

# 1 An integrated personal and population-based Egyptian genome 2 reference

3

4 Inken Wohlers, Axel Künstner, Matthias Munz, Michael Olbrich, Anke Fähnrich, Caixia Ma,  
5 Misa Hirose, Shaaban El-Mosallamy, Mohamed Salama, Hauke Busch\* & Saleh Ibrahim\*

6 \* These authors contributed equally to this work

7

8

## 9 Abstract

10

11 The human genome is composed of chromosomal DNA sequences consisting of bases A, C,  
12 G and T – the blueprint to implement the molecular functions that are the basis of every  
13 individual’s life. Deciphering the first human genome was a consortium effort that took more  
14 than a decade and considerable cost. With the latest technological advances, determining an  
15 individual’s entire personal genome with manageable cost and effort has come within reach.  
16 Although the benefits of the all-encompassing genetic information that entire genomes  
17 provide are manifold, only a small number of *de novo* assembled human genomes have been  
18 reported to date <sup>1-3</sup>, and few have been complemented with population-based genetic variation  
19 <sup>4</sup>, which is particularly important for North Africans who are not represented in current  
20 genome-wide data sets <sup>5-7</sup>. Here, we combine long- and short-read whole-genome next-  
21 generation sequencing data with recent assembly approaches into the first *de novo* assembly  
22 of the genome of an Egyptian individual. The resulting genome assembly demonstrates well-  
23 balanced quality metrics and comes with high-quality variant phasing into maternal and  
24 paternal haplotypes, which are linked to various gene expression changes in blood. To  
25 construct an Egyptian genome reference, we further assayed genome-wide genetic variation

26 occurring in the Egyptian population within a representative cohort of more than 100  
27 Egyptian individuals. We show that differences in allele frequencies and linkage  
28 disequilibrium between Egyptians and Europeans may compromise the transferability of  
29 European ancestry-based genetic disease risk and polygenic scores, substantiating the need for  
30 multi-ethnic genetic studies and corresponding genome references. The Egyptian genome  
31 reference represents a comprehensive population data set based on a high-quality personal  
32 genome. It is a proof of concept to be considered by the many national and international  
33 genome initiatives underway. More importantly, we anticipate that the Egyptian genome  
34 reference will be a valuable resource for precision medicine targeting the Egyptian population  
35 and beyond.

36

37

### 38 **Main**

39 With the advent of personal genomics, population-based genetics as part of an individual's  
40 genome is indispensable for precision medicine. Currently, genomics-based precision  
41 medicine compares the patients' genetic make-up to a reference genome <sup>9</sup>, a genome model  
42 inferred from individuals of mostly European descent, to detect risk mutations that are related  
43 to disease. However, genetic and epidemiologic studies have long recognized the importance  
44 of ancestral origin in conferring genetic risk for disease. Risk alleles and structural variants  
45 (SVs) <sup>10</sup> can be missing from the reference genome or can have different population  
46 frequencies, such that alternative pathways become disease related in patients of different  
47 ancestral origin, which motivates the establishment of national or international multi-ethnic  
48 genome projects <sup>6,7,11</sup>. At present, there are several population-based sequencing efforts that  
49 aim to map specific variants in the 100,000 genome projects in Asia <sup>12</sup> or England <sup>13</sup>.  
50 Furthermore, large-scale sequencing efforts currently explore population, society and history-  
51 specific genomic variations in individuals in Northern and Central Europe <sup>14,15</sup>, North

52 America <sup>7</sup>, Asia <sup>16,17</sup> and, recently, the first sub-Saharan Africans <sup>18,19</sup>. Nonetheless, there is  
53 still little genetic data available for many regions of the world. In particular, North African  
54 individuals are not adequately represented in current genetic data sets, such as the 1000  
55 Genomes <sup>5</sup>, TOPMED <sup>7</sup> or gnomAD <sup>6</sup> databases. Consequently, imminent health disparities  
56 between different world populations have been noted repeatedly for a decade. <sup>20–23</sup>

57

58 In recent years, several high-quality *de novo* human genome assemblies <sup>1–4</sup> and, more  
59 recently, pan-genomes <sup>8</sup> have extended human sequence information and improved the *de*  
60 *facto* reference genome GRCh38 <sup>9</sup>. Nonetheless, it is still prohibitively expensive to obtain  
61 all-embracing genetic information, such as high-quality *de novo* assembled personal genomes  
62 for many individuals. Indeed, previous genetic studies assess only a subset of variants  
63 occurring in the Egyptian population, e.g., single nucleotide polymorphisms (SNPs) on  
64 genotyping arrays <sup>24,25</sup>, variants in exonic regions via exome sequencing <sup>26</sup> or variants  
65 detectable by short-read sequencing <sup>27,28</sup>.

66

67 In this study, we generated a phased *de novo* assembly of an Egyptian individual and  
68 identified single nucleotide variants (SNVs) and SVs from an additional 109 Egyptian  
69 individuals obtained from short-read sequencing. Those were integrated to generate an  
70 Egyptian genome reference. We anticipate that an Egyptian population genome reference will  
71 strengthen precision medicine efforts that eventually benefit nearly 100 million Egyptians,  
72 e.g., by providing allele frequencies (AFs) and linkage disequilibrium (LD) between variants,  
73 information that is necessary for both rare and common disease studies. Likewise, our  
74 genome will be of universal value for research purposes, since it contains both European and  
75 African variant features. Most genome-wide association studies (GWAS) are performed in  
76 Europeans <sup>29</sup>, but genetic disease risk may differ, especially for individuals of African  
77 ancestry <sup>30</sup>. Consequently, an Egyptian genome reference will be well suited to support recent

78 efforts to include Africans in such genetic studies, for example, by serving as a benchmark  
79 data set for SNP array construction and variant imputation or for fine-mapping of disease loci.

80

81 Our Egyptian genome is based on a high-quality human *de novo* assembly for one Egyptian  
82 individual (see workflow in Suppl. Fig. 1). This assembly was generated from PacBio, 10x  
83 Genomics and Illumina paired-end sequencing data at overall 270x genome coverage (Suppl.  
84 Table 1). For this personal genome, we constructed two draft assemblies, one based on long-  
85 read assembly by an established assembler, FALCON<sup>31</sup>, and another based on the assembly by  
86 a novel assembler, WTDBG2<sup>32</sup>, which has a much shorter run time with comparable accuracy  
87 (cf. Suppl. Fig. 1). Both assemblies were polished using short reads and various polishing  
88 tools. For the FALCON-based assembly, scaffolding was performed, whereas we found that  
89 the WTDBG2-based assembly was of comparable accuracy without scaffolding (cf. dot plots in  
90 Suppl. Figs. 3-4). The WTDBG2-based assembly was selected as the base because it performs  
91 comparable or better according to various quality control (QC) measures (Suppl. Table 2).  
92 Where larger gaps outside centromere regions occurred, we complemented this assembly with  
93 sequence from the FALCON-based assembly (Suppl. Table 3) to obtain a final Egyptian meta-  
94 assembly, denoted as EGYPT (for overall assembly strategy, see Suppl. Fig. 1). We found our  
95 assembly to be comparable to the publicly available assemblies of a Korean<sup>33</sup> and a Yoruba  
96 individual (GenBank assembly accession GCA\_001524155.4, unpublished) with respect to  
97 various QC measures<sup>34</sup> (Table 1). Suppl. Fig. 2 compares the assemblies' NA-values, and  
98 Suppl. Figs. 3-7 show dot plots of alignment with reference GRCh38. We performed repeat  
99 annotation and repeat masking for all assemblies (Suppl. Table 4).

100

101 The meta-assembly was complemented with high-quality phasing information (Suppl. Table  
102 5). EGYPT SNVs and small insertions and deletions (indels) called using short-read  
103 sequencing data were phased using high-coverage 10x linked-read sequencing data. This

104 resulted in 98.99% of variants being phased. Furthermore, nearly all (99.41%) of the genes  
105 with lengths less than 100 kb and more than one heterozygous SNP were phased into a single  
106 phase block.

