distinguish particular singleton variants called from more common

209    called variants. For example, DUP27 carried by sample NA12249 could not easily be

210    distinguished from the more frequent DUP2, and DUP24 carried by HG04038 could not be

211    distinguished from DUP8. Other variants, which either had not been unambiguously identified

212    in the 1000 Genomes previously or were identified in other sample cohorts, were given DEL or

213    DUP numbers following on from variants catalogued previously.

214

215    *Analysis of recombination hotspots*

216        Previously published data on hotspot location and type (Pratto et al. 2014) was

217    converted to BED format and intersected with the breakpoint locations in BED format using

218    BEDTools v 2.28.0 (Quinlan and Hall 2010). The statistical significance of the overlap was

219    calculated using the fisher command in BEDTools, which uses a Fisher's exact test on the

220    number of overlaps observed between two BED files.

221

222

223    **Results**

224

225    *Structural variation using sequence read depth analysis*

226

227    Previous work by us and others has shown that unbalanced structural variation - that is,

228    variation that causes a copy number change - can be effectively discovered by measuring the

229    relative depth of sequence reads across the glycophorin region (Leffler et al. 2017; Algady et al.

230    2018). We analysed a total of 6466 genomes from four datasets spanning the globe - the 1000

231    Genomes phase 3 project set, the Gambian variation project, the Simons diversity project, and

232    the Brazilian genomes project. We took a different sequence read depth approach to that

233    previously used, counting the reads that map to the glycophorin repeat region and dividing by

234  the number of reads mapping to a nearby non-structurally variable region to normalise for read
235  depth.  By analysing each cohort of diploid genomes as a group, we could identify outliers
236  where a higher value indicated a potential duplication or more complex gain of sequence, and
237  lower values indicated a potential deletion (Supplementary figure 1). Sequence read depth was
238  analysed in 5kb windows across each of the outlying diploid genomes to identify and classify
239  the structural variant.
240
241  Since structural variant calling had been previously done on the 1000 Genomes project cohort,
242  this provided a useful comparison to assess our approach. We analysed 2490 samples from this
243  cohort and identified five distinct deletions carried by 88 individuals that were previously
244  identified, and 16 distinct duplications carried by 34 individuals (table 1) that were also
245  previously identified. We also identified a new duplication variants, termed DUP29 (a
246  duplication of GYPB), that had not been robustly identified previously in that cohort. However,
247  as expected, smaller duplications, most notably DUP1, were not detected by our approach. We
248  extended our analysis to 390 Gambian diploid genomes and identified 51 samples with DEL1 or
249  DEL2 variants, and DEL16, subsequently characterized in the Brazilian cohort below. Two
250  samples were heterozygous for the duplication DUP5.
251
252  Both 1000 Genomes and Gambian Genome Diversity samples have been sequenced to low
253  depth. High depth sequencing will allow more robust identification of structural variation by
254  improving the signal/noise ratio of sequence read depth analysis. We analysed the publically
255  available high-depth data from the Simons Diversity Project for glycophorin variation. From the
256  273 individuals, 4 different deletion types were carried by 13 individuals, and 3 different
257  duplication types were carried by 5 individuals. A novel deletion, DEL15 was identified which
258  deleted part of *GYPB* and part of *GYPE* in an individual from Bergamo in Italy, and a novel
259  duplication was observed in three individuals from Papua New Guinea, termed DUP30 and
260  duplicating the *GYPB* gene. Another duplication variant (DUP8), which is the largest variant
261  found so far, duplicated 240kb, creating an extra full length *GYPB* gene and a *GYPE-GYPA*
262  fusion gene (Table 1).
263
264  A further 1324 samples sequenced to high coverage diploid genomes from Brazil were
265  analysed, which, given the extensive admixture from Africa in the Brazilian population, are
266  likely to be enriched for glycophorin variants from Africa. Four new duplication variants
267  (DUP33-DUP35) and three new deletion variants were found (DEL16, DEL17, DEL18), two of
268  which of which delete the *GYPB* gene (Table 1).
269
270  *Fibre-FISH analysis of structural variants*
271

272 Sequence read depth analysis shows copy number gain and loss with respect to the reference
273 genome to which the sequence reads are mapped, but it does not determine the physical
274 structure of the structural variant. For the more common glycophorin structural variants, we
275 used fibre-FISH in order to determine the physical structure. In all cases, a set of multiplex FISH
276 probes, with each probe being visualised with a unique fluorochrome, was used so that the
277 orientation and placement of the repeats could be identified (Figure 1). The repeated nature of
278 the glycophorin region means that the green and red probes from the *GYPB* repeat cross-
279 hybridise with the other repeats, with the *GYPA* repeat distinguishable from the *GYPB* and
280 *GYPE* repeats by a 16kb insertion resulting in a small gap of signal in the green probe (Figure 1).
281
282 For most variants the fibre-FISH results confirmed the structure previously predicted (Leffler et
283 al. 2017) and expected if the variants had been generated by non-allelic homologous
284 recombination between the glycophorin repeats (Figures 2 and 3). However, three variants
285 showed a complex structure that could not be easily predicted from the sequence read depth
286 analysis. The DUP4 variant shows a complex structure and has been described previously
287 (Algady et al. 2018). Two other structural variants (DUP5 and DUP26) also showed complex
288 patterns of gains or losses, and fibre-FISH clearly shows the physical structure of the variant,
289 including inversions.
290
291 The more frequent of these two complex structural variants, DUP5, seems to be restricted to
292 Gambia, as it is found once in the GWD population from the 1000 Genomes project and twice in
293 the Jola population from the Gambian Genome diversity project (Table 1). Sequence read depth
294 analysis suggests that DUP5 has two extra copies of *GYPE* and an extra copy of *GYPB*, with an
295 additional duplication distal of *GYPA* outside the glycophorin repeated region (Figure 4a).
296 Fibre-FISH analysis on cells from an individual carrying the DUP5 variant (HG02585)
297 confirmed heterozygosity of the variant, with one allele being the reference allele, and revealed,
298 for the first time, that the variant allele presents a complex pattern of duplication and
299 rearrangement, with part of the fosmid (pseudocoloured in white) mapping distal to *GYPA*
300 being translocated into the glycophorin repeated region, adjacent to the green-coloured fosmid
301 (Figure 4b). Alternative fibre-FISH analysis using an additional fosmid probe mapping distally,
302 and labelled in red, confirmed this (Figure 4c). The pattern of FISH signals occurring distally to
303 the translocation suggests that the immediately adjacent glycophorin repeat is inverted. To
304 distinguish the distal end of the *GYPB* repeat from the distal end of the *GYPE* repeat, a pink-
305 coloured probe from a short *GYPE*-repeat-specific PCR product was also used for fibre-FISH,
306 and clearly shows only a single copy of the distal end of the *GYPB* repeat in the DUP5 variant,
307 at the same position as the reference. The predicted breakpoint between the non-duplicated
308 sequence distal to *GYPA* and duplicated sequence within the duplicated region was amplified
309 by PCR and Sanger sequenced, confirming that the non-duplicated sequence was fused to an