107

108 Based on the personal Egyptian genome, we constructed an Egyptian population genome by  
109 considering genome-wide SNV AFs in 109 additional Egyptians (Suppl. Table 6). This  
110 approach enabled the characterization of the major allele (i.e., the allele with highest AF) in  
111 the given Egyptian cohort. To accomplish this, we called variants using short-read data of 13  
112 Egyptians sequenced at high coverage and 97 Egyptians sequenced at low coverage. Although  
113 sequence coverage affects variant-based statistics (Suppl. Fig. 8), due to combined  
114 genotyping, most variants could also be called reliably in low coverage samples (Suppl. Fig.  
115 9). Altogether, we called a total of 19,758,992 SNVs and small indels (Suppl. Fig. 10) in all  
116 110 Egyptian individuals (Fig. 1). The number of called variants per individual varied  
117 between 2,901,883 to 3,934,367 and was correlated with sequencing depth (see Suppl. Figs.  
118 8-9). This relationship was particularly pronounced for low coverage samples. The majority  
119 of variants were intergenic (53.5%) or intronic (37.2%) (Suppl. Fig. 11). Only approximately  
120 0.7% of the variants were located within coding exons, of which 54.4% were non-  
121 synonymous and thus cause a change in protein sequence and possibly structure (Suppl. Fig.  
122 12).

123

124 Using short-read sequencing data of 110 Egyptians, we called 121,141 SVs, most of which  
125 were deletions but also included inversions, duplications, insertions and translocations of  
126 various orders of magnitude (Fig. 1, Suppl. Fig. 13-14). Similar to SNVs, the number of SV  
127 calls also varied between individuals (Suppl. Fig. 15) and is slightly affected by coverage  
128 (Suppl. Fig. 16). After merging overlapping SV calls, we obtained an average of 2,773 SVs  
129 per Egyptian individual (Suppl. Table 7, Suppl. Figs. 17-19).

130

131 To characterize the Egyptian population with respect to European and African populations  
132 that have been genotyped within the 1000 Genomes Project <sup>5</sup> (Suppl. Table 8), we used SNVs  
133 and short indels for a genotype-based principal component analysis. According to this  
134 analysis, Egyptians are a genetically homogenous population compared to other populations,  
135 sharing genetic components with both Europeans and sub-Saharan Africans (see Fig. 2 and  
136 Suppl. Figs. 20-32). Thus far, there are no North African populations with high-quality  
137 whole-genome sequencing-based genotype data available, and in the European and sub-  
138 Saharan African populations reported by the 1000 Genomes Project, Egyptians are closest to  
139 the European Tuscany population (see Fig. 2 and Suppl. Figs. 20-32), which has been  
140 previously proposed through genetic studies of ancient Egyptian mummies <sup>35</sup>.

141

142 The mixed European and African ancestry of Egyptians is further supported by mitochondrial  
143 haplogroup assessment from the literature <sup>27</sup> and our own mtDNA sequencing data (overall  
144 n=327). We found that Egyptians have haplogroups most frequently found in Europeans (e.g.,  
145 H, V, T, J, etc.; more than 60%), and many Egyptians also have African (e.g., L with 24.8%)  
146 or Asian/East Asian haplogroups (e.g., M with 6.7%). This indicates that Egyptian genomes  
147 contain genetic components from various major human populations (Suppl. Fig. 33), as has  
148 been shown recently for a few Egyptian individuals in a study that performed genome-wide  
149 admixture analysis of populations from the Arabian Peninsula using SNP arrays <sup>25</sup>.

150

151 In total, we identified 6,599,037 common Egyptian SNVs (minor allele frequency (MAF) >  
152 5%, genotypes in a minimum of 100 individuals), of which 1,198 are population-specific; i.e.,  
153 they are either rare (MAF < 1%) or not detected in any other population in the 1000 Genomes  
154 <sup>5</sup> and gnomAD databases <sup>6</sup> as well as TOPMed <sup>7</sup> (Suppl. Table 9). These numbers are  
155 comparable to population-specific variant numbers reported previously for 1000 Genomes

156 populations<sup>36</sup>. Four SNVs likely have a molecular impact (Suppl. Table 10), indicated by a  
157 CADD<sup>37</sup> deleteriousness score greater than 20. SNP rs143563851 (CADD 24.2) has recently  
158 been identified in 1% of individuals of a cohort of 211 Palestinians in a study that performed  
159 targeted sequencing of blood group antigen synthase GBGT1<sup>38</sup>. SNP rs143614333 (missense  
160 variant in gene CR2, CADD 23.6) is in ClinVar<sup>39</sup>, with three submitters reporting that the  
161 variant is of uncertain clinical significance. Additionally, we obtained 49 variants with no  
162 dbSNP<sup>40</sup> rsID (Suppl. Table 11). These numbers of population-specific SNPs, of which some  
163 likely have an immediate impact on clinical characteristics and diagnostics, indicate  
164 insufficient coverage of the genetic diversity of the world's population for precision medicine  
165 and thus the need for local genome references. To detect a putative genetic contribution of  
166 Egyptian population-specific SNPs towards molecular pathways, phenotypes or disease, we  
167 performed gene set enrichment analysis for all 461 protein-coding genes that were annotated  
168 to population-specific SNPs by Ensembl VEP<sup>41</sup>. Enrichr, a gene list enrichment tool  
169 incorporating 153 gene sets and pathway databases<sup>42</sup>, reports that genes from obesity-related  
170 traits of the GWAS catalog 2019 collection are over-represented (adj. p-value: 1.02E-6; 49 of  
171 804 genes), which might hint at population-specific metabolism regulation that is linked to  
172 body weight.

173

174 Variants that are not protein coding may have a regulatory effect that affects gene and  
175 eventually protein expression. Using blood expression data obtained from RNA sequencing  
176 for the EGYPT assembly individual in conjunction with 10x sequencing-based phased variant  
177 data, we identified genes whose expression differs between maternal and paternal haplotypes  
178 (see Suppl. Fig. 34 for the analysis overview and Suppl. Figs. 35-36 for the results). We  
179 report 1,180 such genes (Suppl. Table 12). Of these, variants contained in haplotypes of 683  
180 genes (58%) have previously reported expression quantitative trait loci (eQTLs) in blood  
181 according to Qtlizer<sup>43</sup>, for 380 genes supported by multiple studies. For 370 genes (31%), the

182 strongest associated blood eQTL SNV is haplotypically expressed, and for 131 genes, the best  
183 eQTL has been previously reported by multiple studies. Concordance of haplotypic  
184 expression with eQTLs indicates that a common variant may affect gene expression;  
185 discordance hints towards a rare variant.