310    inverted *GYPB* repeat sequence (Figure 4d). The model suggested by the fibre-FISH and
311    breakpoint analysis is consistent with the overall pattern of sequence depth changes observed
312    (Figure 4a). The sequence outside the glycophorin repeat corresponds to an ERV-MaLR long
313    terminal retroviral element, but the sequence inside the glycophorin repeat sequence is not,
314    suggesting that non-allelic homologous recombination was not the mechanism for formation of
315    this breakpoint. However, there is a 4bp microhomology (GTGT) between the two sequences,
316    suggesting that microhomology-mediated end joining could be a mechanism for formation of
317    this variant.
318
319    The DUP26 variant was observed once, in sample HG03729, an Indian Telugu individual from
320    the United Kingdom, sequenced as part of the 1000 Genomes project. Sequence read depth
321    analysis predicts an extra copy of the glycophorin repeat, partly derived from the *GYPB* repeat
322    and partly from the *GYPA* repeat (Figure 4e). The fibre-FISH shows an extra repeat element that
323    is *GYPB*-like at the proximal end and *GYPA*-like at the distal end, and carries the *GYPA* gene.
324    This structure is unlikely to have been generated by a straightforward single NAHR event, and
325    we were unable to resolve the breakpoint at high resolution.
326
327
328    *Breakpoint analysis of structural variants*
329
330    Defining the precise breakpoint of the variants can allow a more accurate prediction of potential
331    phenotypic effects of each variant by assessing, for example, whether a glycophorin fusion gene
332    is formed or whether key regulatory sequences are deleted. We used two approaches to define
333    breakpoints. The first approach identified the two 5kb windows that spanned the change in
334    sequence read depth at both ends of the deletion or duplication, and by designing PCR primers
335    to specifically amplify across the junction fragment (Figure 5a,b), variant-specific PCR
336    amplification produces an amplicon that can be sequenced (Figure 5c). After Sanger sequencing
337    the amplicons, the breakpoint can be shown to be where a switch occurs between paralogous
338    sequence variants (PSVs) that map to different glycophorin repeats (Figure 5d), supporting the
339    model that a NAHR mechanism is responsible for generating these structural variants (Figure
340    5e). The second approach makes use of high depth sequencing. The two 5kb windows spanning
341    the change in sequence read depth are again identified and sequence read depth calculated in
342    1kb windows to further refine the breakpoint. The sequence alignment spanning the two 1kb
343    windows is examined manually for paired sequence reads where the gap between the aligned
344    pairs is consistent with the size of the variant, or where both sequence pairs align but one aligns
345    with multiple sequence mismatches.
346

347    With the exception of DEL4, DUP7 and DUP26, where only low-coverage sequence was

348    available, all other breakpoints could be localised to between 10 kb and 1 bp. For most variants,

349    the breakpoints occur between genes resulting in loss or gain of whole genes, and therefore

350    likely to show gene dosage effect. It is known that DUP4 results in a *GYPB-GYPA* fusion gene

351    that codes for the Dantu blood group, and a fusion gene is also predicted for DUP2, DUP8 and

352    DEL15. The DUP2 variant generates a *GYPB-GYPA* fusion gene comprising exons 1-2 of *GYPB*

353    and exons 4-7 of *GYPA* corresponding to the St$^a$ (GP.Sch) blood group (Anstee et al. 1982;

354    Daniels 2008). Breakpoint analysis of NA12249, the sample carrying the DUP27 variant, showed

355    that DUP27 breakpoint is in the same intron as DUP2, although the exact breakpoint is complex

356    and it is unclear whether DUP27 is exactly the same variant as DUP2.

357

358    The DUP8 variant is predicted to generate a fusion gene consisting of exon 1 of *GYPE* and exons

359    2-7 of *GYPA*, and the DEL15 variant is predicted to combine the first two exons of *GYPB* with

360    the final three exons of *GYPE*. It is unlikely that DUP8 has a phenotype, given the involvement

361    of the 5′ end of *GYPE*, which is not expressed. However, the DEL15 variant is predicted to

362    generate a *GYPB-GYPE* peptide, similar to the rare U- blood group phenotype which has a

363    breakpoint between exon 1 of *GYPB* and exon 2 of *GYPE*, resulting in a lack of expression of

364    glycophorin B in homozygotes (Rahuel et al. 1991). Other variants involve breakpoints within

365    1kb of a gene coding region and could potentially affect expression levels of the neighbouring

366    gene.

367

368    *Mechanism of formation of structural variants*

369

370    The pattern of deletions and duplications observed is consistent with a simple NAHR

371    mechanism of formation for the variants (Figure 5e), with the exception of DUP5 and DUP26.

372    We investigated whether the breakpoints we had found co-localised with known meiotic

373    recombination hotspots previously determined by anti-DMC1 ChIP-Seq of the testes of five

374    males (Pratto et al. 2014). Importantly, the recombination hotspot dataset mapped hotspots in

375    individuals carrying different alleles of the highly-variable PRDM9 protein, a key determinant

376    of recombination hotspot activity, with different alleles activating different recombination

377    hotspots. The glycophorin region contains one hotspot regulated by the PRDM9 A allele,

378    common in Europeans (allele frequency 0.84), and the PRDM9 C allele, common in sub-Saharan

379    Africans (allele frequency 0.13). In our data we found no breakpoints coincident with the

380    PRDM9 A allele hotspot but 4 breakpoints coincident with the PRDM9 C allele hotspot (Figure

381    6). The overlap between the PRDM9 C allele hotspot and the structural variant breakpoints is

382    statistically significant (two-tailed Fisher's exact test, p=0.012) and reflects the observation that

383    there are more different rare structural variants in sub-Saharan African populations, with high

384  frequencies of the C allele, than in European populations where the C allele is almost absent
385  (allele frequency 0.01) (Berg et al. 2011).
386
387

388  **Discussion**
389

390  We have characterised a number of structural variants at the human glycophorin locus. These
391  are almost always large deletions or duplications involving the loss or gain of one or
392  occasionally two glycophorin repeat regions of about 120kb. These losses and gains have been
393  generated by non-allelic homologous recombination (NAHR) mutational events, with particular
394  involvement of the PRDM9 C allele, which is at appreciable frequencies in African populations
395  and directs high recombination rates at its cognate recombination hotspots. A more complex
396  variant, termed DUP5, was also characterised, and was shown to be an inversion-duplication
397  generated by at least 1 microhomology-mediated end-joining event involving DNA sequence
398  outside the glycophorin repeat region. Similarly, DUP26 is a complex variant that is unlikely to
399  have been generated by a single NAHR event.
400

401  Only DEL1, DEL2 and DUP2 are frequent variants. Both DEL1 and DEL2 delete the *GYPB* gene
402  and it is tempting to speculate that their high frequency in African populations and populations
403  with African admixture is due to selection. However, the absence of evidence for any protective
404  effect against malaria argues against malaria being the cause of this selection, so this remains
405  speculation. DUP2 is at notable frequencies particularly in East Asia, and is predicted to
406  generate a *GYPB-GYPA* fusion gene corresponding to the St$^a$ blood group, which is known to be
407  at appreciable frequencies in East Asia (Madden et al. 1964). In this region, malaria infections
408  are caused by *Plasmodium falciparum* as well as *Plasmodium vivax*; alternatively, this fusion gene
409  may facilitate glycophorins acting as a decoy receptor for other pathogens, such as hepatitis A
410  virus (Sanchez et al. 2004). Previous work suggests that DUP2 has arisen on multiple haplotype
411  backgrounds (Leffler et al. 2017), which suggests a large East Asian population panel is need for
412  future accurate imputation.
413

414  Other variants seem either to be geographically localised (for example DUP5) or very rare and
415  detected as singletons in our dataset. This suggests that analysis of other large genomic datasets
416  will discover further unique glycophorin structural variants, and that much glycophorin
417  structural variation is individually rare but collectively more frequent, leading to challenges in
418  imputing glycophorin SV from SNP GWAS data.
419

420  In contrast to other studies, we used a three-step approach to determine copy number. We used
421  read counts over the whole glycophorin region to detect samples with duplications (more than

422   expected number of reads) and deletions (fewer than expected number of reads). We then used

423   window-based analysis of sequence read depth and paralogue-specific allele-specific PCR and

424   Sanger sequencing to refine copy number breakpoints. Finally, we validated the structure of

425   selected variants using fibre-FISH. Our approach has the advantage that it does not rely on a

426   sudden change in sequence read depth for CNV detection by a HMM, which may be

427   compromised by poor mappability of some sequence reads in the breakpoint region and

428   assumptions about the absence of somatic variation, with the consequence that the expected

429   copy number reflecting an integer value. However, our approach cannot detect small copy

430   number changes because the increase or decrease in mapped reads is very small as a proportion

431   of the total number of mapped reads at the glycophorin region.

432

433   Previous work has shown that the DUP4 variant carried by the sample HG02554 shows somatic

434   mosaicism, leading to the suggestion that somatic mosaicism may be a feature of glycophorin

435   structural variants (Algady et al. 2018). In this study, our fiber-FISH analyses identified no other

436   potential somatic variants at the glycophorin locus, showing that it is not a common feature of

437   1000 Genomes lymphoblastoid cell lines nor of non-DUP4 variants. This suggests that somatic

438   mosaicism is either restricted to DUP4 variants in general or restricted to the particular DUP4

439   sample HG02554, although a more thorough investigation of high coverage genome sequences

440   will be needed to address this issue.

441

442   In conclusion, we identify 9 new structural variants at the glycophorin locus, characterise

443   breakpoints and mutational mechanisms for known and novel structural variants, and show

444   that recombination hotspot activity has influenced the nature of the structural variants

445   observed. For some of the variants, targeted high coverage sequence using very long reads

446   analysis will help refine some of the breakpoints. Further efforts are needed to characterise the

447   phenotypic effects of particular variants involving gain, loss and fusion of genes

448

449

450   **Acknowledgements**

454

455

456  **Figure Legends**

457

458

459

460

461

462

463

464

465

466

467

468

469  **Table 1**         **Glycophorin structural variants identified in this study**

| Variant | Proximal breakpoint hg19 | Distal breakpoint hg19 | Variant size (kb) | Breakpoint maximum region (kb) | Index Sample | Genes Involved | Method | In Leffler |
|---|---|---|---|---|---|---|---|---|
| DEL1 | chr4:144835143 -144835279 | chr4:144945375- 144945517 | 110 | 0.143 | NA19223 | *GYPB* | *PCR-Sanger* | *Yes* |
| DEL2 | chr4:144912872 -144913001 | chr4:145016127- 145016256 | 103 | 0.130 | NA19144 | *GYPB* | *PCR-Sanger* | *Yes* |
| DEL4 | chr4:144750739 -144760739 | chr4:144950739- 144960739 | 200 | 10 | HG01986 | *GYPB,GYPE* | *1000G Seq* | *Yes* |
| DEL6 | chr4:144780045 -144780137 | chr4:145004120- 145004212 | 224 | 0.093 | HG04039 | *GYPE and GYPB* | *PCR-Sanger* | *Yes* |
| DEL7 | chr4:144780111 -144780497 | chr4:144900945- 144901334 | 121 | 0.390 | HG02716 | *GYPE* | *PCR-Sanger* | *Yes* |
| DEL13 | chr4:144925739 -144935739 | chr4:145035739- 145045739 | 110 | 10 | NA20867 | *GYPA-B fusion* | *1000G Seq* | *Yes* |
| DEL15 | chr4:144800739 144802739 | chr4:144920739- 144922739 | 119 | 2 | HGDP0117 2 | *GYPB/E fusion gene* | *Deep seq* | *No* |
| DEL16 | chr4:144752739 -144754739 | chr4:144952739- 144954739 | 200 | 2 | BR1296010 301 | *GYPE and GYPB* | *Deep seq* | *No* |
| DEL17 | chr4:144882739 -144987739 | chr4:144984739-- 144987739 | 103 | 3 | BR1183605 501 | *GYPB* | *Deep seq* | *No* |
| DEL18 | chr4:144755739 -144757739 | chr4:144875739- 144878739 | 123 | 2 | BR1099223 302 | *GYPE* | *Deep seq* | *No* |

| DUP2 | chr4: 145039739- | chr4: 144919739-144921739 | 120 | 2 | NA18593 | *GYPB/A fusion* | *PCR-Sanger* | *Yes* |
|---|---|---|---|---|---|---|---|---|
| DUP3 | chr4:145004465 -145004526 | chr4:144780388-144780449 | 224 | 0.062 | NA19360 | *GYPB, GYPE* | *PCR-Sanger* | *Yes* |
| DUP4 | Multiple | Multiple | n/a | n/a | HG02554 | *GYPB/A fusion gene,* | *Leffler et al.* | *Yes* |
| DUP5 | Multiple, including chr4:145113700 | Multiple, including chr4:144936865 | n/a | 0.001 | HG02585 | *GYPB, GYPE* | *PCR-Sanger* | *Yes* |
| DUP7 | chr4:144895000 -144905000 | chr4:144775000-144785000 | 120 | 10 | HG02679 | *GYPB* | *1000G Seq* | *Yes* |
| DUP8 | chr4:145045739-145048739 | chr4:144808739-144810739 | 240 | 3 | I1_S_Irula1, HG03837 | *GYPB, GYPE/A fusion* | *Deep Seq* | *Yes* |
| DUP14 | chr4:144853613 -144853688 | chr4:144723019-144723094 | 131 | 0.076 | NA18646 | *GYPE* | *PCR-Sanger* | *Yes* |
| DUP22 | chr4:144926739 -144929739 | chr4:144881739-144884739 | 45 | 3 | BR210800138, | *GYPB (partial)* | *DeepSeq* | *Yes* |
| DUP26 | chr4:145065739 -145075739 | chr4:144830739-144840739 | 155 | 10 | HG03729 | *GYPA* | *1000G Seq* | *Yes* |
| DUP27 | chr4: 145039739- | chr4: 144919739-144921739 | 120 | 2 | NA12249 | *GYPB/A fusion* | *PCR-Sanger* | *Yes* |
| DUP29 | chr4:144939393 -144939452 | chr4:144825584-144825643 | 114 | 0.060 | HG03686 | *GYPE and GYPB* | *PCR-Sanger* | *No* |
| DUP30 | chr4 144989739- | chr4 144885739-144887739 | 102 | 2 | HGDP00543 | *GYPB* | *Deep seq* | *No* |
| DUP33 | chr4:144959739 -144962739 | chr4:144849739-144851739 | 111 | 3 | BR54409051 | *GYPB* | *DeepSeq* | *No* |
| DUP34 | chr4:145002739 -145004739 | chr4:144900739-144902739 | 102 | 2 | BR1086675791 | *GYPB* | *DeepSeq* | *No* |
| DUP35 | chr4:144878739 -144880739 | chr4:144758739-144760739 | 120 | 2 | BR981404021 | *GYPE* | *DeepSeq* | *No* |

470

471 Notes: DUP19 (NA19223), DUP25 (HG02031), DUP28 (NA19084) no clear 5kb window pattern,

472 DEL4 and DEL16, and DUP2 and DUP27 share overlapping breakpoint regions and may be the

473 same variants. DUP23 (HG02491) and DUP24 (hg03837), identified by Leffler et al, share

474 population and breakpoint regions with DUP8 and are classified as DUP8.