186

187 We investigated the impact of Egyptian ancestry on disease risk by integrating Egyptian  
188 variant data with the GWAS catalog <sup>44</sup>, a curated database of GWAS. According to the  
189 GWAS catalog, most published GWAS are performed in Europeans <sup>29</sup>, and only a single  
190 study has been performed in Egyptians <sup>45</sup> (by one of our groups). Furthermore, only 2% of  
191 individuals included in GWAS are of African ancestry <sup>29</sup>. AFs, LD and genetic architecture  
192 can differ between populations, such that results from European GWAS cannot necessarily be  
193 transferred <sup>30</sup>. This lack of transferability also compromises the prediction of an individual's  
194 traits and disease risk using polygenic scores: such scores are estimated to be approximately  
195 one-third as informative in African individuals compared to Europeans <sup>46</sup>. From the GWAS  
196 catalog, we constructed a set of 4,008 different, replicated, high-quality tag SNPs (i.e., one  
197 strongest associated SNP per locus) from European ancestry GWAS for 584 traits and  
198 diseases. We compared the tag SNPs' AFs and proxy SNPs in the Egyptian cohort (n=110)  
199 and Europeans from 1000 Genomes (n=503) (Suppl. Table 13). Egyptian AFs of tag SNPs are  
200 comparable to European AFs, with a tendency to be lower (Fig. 3a). There are variants  
201 common in Europeans (AF>5%) but rare in Egyptians (AF<5%) (Suppl. Fig. 37). A total of  
202 261 tag SNPs are not present in the Egyptian cohort (~7%), clearly indicating a need to  
203 perform GWAS in non-European populations to further elucidate disease risk conferred by  
204 these loci. We investigated differences in LD structure using an approach that is used for fine-  
205 mapping of GWAS data, which identifies proxy variants (illustrated in Fig. 3c). Proxy  
206 variants are variants correlated with the tag GWAS SNP, i.e., in high LD (here,  $R^2>0.8$ ). The  
207 post-GWAS challenge is the identification of a causal variant from a set of variants in LD (tag



208 SNP and proxy variants). We found that the number of proxy variants was much lower in the  
209 Egyptian cohort (Fig. 3b), likely due to shorter haplotype blocks known from African  
210 populations. This indicates that LD differences between Egyptians and Europeans may  
211 compromise GWAS transferability and European ancestry-based polygenic scores. However,  
212 Egyptian proxy variants are usually included in the larger set of European proxy variants (Fig.  
213 3d). An example is variant rs2075650 (a locus sometimes attributed to gene TOMM40),  
214 which has been linked to Alzheimer's disease in seven GWASs (cf. Suppl. Fig. 38). This tag  
215 SNP has seven proxy variants in Europeans but only two proxy variants in Egyptians. One  
216 European proxy, rs72352238, has also been reported as a GWAS tag SNP, but it is not a  
217 proxy of rs2075650 in Egyptians and may thus fail replication and transfer of GWAS results  
218 from the European to the Egyptian population.

219  
220 With our Egyptian genome reference, it will be possible to perform comprehensive integrated  
221 genome and transcriptome comparisons for Egyptian individuals and shed light on personal as  
222 well as population-wide common genetic variants. Fig. 4 visualizes the various types of data  
223 of this resource in the integrative genomics viewer <sup>47</sup> (IGV). Here, we selected the DNA  
224 repair-associated gene BRCA2, which, if mutated, is linked to breast cancer and other cancer  
225 types <sup>48</sup>. IGV depicts the sample coverage based on sequencing data from PacBio, 10x  
226 Genomics and Illumina (whole genome as well as RNA) for the personal EGYPT genome  
227 together with common Egyptian SNPs. Variants previously assessed in a breast cancer  
228 GWAS <sup>48</sup> are displayed as Manhattan plot; note the three significant GWAS SNPs between  
229 positions 32,390 and 32,400 kb. The bottom compares the identified SNVs and indels from  
230 the Korean and Yoruba *de novo* assembly with our *de novo* EGYPT assembly. Visual  
231 inspection already yields significantly different variants. This integrative view sheds light on  
232 both small and structural variations at the personal and population-based genome levels.

233

234 In conclusion, we constructed the first Egyptian – and North African – genome reference,  
235 which is an essential step towards a comprehensive, genome-wide knowledge base of the  
236 world’s genetic variations. The wealth of information it provides can be immediately utilized  
237 to study in-depth personal genomics and common Egyptian genetics and its impact on  
238 molecular phenotypes and disease. This reference will pave the way towards a better  
239 understanding of the Egyptian, African and global genomic landscape for precision medicine.

240

## 241 **Methods**

242

### 243 **Sample acquisition**

244 Samples were acquired from 10 Egyptian individuals. For nine individuals, high-coverage  
245 Illumina short-read data were generated. For the assembly individual, high-coverage short-  
246 read data were generated as well as high-coverage PacBio data and 10x data. Furthermore, we  
247 used public Illumina short-read data from 100 Egyptian individuals from Pagani *et al.* <sup>27</sup>. See  
248 Supplementary Tables 1 and 6 for an overview of the individuals and the corresponding raw  
249 and result data generated in this study.

250

### 251 **PacBio data generation**

252 For PacBio library preparation, the SMRTbell DNA libraries were constructed following the  
253 manufacturer’s instructions (Pacific Bioscience, [www.pacb.com](http://www.pacb.com)). The SMRTbell DNA  
254 libraries were sequenced on the PacBio Sequel and generated 298.2 GB of data.

255 Sequencing data from five PacBio libraries were generated at overall 99x genome coverage.

256

### 257 **Illumina short-read data generation**

258 For 350 bp library construction, the genomic DNA was sheared, and fragments with sizes of  
259 approximately 350 bp were purified from agarose gels. The fragments were ligated to

260 adaptors and amplified using PCR. The generated libraries were then sequenced on the  
261 Illumina HiSeq X Ten using PE150 and generated 312.8 GB of data.

262 For the assembly individual, sequencing data from five libraries was generated at overall 90x  
263 genome coverage. For nine additional individuals, one library each was generated, amounting  
264 to an overall 305x coverage of sequencing data. For the 100 individuals of Pagani *et al.*<sup>27</sup>,  
265 three were sequenced at high coverage (30x) and 97 at low coverage (8x). The average  
266 coverage over SNV positions for all 110 samples is provided in Supplementary Table 6.

267

### 268 **RNA sequencing data generation**

269 For RNA sequencing, ribosomal RNA was removed from total RNA, double-stranded cDNA  
270 was synthesized, and then adaptors were ligated. The second strand of cDNA was then  
271 degraded to generate a directional library. The generated libraries with insert sizes of 250-300  
272 bp were selected and amplified and then sequenced on the Illumina HiSeq using PE150.  
273 Overall, 64,875,631 150 bp paired-end sequencing reads were generated.

274

### 275 **10x sequencing data generation**

276 For 10x genomic sequencing, the Chromium Controller was used for DNA indexing and  
277 barcoding according to the manufacturer's instructions (10x Genomics,  
278 [www.10xgenomics.com](http://www.10xgenomics.com)). The generated fragments were sheared, and then adaptors were  
279 ligated. The generated libraries were sequenced on the Illumina HiSeq X Ten using PE150  
280 and generated 272.7 GB of data. Sequencing data from four 10x libraries was generated at  
281 overall 80x genome coverage.

282

### 283 **Construction of draft *de novo* assemblies and meta-assembly**

284 We used WTDBG2<sup>32</sup> for human *de novo* assembly followed by its accompanying polishing  
285 tool WTPOA-CNS with PacBio reads and in a subsequent polishing run with Illumina short

286 reads. This assembly was further polished using PILON<sup>49</sup> with short-read data (cf. Suppl.  
287 Methods: *WTDBG2-based assembly*).

288 An alternative assembly was generated by using FALCON<sup>50</sup>, QUIVER<sup>51</sup>, SSPACE-  
289 LONGREAD<sup>52</sup>, PBJELLY<sup>53</sup>, FRAGSCAFF<sup>54</sup> and PILON<sup>49</sup> (cf. Suppl. Methods:  
290 *FALCON-based assembly*).

291 Proceeding from the WTDBG2-based assembly, we constructed a meta-assembly. Regions  
292 larger than 800 kb that were not covered by this base assembly and were not located within  
293 centromere regions were extracted from the alternative FALCON-based assembly (Suppl.  
294 Table 3). See Suppl. Fig. 1 for an overview of our assembly strategy, including meta-  
295 assembly construction (cf. Suppl. Methods: *Meta-assembly construction*).

296 Assembly quality and characteristics were assessed with QUAST-LG<sup>55</sup> (cf. Suppl. Methods:  
297 *Assembly comparison and QC*). The extraction of coordinates for meta-assembly construction  
298 was performed using QUAST-LG output.