475

476

477   **Table 2**              **Frequency of structural variants observed more than once across the**

478   **cohorts analysed**

479

| | 1000 Genomes | | | | | Gambian | Simons | Brazilian |
|---|---|---|---|---|---|---|---|---|
| Population | EUR | AFR | SAS | EAS | AMR | AFR | ALL | AMR |
| Total chromosomes | 600 | 640 | 386 | 606 | 258 | 782 | 546 | 2648 |
| DEL1 | 0 | 53 | 0 | 1 | 1 | 55 | 7 | 19 |
| DEL2 | 0 | 26 | 0 | 0 | 2 | 2 | 4 | 12 |
| DEL4/16 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 3 |
| DEL6 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| DEL7 | 0 | 2 | 0 | 0 | 0 | 0 | 1 | 0 |
| DUP2/27 | 0 | 1 | 1 | 11 | 1 | 0 | 0 | 7 |
| DUP3 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 |
| DUP5 | 0 | 1 | 0 | 0 | 0 | 2 | 0 | 0 |
| DUP7 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 2 |
| DUP8 | 0 | 0 | 4 | 0 | 0 | 0 | 1 | 2 |
| DUP29 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| DUP22 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |
| DUP30 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 |
| DUP35 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |

480

481

482

483

484

485 **Supplementary table 1**

486

| Variant | Primer name | Primer Sequences 5′ - 3′ | Annealing temperature °C |
|---------|-------------|--------------------------|--------------------------|
| DEL1 | *GYP*_DEL1_F | CCAGTTGCCTCTAAGTCCAT[C] | 65 |
| | *GYP*_DEL1_R | GCAGTGCACACCCTGG[A] | |
| DEL2 | *GYP*_DEL2_F | AGGCAAAAGCTGAGGTCTT[C] | 65 |
| | *GYP*_DEL2_R | CAGCCTCTGGTAACCACTGTTA[C] | |
| DEL6 | *GYP*_DEL6_F | GAAGAAAGAGCTAATTCCAT[G] | 63 |
| | *GYP*_DEL6_R | AGTTGGAACTTGCAAACTTA[G] | |
| DEL7 | *GYP*_DEL7_F | ATCCTGCACTAGAAATTCCTCCCA[C] | 65 |
| | *GYP*_DEL7_R | GATCAGAAAAGCAAAATGGGGC[A] | |
| DEL13 | *GYP*_DEL13_F | CCCTCACCCACAGAAAGAAC[C] | 62 |
| | *GYP*_DEL13_R | GGAAGGTTTTAGAAGTCTTCAGTTG[G] | |
| DUP2 | *GYP*_DUP2_F | CAGAGAAATGATGGGCAAGTTGT[A] | 62 |
| | *GYP*_DUP2_R | ACTGCGTGGACATAGAGCGTAT[T] | |
| DUP3 | *GYP*_DUP3_F | CAAATGAAGTCAAACATCTTC[A] | 63.5 |
| | *GYP*_DUP3_R | CTTGAGACACTCCTTTATATGCTA[C] | |
| DUP5 | *GYP*_DUP5_F | AGCTTGGATGAGATAAATGTCC[T] | 65 |
| | *GYP*_DUP5_R | ATTGGATTCTGATGTGCGG[C] | |
| DUP14 | *GYP*_DUP14_F | GTCTTTAAAGTATTGTTTCGTGC[A] | 65 |
| | *GYP*_DUP14_R | AGGTTAATCTAAAACTTTAGAGCAA[C] | |
| DUP29 | *GYP*_DUP29_F | GCTGCCAGATCAATAGC[G] | 64 |
| | *GYP*_DUP29_R | TAGTAGTATAAACCACAGTGCCTC[A] | |

487

488 Nucleotides that are linked nucleic acids are shown in square brackets.

489

# References

Algady W, Louzada S, Carpenter D, Brajer P, Farnert A, Rooth I, Ngasala B, Yang F, Shaw MA, Hollox EJ. 2018. The Malaria-Protective Human Glycophorin Structural Variant DUP4 Shows Somatic Mosaicism and Association with Hemoglobin Levels. *Am J Hum Genet* **103**: 769-776.

Anstee DJ, Mawby WJ, Parsons SF, Tanner MJ, Giles CM. 1982. A novel hybrid sialoglycoprotein in Sta positive human erythrocytes. *J Immunogenet* **9**: 51-55.

Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, Abecasis GR. 2015. A global reference for human genetic variation. *Nature* **526**: 68-74.

Band G, Le QS, Clarke GM, Kivinen K, Hubbart C, Jeffreys AE, Rowlands K, Leffler EM, Jallow M, Conway DJ et al. 2019. New insights into malaria susceptibility from the genomes of 17,000 individuals from Africa, Asia, and Oceania. *bioRxiv* doi:10.1101/535898: 535898.

Barbosa AR, Souza JM, Lebrao ML, Laurenti R, Marucci Mde F. 2005. Functional limitations of Brazilian elderly by age and gender differences: data from SABE Survey. *Cad Saude Publica* **21**: 1177-1185.

Berg IL, Neumann R, Sarbajna S, Odenthal-Hesse L, Butler NJ, Jeffreys AJ. 2011. Variants of the protein PRDM9 differentially regulate a set of human meiotic recombination hotspots highly active in African populations. *Proc Natl Acad Sci U S A* **108**: 12378-12383.

Cantsilieris S, Western PS, Baird PN, White SJ. 2014. Technical considerations for genotyping multi-allelic copy number variation (CNV), in regions of segmental duplication. *BMC genomics* **15**: 329.

Clarke L, Fairley S, Zheng-Bradley X, Streeter I, Perry E, Lowy E, Tasse AM, Flicek P. 2017. The international Genome sample resource (IGSR): A worldwide collection of genome variation incorporating the 1000 Genomes Project data. *Nucleic Acids Res* **45**: D854-D859.

Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, Aerts J, Andrews TD, Barnes C, Campbell P et al. 2009. Origins and functional impact of copy number variation in the human genome. *Nature* **464**: 704-712.

Daniels G. 2008. *Human blood groups*. John Wiley & Sons.

Gribble SM, Wiseman FK, Clayton S, Prigmore E, Langley E, Yang F, Maguire S, Fu B, Rajan D, Sheppard O. 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.

Handsaker RE, Van Doren V, Berman JR, Genovese G, Kashin S, Boettger LM, McCarroll SA. 2015. Large multiallelic copy number variations in humans. *Nat Genet* **47**: 296-303.

Hehir-Kwa JY, Marschall T, Kloosterman WP, Francioli LC, Baaijens JA, Dijkstra LJ, Abdellaoui A, Koval V, Thung DT, Wardenaar R et al. 2016. A high-quality human reference panel reveals the complexity and distribution of genomic structural variants. *Nat Commun* **7**: 12989.

Hollox EJ, Hoh B-P. 2014. Human gene copy number variation and infectious disease. *Human Genetics*: 1-17.

Huddleston J, Eichler EE. 2016. An Incomplete Understanding of Human Genetic Variation. *Genetics* **202**: 1251-1254.

Jeffreys AJ, Neumann R, Wilson V. 1990. Repeat unit sequence variation in minisatellites: a novel source of DNA polymorphism for studying variation and mutation by single molecule analysis. *Cell* **60**: 473-485.

Kariuki SN, Marin-Menendez A, Introini V, Ravenhill BJ, Lin Y-C, Macharia A, Makale J, Tendwa M, Nyamu W, Kotar J et al. 2018. Red blood cell tension controls Plasmodium

540    falciparum invasion and protects against severe malaria in the Dantu blood group.
541        *bioRxiv* doi:10.1101/475442: 475442.
542  Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:
543        improvements in performance and usability. *Mol Biol Evol* **30**: 772-780.
544  Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler, David. 2002.
545        The Human Genome Browser at UCSC. *Genome Research* **12**: 996-1006.
546  Lan T, Lin H, Zhu W, Laurent T, Yang M, Liu X, Wang J, Wang J, Yang H, Xu X et al. 2017.
547        Deep whole-genome sequencing of 90 Han Chinese genomes. *Gigascience* **6**: 1-7.
548  Latorra D, Campbell K, Wolter A, Hurley JM. 2003. Enhanced allele-specific PCR discrimination
549        in SNP genotyping using 3' locked nucleic acid (LNA) primers. *Hum Mutat* **22**: 79-85.
550  Leffler EM, Band G, Busby GBJ, Kivinen K, Le QS, Clarke GM, Bojang KA, Conway DJ, Jallow
551        M, Sisay-Joof F et al. 2017. Resistance to malaria through structural variation of red
552        blood cell invasion receptors. *Science* **356**.
553  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R.
554        2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078-
555        2079.
556  Li W, Cowley A, Uludag M, Gur T, McWilliam H, Squizzato S, Park YM, Buso N, Lopez R. 2015.
557        The EMBL-EBI bioinformatics web and programmatic tools framework. *Nucleic Acids*
558        *Res* **43**: W580-584.
559  Louzada S, Komatsu J, Yang F. 2017. Fluorescence in situ hybridization onto DNA fibres
560        generated using molecular combing. In *Fluorescence In Situ Hybridization (FISH)*
561        *Application Guide*,  (ed. T Liehr, B Heidelberg), pp. 275-293. Springer-Verlag.
562  Ma H, Difazio S. 2008. An efficient method for purification of PCR products for sequencing.
563        *Biotechniques* **44**: 921-923.
564  MacDonald JR, Ziman R, Yuen RK, Feuk L, Scherer SW. 2014. The Database of Genomic
565        Variants: a curated collection of structural variation in the human genome. *Nucleic Acids*
566        *Res* **42**: D986-992.
567  Madden HJ, Cleghorn TE, Allen FH, Jr., Rosenfield RE, Mackeprang M. 1964. A NOTE ON
568        THE RELATIVELY HIGH FREQUENCY OF ST-A ON THE RED BLOOD CELLS OF
569        ORIENTALS, AND REPORT OF A THIRD EXAMPLE OF ANTI-ST-A. *Vox Sang* **9**: 502-
570        504.
571  Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, Zhao M, Chennagiri N, Nordenfelt
572        S, Tandon A et al. 2016. The Simons Genome Diversity Project: 300 genomes from 142
573        diverse populations. *Nature* **538**: 201-206.
574  Ndila CM, Uyoga S, Macharia AW, Nyutu G, Peshu N, Ojal J, Shebe M, Awuondo KO, Mturi N,
575        Tsofa B et al. 2018. Human candidate gene polymorphisms and risk of severe malaria in
576        children in Kilifi, Kenya: a case-control association study. *Lancet Haematol* **5**: e333-
577        e345.
578  Onda M, Kudo S, Fukuda M. 1994. Genomic organization of glycophorin A gene family revealed
579        by yeast artificial chromosomes containing human genomic DNA. *J Biol Chem* **269**:
580        13013-13020.
581  Pratto F, Brick K, Khil P, Smagulova F, Petukhova GV, Camerini-Otero RD. 2014. DNA
582        recombination. Recombination initiation maps of individual human genomes. *Science*
583        **346**: 1256442.
584  Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
585        features. *Bioinformatics* **26**: 841-842.
586  Raffield LM, Ulirsch JC, Naik RP, Lessard S, Handsaker RE, Jain D, Kang HM, Pankratz N,
587        Auer PL, Bao EL et al. 2018. Common alpha-globin variants modify hematologic and
588        other clinical phenotypes in sickle cell trait and disease. *PLoS Genet* **14**: e1007293.
589  Rahuel C, London J, Vignal A, Ballas SK, Cartron JP. 1991. Erythrocyte glycophorin B
590        deficiency may occur by two distinct gene alterations. *Am J Hematol* **37**: 57-58.

591    Sanchez G, Aragones L, Costafreda MI, Ribes E, Bosch A, Pinto RM. 2004. Capsid region
592        involved in hepatitis A virus binding to glycophorin A of the erythrocyte membrane. *J
593        Virol* **78**: 9807-9813.
594    Sekar A, Bialas AR, de Rivera H, Davis A, Hammond TR, Kamitaki N, Tooley K, Presumey J,
595        Baum M, Van Doren V et al. 2016. Schizophrenia risk from complex variation of
596        complement component 4. *Nature* **530**: 177-183.
597    Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y, Ye K,
598        Jun G, Hsi-Yang Fritz M et al. 2015. An integrated map of structural variation in 2,504
599        human genomes. *Nature* **526**: 75-81.
600    Telenti A, Pierce LC, Biggs WH, di Iulio J, Wong EH, Fabani MM, Kirkness EF, Moustafa A,
601        Shah N, Xie C et al. 2016. Deep sequencing of 10,000 human genomes. *Proc Natl Acad
602        Sci U S A* **113**: 11901-11906.
603    Thorvaldsdóttir H, Robinson JT, Mesirov JP. 2012. Integrative Genomics Viewer (IGV): high-
604        performance genomics data visualization and exploration. *Briefings in Bioinformatics* **14**:
605        178-192.
606    Usher CL, McCarroll SA. 2015. Complex and multi-allelic copy number variation in human
607        disease. *Brief Funct Genomics* **14**: 329-338.
608    Vignal A, Rahuel C, London J, Cherif Zahar B, Schaff S, Hattab C, Okubo Y, Cartron JP. 1990.
609        A novel gene member of the human glycophorin A and B gene family. Molecular cloning
610        and expression. *Eur J Biochem* **191**: 619-625.
611    Zarrei M, MacDonald JR, Merico D, Scherer SW. 2015. A copy number variation map of the
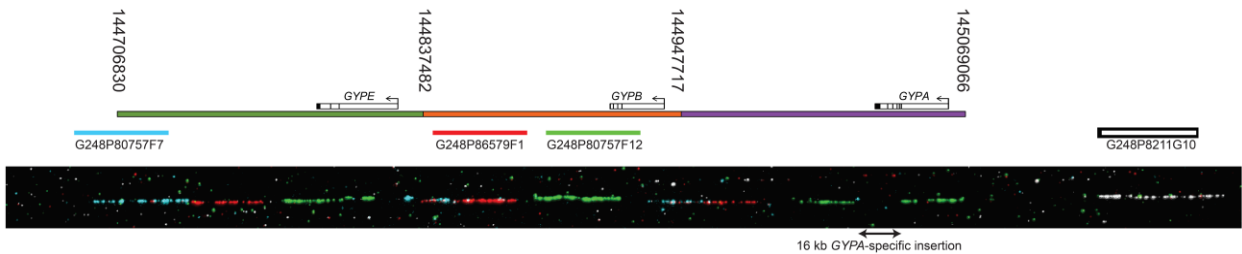612        human genome. *Nat Rev Genet* **16**: 172-183.