299

### 300 **Repeatmasking**

301 Repeatmasking was performed by using REPEATMASKER<sup>56</sup> with RepBase version 3.0  
302 (Repeatmasker Edition 20181026) and Dfam\_consensus (<http://www.dfam-consensus.org>)  
303 (cf. Suppl. Methods: *Repeat annotation*).

304

### 305 **Phasing**

306 Phasing was performed for the assembly individual's SNVs and short indels obtained from  
307 combined genotyping with the other Egyptian individuals, i.e., based on short-read data.  
308 These variants were phased using 10x data and the 10x Genomics LONGRANGER WGS  
309 pipeline with four 10x libraries provided for one combined phasing. See Supplementary  
310 Methods *Variant phasing* for details.

311

## 312 **SNVs and small indels**

313 Calling of SNVs and small indels was performed with GATK 3.8<sup>57</sup> using the parameters of  
314 the best practice workflow. Reads in each read group were trimmed using Trimmomatic<sup>58</sup>  
315 and subsequently mapped against reference genome hg38 using BWA-MEM<sup>59</sup> version 0.7.17.  
316 Then, the alignments for all read groups were merged sample-wise and marked for duplicates.  
317 After the base recalibration, we performed variant calling using HaplotypeCaller to  
318 obtain GVCF files. These files were input into GenotypeGVCFs to perform joint genotyping.  
319 Finally, the variants in the outputted VCF file were recalibrated, and only those variants that  
320 were flagged as “PASS” were kept for further analyses. We used FastQC<sup>60</sup>, Picard  
321 Tools<sup>61</sup> and verifyBamId<sup>62</sup> for QC (cf. Suppl. Methods: *Small variant QC*).

322

## 323 **Variant annotation**

324 Variant annotation was performed using ANNOVAR<sup>63</sup> and VEP<sup>41</sup> (cf. Suppl. Methods: *Small*  
325 *variant annotation*)

326

## 327 **Structural variants**

328 SVs were called using DELLY2<sup>64</sup> with default parameters as described on the DELLY2  
329 website for germline SV calling (<https://github.com/dellytools/delly>) (cf. Suppl. Methods:  
330 *Structural variant QC*). Overlapping SV calls in the same individual were collapsed by the  
331 use of custom scripts. See Supplementary Methods *Collapsing structural variants* for details.

332

## 333 **Genotype-based principal components**

334 1000 Genomes phase 3 variant data were obtained for all European and African individuals  
335 and merged with the Egyptian variant data. Variants were excluded if their MAF was less  
336 than 5% among individuals in the 1000 Genomes database, they violated Hardy-Weinberg  
337 equilibrium, or they were multi-allelic or within regions of high LD and/or of known

338 inversions. LD pruning was performed, and the remaining SNPs passed on to the SMARTPCA  
339 program <sup>65</sup> of the EIGENSOFT package for PC computation. See Supplementary Methods  
340 *Genotype-based principal components* for details.

341

### 342 **Mitochondrial haplogroups**

343 Haplogroup assignment was performed for 227 individuals using HAPLOGREP 2 <sup>66</sup>.  
344 Furthermore, mitochondrial haplogroups were obtained from Pagani *et al.* <sup>27</sup> for 100  
345 individuals. See Supplementary Methods *Mitochondrial haplogroups* for details.

346

### 347 **Population-specific variants**

348 Our set of common Egyptian SNVs comprises variants with genotypes in a minimum of 100  
349 individuals whose alternative allele has a frequency of more than 5%. Those common  
350 Egyptian SNVs that are otherwise rare, i.e., have an AF of less than 1% in the 1000 Genomes,  
351 and gnomAD populations as well as in TOPMed were considered Egyptian-specific. AFs  
352 were annotated using the Ensembl API. Furthermore, a list of Egyptian common variants  
353 without dbSNP rsID was compiled, see Supplementary Methods *Small variant annotation* for  
354 details.

355

### 356 **Haplotypic expression analysis**

357 RNA sequencing reads were mapped and quantified using STAR (Version 2.6.1.c) <sup>67</sup>.  
358 Haplotypic expression analysis was performed by using PHASER and PHASER GENE AE  
359 (version 1.1.1) <sup>68</sup> with Ensembl version 95 annotation on the 10x-phased haplotypes using  
360 default parameters. See Supplementary Methods *Haplotypic expression* for details.

361

362

363

## 364 **GWAS catalog data integration**

365 GWAS catalog associations for GWAS of European ancestry were split into trait-specific data  
366 sets using Experimental Factor Ontology (EFO) terms. For every trait, a locus was defined as  
367 an associated variant +/- 1 MB, and only loci that were replicated were retained. For proxy  
368 computation, we used our Egyptian cohort (n=110) and the European individuals of 1000  
369 Genomes (n=503). For details, see Supplementary Methods *Data integration with the GWAS*  
370 *catalog*.

371

## 372 **Integrative genomics view**

373 We implemented a workflow to extract all Egyptian genome reference data for view in the  
374 IGV<sup>47</sup>. This includes all sequencing data mapped to GRCh38 (cf. Suppl. Methods *Sequencing*  
375 *read mapping to GRCh38*) as well as all assembly differences (cf. Suppl. Methods *Alignment*  
376 *to GRCh38 and Assembly-based variant identification*) and all Egyptian variant data. See  
377 Supplementary Methods *Gene-centric integrative data views* for details.

378

## 379 **Ethics statement**

380 This study was approved by the Mansoura Faculty of Medicine Institutional Review Board  
381 (MFM-IRB) Approval Number RP/15.06.62. All subjects gave written informed consent in  
382 accordance with the Declaration of Helsinki. This study and its results are in accordance with  
383 the [Jena Declaration](https://www.uni-jena.de/unijenamedia/Universit%C3%A4t/Abteilung+Hochschulkommunikation/Presse/Jenae+r+Erkl%C3%A4rung/Jenaer_Erklaerung_EN.pdf) ([https://www.uni-](https://www.uni-jena.de/unijenamedia/Universit%C3%A4t/Abteilung+Hochschulkommunikation/Presse/Jenae+r+Erkl%C3%A4rung/Jenaer_Erklaerung_EN.pdf)  
384 [jena.de/unijenamedia/Universit%C3%A4t/Abteilung+Hochschulkommunikation/Presse/Jenae](https://www.uni-jena.de/unijenamedia/Universit%C3%A4t/Abteilung+Hochschulkommunikation/Presse/Jenae+r+Erkl%C3%A4rung/Jenaer_Erklaerung_EN.pdf)  
385 [r+Erkl%C3%A4rung/Jenaer\\_Erklaerung\\_EN.pdf](https://www.uni-jena.de/unijenamedia/Universit%C3%A4t/Abteilung+Hochschulkommunikation/Presse/Jenae+r+Erkl%C3%A4rung/Jenaer_Erklaerung_EN.pdf)).

386

## 387 **References**

- 388 1. Cao, H. *et al.* De novo assembly of a haplotype-resolved human genome. *Nat. Biotechnol.*  
389 **33**, 617–622 (2015).

- 390 2. Seo, J.-S. *et al.* De novo assembly and phasing of a Korean human genome. *Nature* **538**,  
391 243–247 (2016).
- 392 3. Shi, L. *et al.* Long-read sequencing and de novo assembly of a Chinese genome. *Nat*  
393 *Commun* **7**, 12065 (2016).
- 394 4. Cho, Y. S. *et al.* An ethnically relevant consensus Korean reference genome is a step  
395 towards personal reference genomes. *Nat Commun* **7**, 13637 (2016).
- 396 5. The 1000 Genomes Project Consortium. A global reference for human genetic variation.  
397 *Nature* **526**, 68–74 (2015).
- 398 6. Karczewski, K. J. *et al.* Variation across 141,456 human exomes and genomes reveals the  
399 spectrum of loss-of-function intolerance across human protein-coding genes. *bioRxiv*  
400 531210 (2019) doi:10.1101/531210.
- 401 7. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed  
402 Program. *bioRxiv* 563866 (2019) doi:10.1101/563866.
- 403 8. Schneider, V. A. *et al.* Evaluation of GRCh38 and de novo haploid genome assemblies  
404 demonstrates the enduring quality of the reference assembly. *Genome Res.* **27**, 849–864  
405 (2017).
- 406 9. Levy-Sakin, M. *et al.* Genome maps across 26 human populations reveal population-  
407 specific patterns of structural variation. *Nature Communications* **10**, 1025 (2019).
- 408 10. Stark, Z. *et al.* Integrating Genomics into Healthcare: A Global Responsibility. *Am. J.*  
409 *Hum. Genet.* **104**, 13–20 (2019).
- 410 11. GenomeAsia 100k. *GenomeAsia 100k* <http://www.genomeasia100k.com/>.
- 411 12. Turnbull, C. *et al.* The 100 000 Genomes Project: bringing whole genome sequencing to  
412 the NHS. *BMJ* **361**, k1687 (2018).
- 413 13. Maretty, L. *et al.* Sequencing and de novo assembly of 150 genomes from Denmark as a  
414 population reference. *Nature* **548**, 87–91 (2017).



- 415 14. Leslie, S. *et al.* The fine-scale genetic structure of the British population. *Nature* **519**,  
416 309–314 (2015).
- 417 15. Chiang, C. W. K., Mangul, S., Robles, C. & Sankararaman, S. A Comprehensive Map of  
418 Genetic Variation in the World’s Largest Ethnic Group-Han Chinese. *Mol. Biol. Evol.* **35**,  
419 2736–2750 (2018).
- 420 16. Bai, H. *et al.* Whole-genome sequencing of 175 Mongolians uncovers population-specific  
421 genetic architecture and gene flow throughout North and East Asia. *Nat. Genet.* (2018)  
422 doi:10.1038/s41588-018-0250-5.
- 423 17. Sherman, R. M. *et al.* Assembly of a pan-genome from deep sequencing of 910 humans of  
424 African descent. *Nat. Genet.* (2018) doi:10.1038/s41588-018-0273-y.
- 425 18. Choudhury, A. *et al.* Whole-genome sequencing for an enhanced understanding of genetic  
426 variation among South Africans. *Nat Commun* **8**, 2062 (2017).
- 427 19. Bustamante, C. D., Burchard, E. G. & De la Vega, F. M. Genomics for the world. *Nature*  
428 **475**, 163–165 (2011).
- 429 20. Popejoy, A. B. & Fullerton, S. M. Genomics is failing on diversity. *Nature* **538**, 161–164  
430 (2016).
- 431 21. Wojcik, G. L. *et al.* Genetic analyses of diverse populations improves discovery for  
432 complex traits. *Nature* **570**, 514–518 (2019).
- 433 22. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health  
434 disparities. *Nat. Genet.* **51**, 584–591 (2019).
- 435 23. Sherman, R. M. *et al.* Assembly of a pan-genome from deep sequencing of 910 humans of  
436 African descent. *Nature Genetics* **51**, 30 (2019).
- 437 24. Henn, B. M. *et al.* Genomic ancestry of North Africans supports back-to-Africa  
438 migrations. *PLoS Genet.* **8**, e1002397 (2012).

- 439 25. Fernandes, V. *et al.* Genome-Wide Characterization of Arabian Peninsula Populations:  
440 Shedding Light on the History of a Fundamental Bridge between Continents. *Mol. Biol.*  
441 *Evol.* **36**, 575–586 (2019).
- 442 26. Scott, E. M. *et al.* Characterization of Greater Middle Eastern genetic variation for  
443 enhanced disease gene discovery. *Nat. Genet.* **48**, 1071–1076 (2016).
- 444 27. Pagani, L. *et al.* Tracing the Route of Modern Humans out of Africa by Using 225 Human  
445 Genome Sequences from Ethiopians and Egyptians. *Am J Hum Genet* **96**, 986–991  
446 (2015).
- 447 28. ElHefnawi, M. *et al.* Whole genome sequencing and bioinformatics analysis of two  
448 Egyptian genomes. *Gene* **668**, 129–134 (2018).
- 449 29. Sirugo, G., Williams, S. M. & Tishkoff, S. A. The Missing Diversity in Human Genetic  
450 Studies. *Cell* **177**, 1080 (2019).
- 451 30. Kim, M. S., Patel, K. P., Teng, A. K., Berens, A. J. & Lachance, J. Genetic disease risks  
452 can be misestimated across global populations. *Genome Biol.* **19**, 179 (2018).
- 453 31. Chin, C.-S. *et al.* Phased diploid genome assembly with single-molecule real-time  
454 sequencing. *Nature Methods* **13**, 1050–1054 (2016).
- 455 32. Ruan, J. & Li, H. Fast and accurate long-read assembly with wtdbg2. *bioRxiv* 530972  
456 (2019) doi:10.1101/530972.
- 457 33. Seo, J.-S. *et al.* De novo assembly and phasing of a Korean human genome. *Nature* **538**,  
458 243–247 (2016).
- 459 34. Mikheenko, A., Prjibelski, A., Saveliev, V., Antipov, D. & Gurevich, A. Versatile  
460 genome assembly evaluation with QUAST-LG. *Bioinformatics* **34**, i142–i150 (2018).
- 461 35. Schuenemann, V. J. *et al.* Ancient Egyptian mummy genomes suggest an increase of Sub-  
462 Saharan African ancestry in post-Roman periods. *Nature Communications* **8**,  
463 ncomms15694 (2017).

- 464 36. Choudhury, A. *et al.* Population-specific common SNPs reflect demographic histories and  
465 highlight regions of genomic plasticity with functional relevance. *BMC Genomics* **15**, 437  
466 (2014).
- 467 37. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD: predicting  
468 the deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* **47**,  
469 D886–D894 (2019).
- 470 38. Abusibaa, W. A. *et al.* Expression of the GBGT1 Gene and the Forssman Antigen in Red  
471 Blood Cells in a Palestinian Population. *Transfusion Medicine and Hemotherapy* (2019)  
472 doi:10.1159/000497288.
- 473 39. Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and supporting  
474 evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
- 475 40. Sherry, S. T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**,  
476 308–311 (2001).
- 477 41. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biology* **17**, 122  
478 (2016).
- 479 42. Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web server  
480 2016 update. *Nucleic Acids Res.* **44**, W90-97 (2016).
- 481 43. Munz, M. *et al.* Qtlizer: comprehensive QTL annotation of GWAS results. *bioRxiv*  
482 495903 (2019) doi:10.1101/495903.
- 483 44. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide  
484 association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**,  
485 D1005–D1012 (2019).
- 486 45. Bejaoui, Y. *et al.* Genome-wide association study of psoriasis in an Egyptian population.  
487 *Exp. Dermatol.* **28**, 623–627 (2019).
- 488 46. Duncan, L. *et al.* Analysis of polygenic risk score usage and performance in diverse  
489 human populations. *Nat Commun* **10**, 3328 (2019).