613

614

**Figure 1**          **Structure of the glycophorin reference allele**

A representation of the reference allele assembled in the GRCh37/hg19 assembly is shown, with the three distinct paralogous ~120kb repeats of the glycophorin region coloured green, orange and purple, carrying *GYPE*, *GYPB* and *GYPA* respectively. Numbers over the start and end of each paralogue represent coordinates in chromosome 4 GRCh37/hg19 assembly. Coloured bars represent fosmids used as probes in fibre-FISH, with the fosmid identification number underneath. An example fibre FISH image of this reference haplotype (from sample HG02585) is shown below.
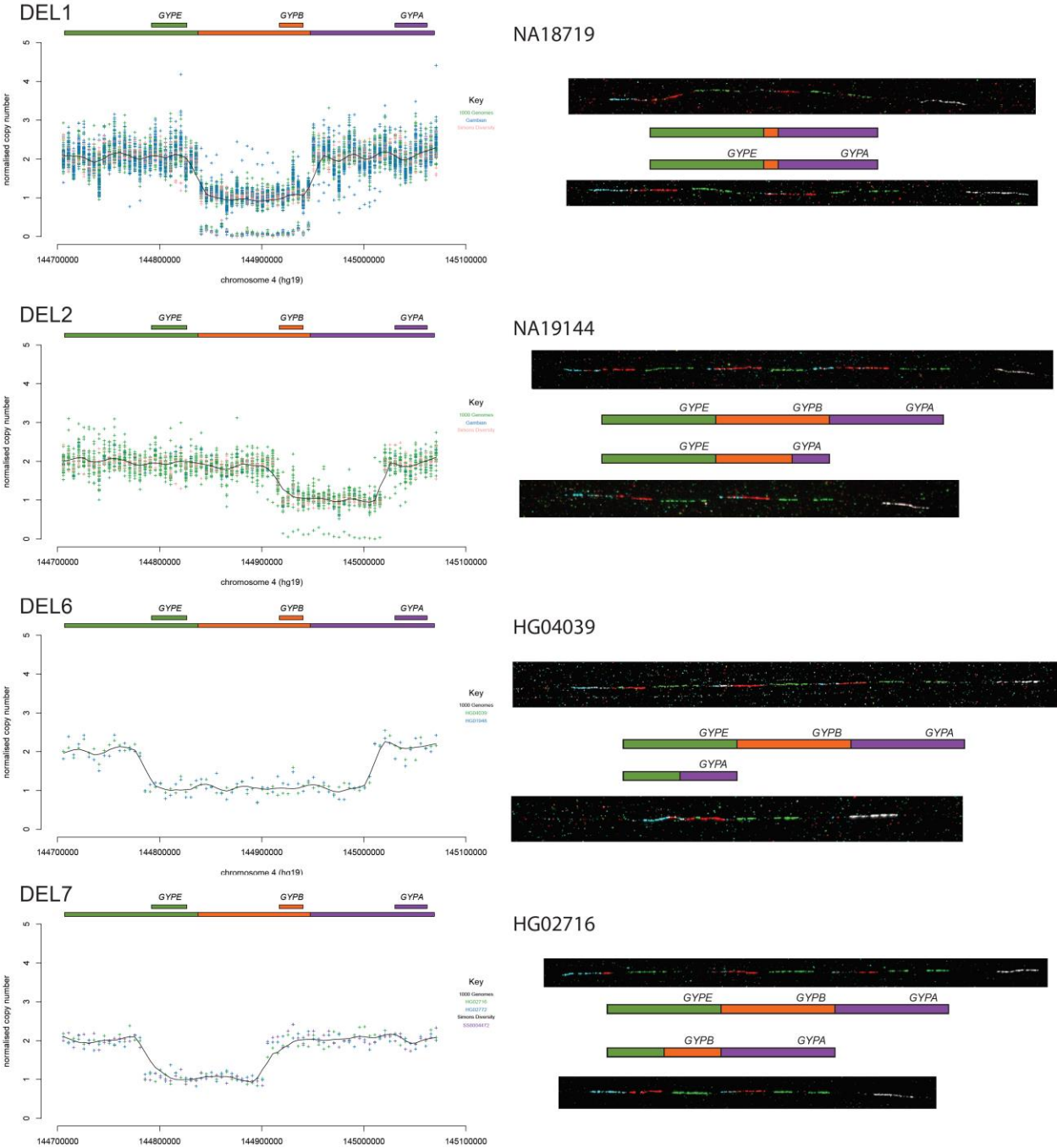
626
627

628

**Figure 2**        **Fibre-FISH validation of four glycophorin deletions**

Sequence read depth (SRD) analysis of selected deletions (DEL1, DEL2, DEL6, DEL7) is shown on the left, with sample, or cohort, coloured according to the key on each plot, together with a Loess best-fit line Note for DEL1 and DEL2 homozygous individuals are detected with a

634     normalised copy number of zero across the deletion. Representative fibre-FISH images from the

635     index sample of each variant are shown on the right, with clones and fluorescent labels as

636     shown in figure 1. All index samples apart from NA18719 are heterozygous, with a

637     representative reference (top) and variant (bottom) allele from that sample shown. A schematic

638     diagram next to the corresponding fibre-FISH image shows the structure of each allele inferred

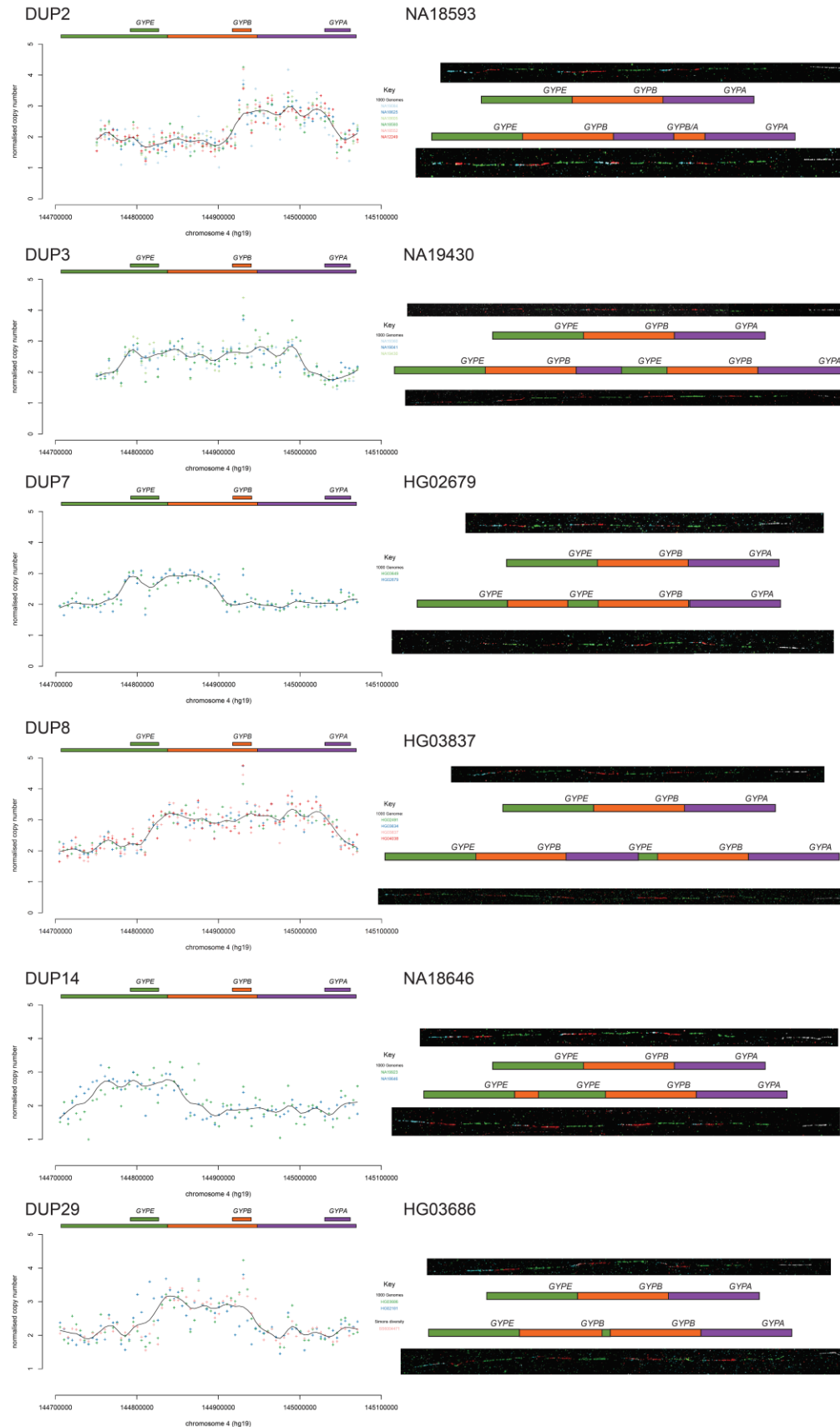639     from the fibre-FISH and SRD analysis.