- 490 47. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).
- 491 48. Michailidou, K. *et al.* Association analysis identifies 65 new breast cancer risk loci.
- 492 *Nature* **551**, 92–94 (2017).
- 493 49. Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection
- 494 and genome assembly improvement. *PLoS ONE* **9**, e112963 (2014).
- 495 50. Chin, C.-S. *et al.* Phased diploid genome assembly with single-molecule real-time
- 496 sequencing. *Nat. Methods* **13**, 1050–1054 (2016).
- 497 51. Chin, C.-S. *et al.* Nonhybrid, finished microbial genome assemblies from long-read
- 498 SMRT sequencing data. *Nat. Methods* **10**, 563–569 (2013).
- 499 52. Boetzer, M. & Pirovano, W. SSPACE-LongRead: scaffolding bacterial draft genomes
- 500 using long read sequence information. *BMC Bioinformatics* **15**, (2014).
- 501 53. English, A. C. *et al.* Mind the Gap: Upgrading Genomes with Pacific Biosciences RS
- 502 Long-Read Sequencing Technology. *PLoS ONE* **7**, e47768 (2012).
- 503 54. Adey, A. *et al.* In vitro, long-range sequence information for de novo genome assembly
- 504 via transposase contiguity. *Genome Research* **24**, 2041–2049 (2014).
- 505 55. Mikheenko, A., Prjibelski, A., Saveliev, V., Antipov, D. & Gurevich, A. Versatile
- 506 genome assembly evaluation with QUAST-LG. *Bioinformatics* **34**, i142–i150 (2018).
- 507 56. SMIT, A. F. A. Repeat-Masker Open-3.0. <http://www.repeatmasker.org> (2004).
- 508 57. McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for
- 509 analyzing next-generation DNA sequencing data. *Genome Research* **20**, 1297–1303
- 510 (2010).
- 511 58. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
- 512 sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
- 513 59. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler
- 514 transform. *Bioinformatics* **26**, 589–595 (2010).

- 515 60. Andrews, S. FASTQC - A quality control tool for high throughput sequence data.  
516 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- 517 61. Picard Toolkit. <http://broadinstitute.github.io/picard/>.
- 518 62. Jun, G. *et al.* Detecting and Estimating Contamination of Human DNA Samples in  
519 Sequencing and Array-Based Genotype Data. *The American Journal of Human Genetics*  
520 **91**, 839–848 (2012).
- 521 63. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic  
522 variants from high-throughput sequencing data. *Nucleic Acids Res* **38**, e164–e164 (2010).
- 523 64. Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end and split-  
524 read analysis. *Bioinformatics* **28**, i333–i339 (2012).
- 525 65. Patterson, N., Price, A. L. & Reich, D. Population structure and eigenanalysis. *PLoS*  
526 *Genet.* **2**, e190 (2006).
- 527 66. Kloss-Brandstätter, A. *et al.* HaploGrep: a fast and reliable algorithm for automatic  
528 classification of mitochondrial DNA haplogroups. *Hum. Mutat.* **32**, 25–32 (2011).
- 529 67. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* bts635 (2012)  
530 doi:10.1093/bioinformatics/bts635.
- 531 68. Castel, S. E., Mohammadi, P., Chung, W. K., Shen, Y. & Lappalainen, T. Rare variant  
532 phasing and haplotypic expression from RNA sequencing with phASER. *Nat Commun* **7**,  
533 12817 (2016).

534

### 535 **Supplementary information**

536 *Supplementary Tables 1-13: An\_Egyptian\_genome\_reference\_supplementary\_tables.xlsx*

537 *Supplementary Methods and Supplementary Figures 1-38:*

538 *An\_Egyptian\_genome\_reference\_supplement.pdf*

539

540

541 **Acknowledgements**

542 We acknowledge support on coordination of the project and assembly work w.r.t the  
543 FALCON-based assembly through Ms. Lu Wang from the Novogene (UK) Company Limited.  
544 HB and IW acknowledge funding by the Deutsche Forschungsgemeinschaft (DFG, German  
545 Research Foundation) under Germany`s Excellence Strategy – EXC 22167-390884018. All  
546 authors acknowledge computational support from the OMICS compute cluster at the  
547 University of Lübeck.

548

549 **Author information**

550

551 *Medical Systems Biology Division, Lübeck Institute of Experimental Dermatology and*  
552 *Institute for Cardiogenetics, University of Lübeck, Lübeck, Germany*

553 Inken Wohlers, Axel Künstner, Matthias Munz, Michael Olbrich, Anke Fähnrich & Hauke  
554 Busch

555

556 *Novogene (UK) Company Limited, Babraham Research Campus, Cambridge, United*  
557 *Kingdom*

558 Caixia Ma

559

560 *Medical Experimental Research Center (MERC), Mansoura University, Mansoura, Egypt*

561 Mohamed Salama & Shaaban El-Mosallamy

562

563 *Institute of Global Health and Human Ecology, The American University in Cairo, Cairo,*  
564 *Egypt*

565 Mohamed Salama

566

567 *Genetics Division, Lübeck Institute of Experimental Dermatology, University of Lübeck,*  
568 *Lübeck, Germany*

569 Misa Hirose & Saleh Ibrahim

570

#### 571 **Contributions**

572 H.B, S.I. and M.S. conceived the study. I.W, A.K, M.M., H.B. and S.I. designed the study.  
573 I.W., A.K., M.M., M.O and A.F. performed data analysis. C.M. constructed the FALCON-  
574 based assembly. M.S. and S.E-M. compiled the Egyptian cohort and provided samples. M.H.  
575 performed mtDNA library preparation and sequencing. I.W., H.B. and S.I. wrote the  
576 manuscript. All authors read and approved the final manuscript.

577

#### 578 **Competing interests**

579 The authors declare no competing interests.

580

#### 581 **Data availability**

582 All summary data of the Egyptian genome reference are available at [www.egyptian-](http://www.egyptian-genome.org)  
583 [genome.org](http://www.egyptian-genome.org). The Egyptian genome reference will be publicly available upon journal  
584 publication.

585

#### 586 **Code availability**

587 Computational tools and parameters used are specified in the Supplementary Methods.  
588 Workflows have been implemented to permit reproducible data analyses by using Snakemake  
589 as workflow management system, Git for version control of workflow code and Conda  
590 (especially Bioconda) for managing software environments.

591

592

593 **Corresponding authors**

594 Correspondence to Hauke Busch or Saleh Ibrahim.

595

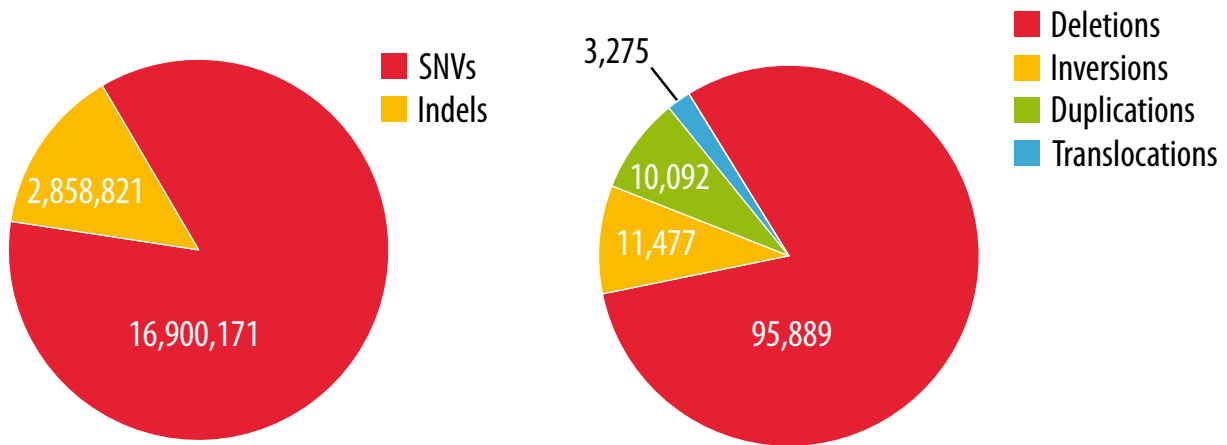


596 *Table 1: Default assembly quality measures according to QUAST-LG. The extended QUAST-LG report is*  
 597 *provided in Suppl. Table 2. Yoruba is a chromosome-level assembly. Best quality for every measure is denoted in*  
 598 *bold.*

Genome statistics	EGYPT	EGYPT_wtdbg2	EGYPT_falcon	AK1	YORUBA
Genome fraction (%)	94.174	92.247	95.924	95.177	<b>95.391</b>
Duplication ratio	1.01	<b>0.999</b>	1.018	1.023	1.088
	20,908	20,613	<b>21,176</b>	21,047	21,077
# genomic features	(3,226 part)	(3,229 part)	<b>(1578 part)</b>	(1,396 part)	(1,721 part)
Largest alignment	<b>75,492,126</b>	<b>75,492,126</b>	56,458,009	58,219,133	65,512,502
Total aligned length	2,800,100,449	2,713,712,375	<b>2,865,356,241</b>	2,829,006,639	2,832,740,986
NGA50	<b>11,187,777</b>	<b>11,187,777</b>	8,226,500	13,028,687	19,529,238
LGA50	71	71	95	66	<b>43</b>
<b>Misassemblies</b>					
# misassemblies	<b>1,276</b>	<b>1,276</b>	3,499	1,952	1,756
Misassembled contigs					
length	<b>2,137,050,584</b>	2,137,050,584	2,851,404,290	2,657,569,650	3,053,643,982
<b>Mismatches</b>					
# mismatches per 100 kbp	139	138.72	143.64	<b>126.92</b>	141.56
# indels per 100 kbp	32.09	<b>31.74</b>	40.06	32.77	46.95
# N's per 100 kbp	<b>0</b>	<b>0</b>	209.01	1285.7	7180.2
<b>Statistics without reference</b>					
# contigs	3,235	3,106	<b>1,615</b>	2,832	1,647
Largest contig	88,566,048	88,566,048	84,324,762	113,921,103	<b>248,986,603</b>
Total length	2,836,714,529	2,750,324,638	2,916,268,178	2,904,207,228	<b>3,088,335,497</b>
Total length (>= 1000 bp)	2,837,367,164	2,750,799,236	2,916,433,762	2,904,207,228	<b>3,088,485,407</b>
Total length (>= 10000 bp)	2,828,723,737	2,742,501,225	2,914,302,309	2,904,207,228	<b>3,086,359,078</b>
Total length (>= 50000 bp)	2,803,817,652	2,718,165,929	2,895,137,452	2,855,011,855	<b>3,059,626,724</b>
<b>K-mer-based statistics</b>					
K-mer-based compl. (%)	86.01	85.15	<b>87.75</b>	87.68	85.82
# k-mer-based misjoins	1,654	1,649	1,786	<b>1,345</b>	1,453

599

600



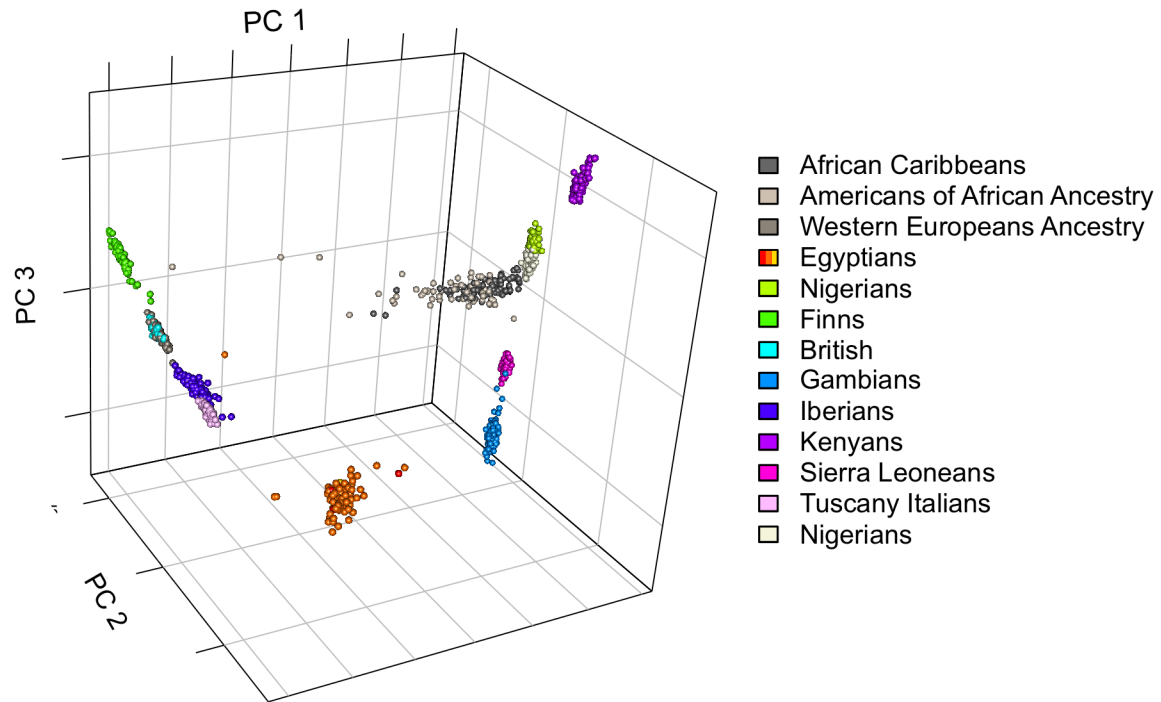
601

602 *Figure 1: Number of various genetic variant types identified in the Egyptian cohort. Left: The number of SNVs*

603 *and indels. Right: The number of SV calls: deletions, inversions, duplications and translocations. Additionally,*

604 *408 insertions have been called.*

605



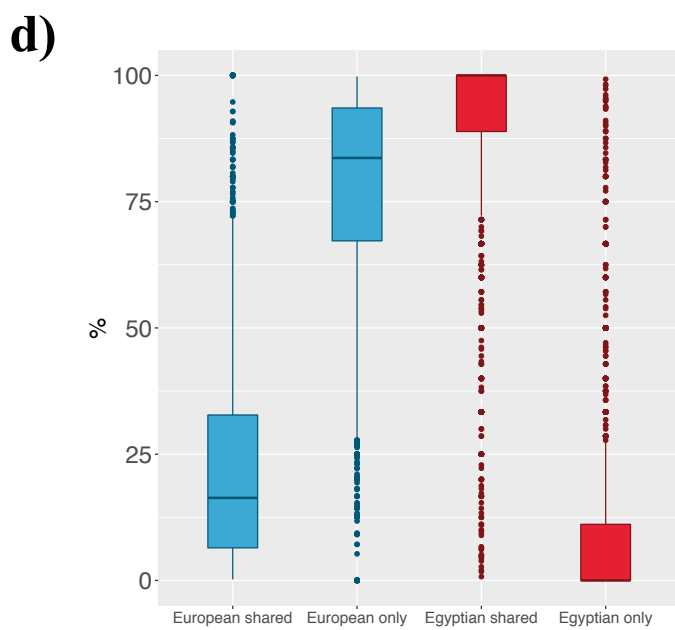
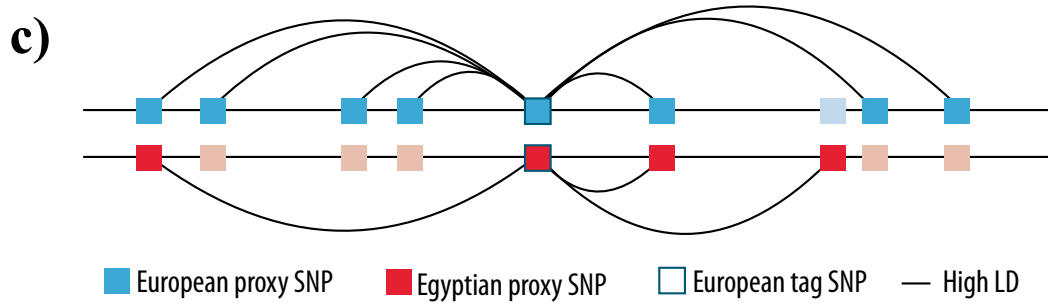
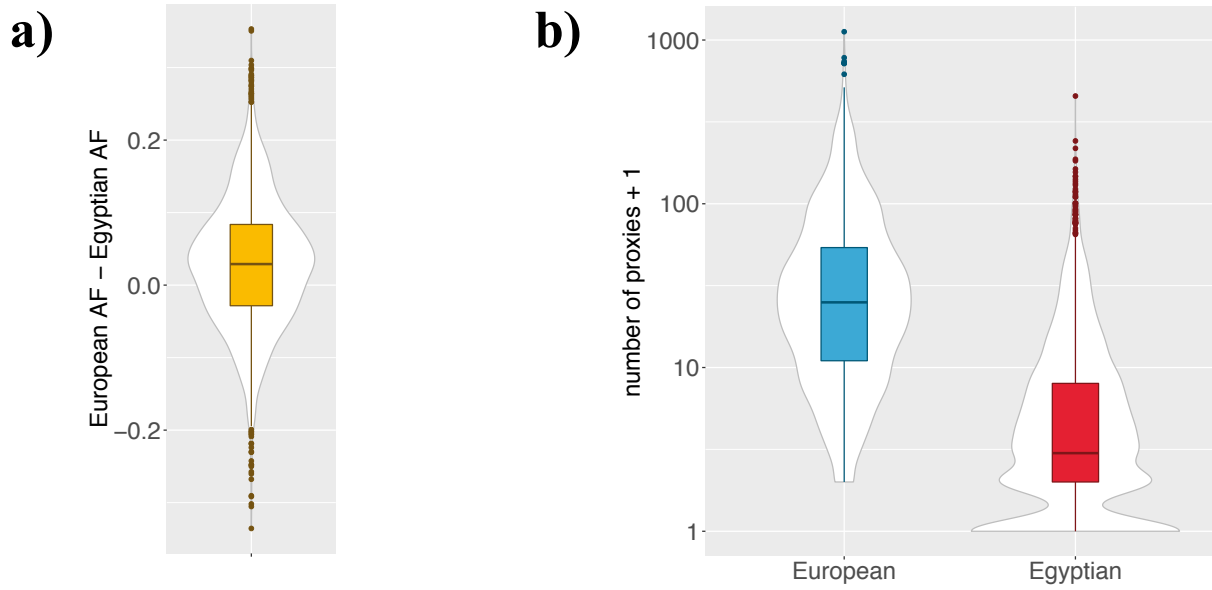
606