640

642

643    **Figure 3        Fibre-FISH validation of six glycophorin duplications**

644

645    Sequence read depth (SRD) analysis of selected duplications (DUP2, DUP3, DUP7, DUP8,

646    DUP14 and DUP29) is shown on the left, with sample, or cohort, coloured according to the key

647    on each plot, together with a Loess best-fit line. Representative fibre-FISH images from the

648    index sample of each variant are shown on the right, with clones and fluorescent labels as

649    shown in figure 1, (with an additional green-labelled PCR product specific to the glycophorin E

650    repeat for HG03686). All index samples are heterozygous, with a representative reference and

651    variant allele from that sample shown. A schematic diagram next to the corresponding fibre-

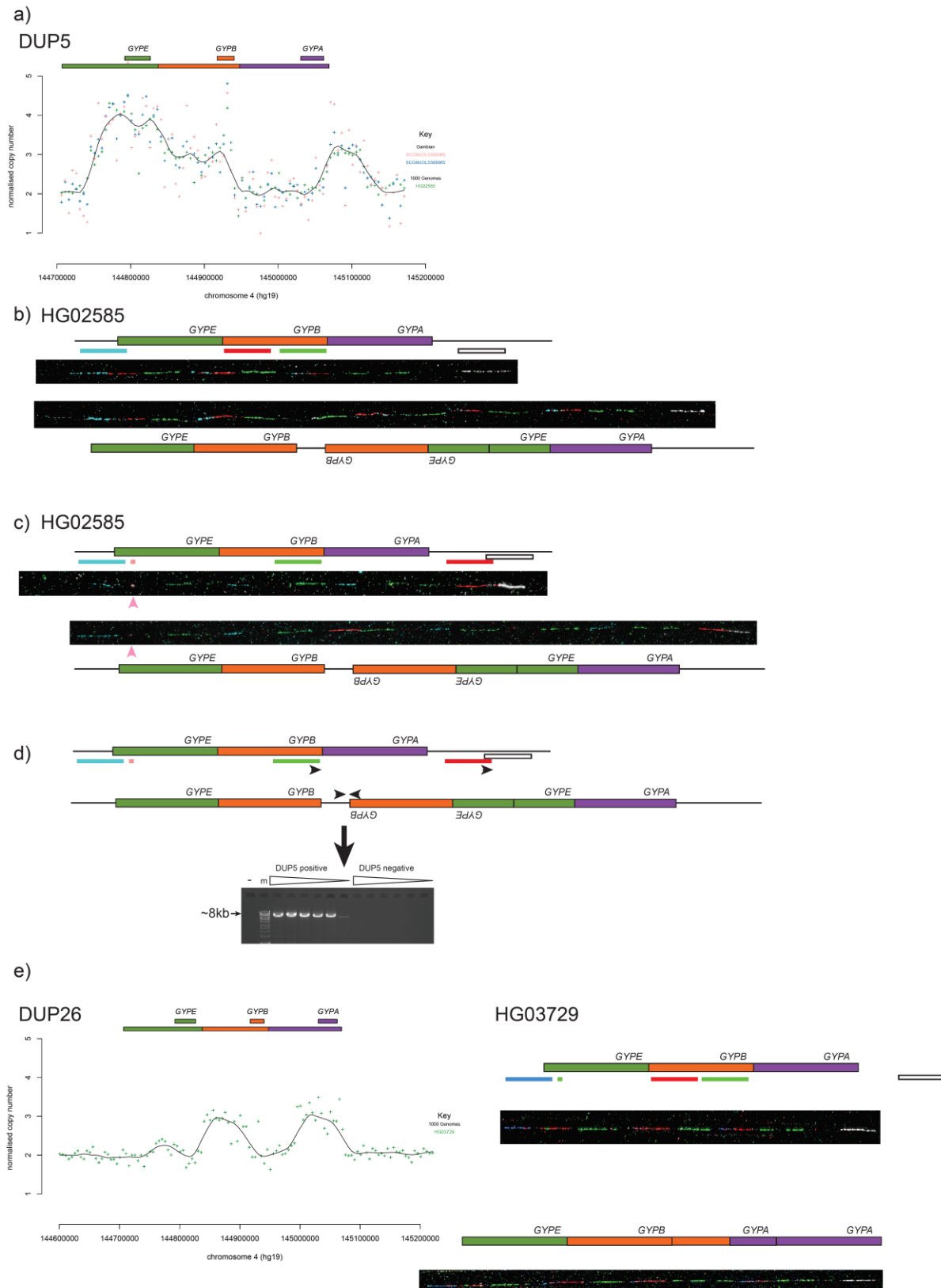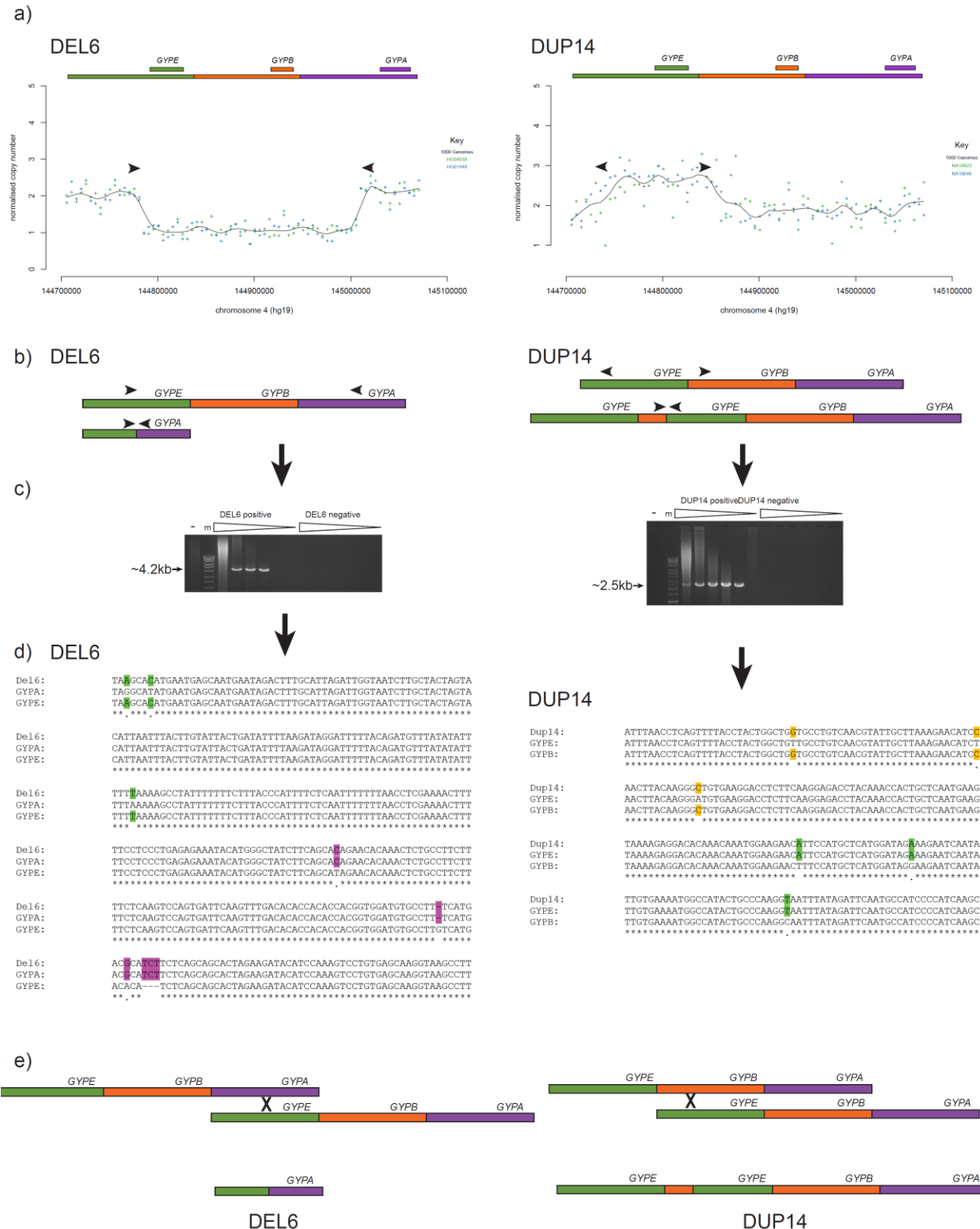652    FISH image shows the structure of each allele inferred from the fibre-FISH and SRD analysis.

653

655     **Figure 4**        **Analysis of DUP5 and DUP26 complex structures**

656

657     a) Sequence read depth (SRD) analysis of three individuals heterozygous for the DUP5
658        variant.
659     b) Representative fibre-FISH images from the DUP5 index sample HG02585. Clones and
660        fluorescent labels as shown in figure 1.
661     c) Representative fibre-FISH images from the DUP5 index sample HG02585. Clones and
662        fluorescent labels as shown in figure 1, except the red probe is fosmid G248P89366H1
663        and the pink probe is the glycophorin E repeat-specific PCR product.
664     d) Schematic showing design of PCR primers for specific amplification (black arrows) on
665        reference and DUP5 structures. The ethidium bromide stained agarose gel shows a ~8kb
666        PCR product generated by these DUP5 specific primers. HG02554 is the DUP5 sample,
667        "-" indicates a negative control with no genomic DNA and the marker, indicated by
668        "m", is Bioline Hyperladder 1kb+. The triangles indicate increasing PCR annealing
669        temperature from 65°C to 67°C.
670     e) Sequence read depth (SRD) analysis (left) and fibre-FISH analysis (right) of the index
671        sample HG03729 heterozygous for DUP26 variant. Fosmid clones for fibre-FISH are as
672        figure 1, except with the addition of the glycophorin E repeat-specific PCR product
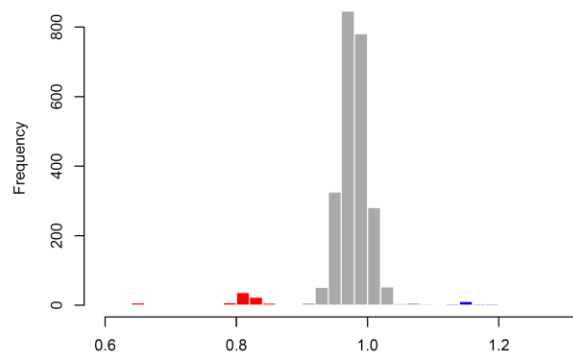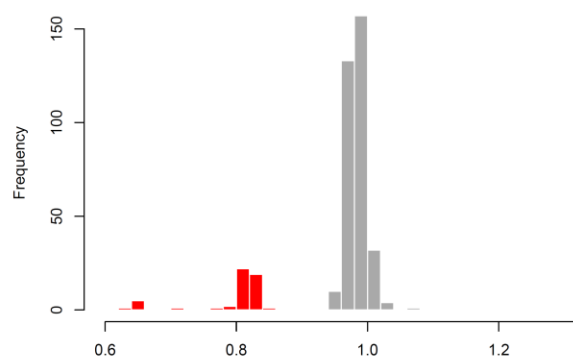673        labelled in green.

674

675
676
677

678     **Figure 5          Examples of refining breakpoints of a deletion (DEL6) and a duplication**
679     **(DUP14)**
680
681     a)  Sequence read depth analysis, indicating position of PCR primers (not to scale)
682     b)  Variant model, showing position of primers on reference and variant
683     c)  Agarose electrophoresis of long PCR products using variant-specific primers indicated
684         in b). "-" indicates a negative control with no genomic DNA and the marker, indicated
685         by "m", is Bioline Hyperladder 1kb+. The triangles indicate increasing PCR annealing
686         temperature from 58°C to 67°C.
687     d)  Multiple sequence alignment of the variant-specific PCR product, with homologous
688         sequence on the *GYPA* repeat and the *GYPE* repeat. *GYPE*-specific variants are in green,
689         *GYPA*-repeat-specific variants are in purple.
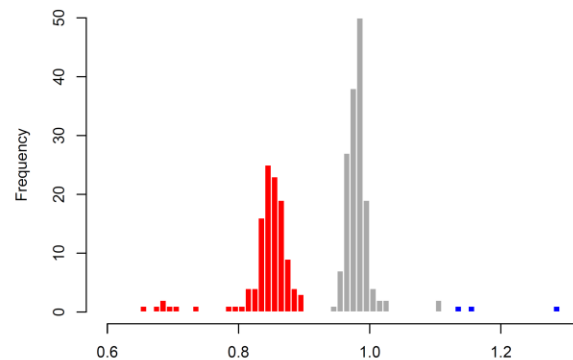690     e)  A model of the generation of the variants by NAHR.
691

692

693

**Figure 6**           **Structural variant breakpoints and meiotic recombination hotspots**

The glycophorin region is shown together with the glycophorin genes. Below are the breakpoint regions for each structural variant, labelled in blue for the distal breakpoint in the variant, and red for the proximal breakpoint in the variant. Meiotic double strand break hotspots, corresponding to recombination hotspots (Pratto et al. 2014) are shown in orange, labelled the PRDM9 allele responsible for activating that hotspot.
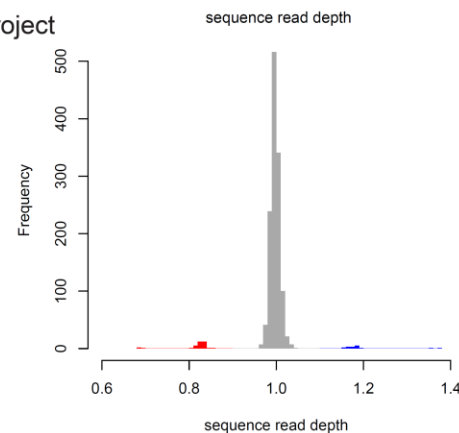
701

a) 1000 Genomes Project

b) Gambian Genomes Project

c) Simons Diversity Project

d) Brazilian Genomes Project



702

703 **Supplementary figure 1** **Histograms of sequence read depths of the glycophorin**
704 **region**

705

706 Histograms of normalised sequence read depths of the four cohorts used for this study, with

707 red indicating putative deletions and blue putative duplications.

708     a) 1000 Genomes Project

709     b) Gambian Genomes Project

710     c) Simons Diversity Project

711     d) Brazilian Genomes Project

712