607 *Figure 2: PCA plot of different populations from the 1000 Genomes Project and 110 Egyptian genomes from*

608 *Pagani et al. as well as from our own study.*

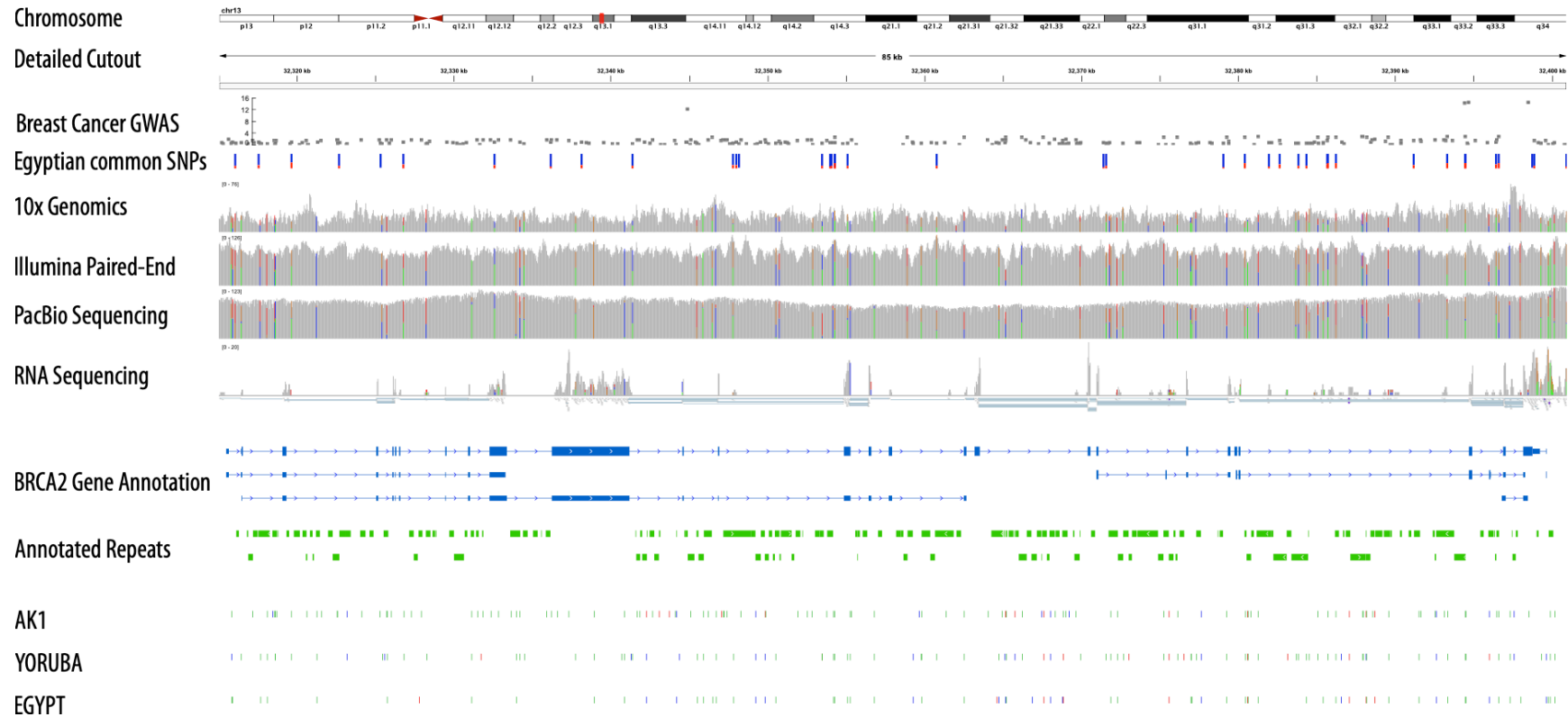
609

610



611

612 *Figure 3: AF and proxy SNP comparisons for 3,698 GWAS tag SNPs called in a minimum of 100 Egyptians. a)*  
613 *AF differences. b) Number of proxies. c) Illustration of the proxy SNP comparison. A European GWAS tag SNP*  
614 *(center) and variants in Europeans (top) and Egyptians (bottom). Lines denote variants in high LD. The tag SNP*  
615 *has 7 proxy variants in Europeans and 3 in Egyptians. Light blue/red variants are no proxy variants in*  
616 *Europeans/Egyptians. Two proxy variants are shared. Thus 2 of 7 European (~29%) and 2 of 3 Egyptian (~67%)*  
617 *variants are shared. Further 5 of 7 European proxies are European-only (~71%) and 1/3 Egyptian proxies are*  
618 *Egyptian-only (~33%). d) European shared: Percentage of European proxy SNPs shared with Egyptian proxy*  
619 *SNPs. European only: Percentage of European proxy SNPs not shared with Egyptian proxies. Egyptian shared /*  
620 *Egyptian only respectively.*  
621



623

624 *Figure 4: Integrative view of Egyptian genome reference data for the gene BRCA2, which is associated with breast cancer. The rows denote from top to bottom: Genome location*

625 *on chromosome 13 of the magnified region for BRCA2 (first and second row); GWAS data for breast cancer risk<sup>48</sup>; Variants that are common in the cohort of 110 Egyptians;*

626 *Read coverage of genetic region based on 10x Genomics, Illumina paired-end and PacBio sequencing data; Coverage and reads of RNA sequencing data; BRCA2 gene*

627 *annotation from Ensembl; Repeats annotated by REPEATMASKER; SNVs and indels identified by comparison of assemblies AK1, YOURUBA and EGYPT with GRCh38. The*

628 *colors denote base substitutions (green), deletions (blue) and insertions (red).*