

1                                   **Construction and Integration of**  
2                                   **Three *De Novo* Japanese Human Genome Assemblies**  
3                                   **toward a Population-Specific Reference**

4  
5  
6   **AUTHORS**

7   Jun Takayama<sup>1,2,3</sup>, Shu Tadaka<sup>2</sup>, Kenji Yano<sup>2,3</sup>, Fumiki Katsuoka<sup>1,2</sup>, Chinatsu Gocho<sup>2</sup>,  
8   Takamitsu Funayama<sup>2</sup>, Satoshi Makino<sup>2</sup>, Yasunobu Okamura<sup>1,2</sup>, Atsuo Kikuchi<sup>4</sup>, Junko  
9   Kawashima<sup>2</sup>, Akihito Otsuki<sup>2</sup>, Jun Yasuda<sup>2</sup>, Shigeo Kure<sup>2,4</sup>, Kengo Kinoshita<sup>1,2,5,\*</sup>,  
10   Masayuki Yamamoto<sup>1,2,\*</sup>, and Gen Tamiya<sup>1,2,3,\*</sup>

11  
12   **AFFILIATIONS**

- 13   **1.** Advanced Research Center for Innovations in Next-Generation Medicine, Tohoku  
14   University, Sendai, Japan.  
15   **2.** Tohoku Medical Megabank Organization, Tohoku University, Sendai, Japan.  
16   **3.** RIKEN Center for Advanced Intelligence Project, Tokyo, Japan.  
17   **4.** Department of Pediatrics, Tohoku University School of Medicine, Sendai, Japan.  
18   **5.** Graduate School of Information Sciences, Tohoku University, Sendai, Japan.

19  
20   **CORRESPONDENCE**

21   Kengo Kinoshita (kengo@ecei.tohoku.ac.jp)  
22   Masayuki Yamamoto (masiyamamoto@med.tohoku.ac.jp)  
23   Gen Tamiya (gtamiya@megabank.tohoku.ac.jp)

25 **ABSTRACT**

26 **The complete sequence of the human genome is used as a reference for next-**  
27 **generation sequencing analyses. However, some ethnic ancestries are under-**  
28 **represented in the international human reference genome (e.g., GRCh37), especially**  
29 **Asian populations, due to a strong bias toward European and African ancestries in**  
30 **a single mosaic haploid genome consisting chiefly of a single donor. Here, we**  
31 **performed *de novo* assembly of the genomes from three Japanese male individuals**  
32 **using >100× PacBio long reads and Bionano optical maps per sample. We integrated**  
33 **the genomes using the major allele for consensus, and anchored the scaffolds using**  
34 **sequence-tagged site markers from conventional genetic and radiation hybrid maps**  
35 **to reconstruct each chromosome sequence. The resulting genome sequence,**  
36 **designated JG1, is highly contiguous, accurate, and carries the major allele in the**  
37 **majority of single nucleotide variant sites for a Japanese population. We adopted**  
38 **JG1 as the reference for confirmatory exome re-analyses of seven Japanese families**  
39 **with rare diseases and found that re-analysis using JG1 reduced false-positive**  
40 **variant calls versus GRCh37 while retaining disease-causing variants. These results**  
41 **suggest that integrating multiple genome assemblies from a single ethnic population**  
42 **can aid next-generation sequencing analyses of individuals originated from the**  
43 **population.**

44

45 **INTRODUCTION**

46 The complete human genome sequence<sup>1,2</sup> has been an invaluable resource for both basic  
47 research in human genetics and clinical diagnosis. The complete genome sequence is  
48 currently used as a reference for mapping the enormous number of short reads generated  
49 using major next-generation sequencing (NGS) techniques<sup>3,4</sup>, is thus also called "the  
50 reference genome". Because the short reads generated in NGS studies are approximately  
51 100–300 bp in length, mapping them to the reference genome is an indispensable step for

52 calling single nucleotide variants (SNVs) and short insertions and deletions (indels) in  
53 the sample individuals. The coordinate system of the reference genome is used for  
54 biological and medical annotations, such as the position or sequence of specific genes, or  
55 sites of causal variants associated with both rare and common diseases. Therefore, the  
56 reference genome is one of the most foundational resources in human genetics, and as  
57 such, it is maintained and continually updated by the Genome Reference Consortium  
58 (GRC). The latest and second-latest versions of the reference genome (GRCh38/hg38 and  
59 GRCh37/hg19, published in 2013 and in 2009, respectively) are nearly complete, and  
60 both are widely used for NGS analyses and genome annotations<sup>5,6</sup>.

61

62 The reference genome was constructed using a hierarchical shotgun sequencing strategy  
63 in which fragmented genomic DNA segments cloned in bacterial (BAC) or P1-derived  
64 (PAC) artificial chromosome libraries are arranged into a correct physical map to  
65 guarantee that the reference genome was haploid (mosaic)<sup>1</sup>. The assembled contigs or  
66 scaffolds were then anchored on each chromosome using information from genetic and  
67 radiation hybrid (RH) maps, which have thousands to tens of thousands of sequence-  
68 tagged site (STS) markers in linkage groups (i.e. chromosomes). It should be noted that  
69 these genetic and RH maps are original information sources used to construct the  
70 reference genome and not derived from the reference genome itself.

71

72 Although the reference genome is a resource of unparalleled value, several of its  
73 characteristics are not ideal for application to NGS analyses, particularly for some  
74 populations<sup>7</sup>. For example, although the reference genome is constructed using genetic  
75 information from multiple donors, each clone comprising the resulting reference genome  
76 is derived from either haploid genome of a particular individual. As such, the reference  
77 genome inevitably harbors rare or even private variants. Over 90,000 rare variants were  
78 used as a reference allele including disease-susceptibility variants for thrombophilia and

79 type 2 diabetes<sup>8,9</sup>. Inclusion of such variants in the reference can lead to erroneous and  
80 confusing results of short read mapping or variant calling<sup>9</sup>. As the NGS analyses typically  
81 assume that the reference allele is the ancestral, healthy, or major allele for any variable  
82 site, the inclusion of such rare alleles may also confuse subsequent interpretations.

83

84 Another possible problem associated with the reference genome is that the samples used  
85 for its construction are biased toward African and European ancestries. For example,  
86 >70% of the reference genome is composed of a BAC library known as RP-11 (aliased  
87 RPCI-11)<sup>1</sup> from a donor with both African and European ancestry<sup>10</sup>. With the exception  
88 of one donor with an Asian background, all of the donors had a European background  
89 resulting in the composition of an Asian haplotype for 4.3% of the reference genome<sup>1,10</sup>.  
90 In addition, one recent study revealed a lack of population-specific sequences in the  
91 reference genome<sup>11</sup>, whereas another discovered thousands of structural variants (SVs) in  
92 world-wide samples<sup>12</sup>. These issues can also complicate short-read mapping and variant  
93 callings.

94

95 Several studies have examined ways to overcome the above-mentioned drawbacks to the  
96 reference genome. Dewey et al<sup>13</sup> proposed modifying the reference genome by  
97 substituting its minor variants with the major variants from African, Asian, or European  
98 populations<sup>13</sup>. The resulting modified reference genome was better-suited for genome  
99 analyses of sample individuals with matched population backgrounds. Several studies<sup>14-</sup>  
100 <sup>17</sup> utilized a genome graph, which is an extended reference genome represented as a graph  
101 harboring known variants. Other studies have proposed the addition of sequences not  
102 included in the reference genome<sup>11,12,18,19</sup>. However, these proposed adjustments are  
103 based largely on variants discovered using the reference genome itself, albeit only  
104 partially, in a circular fashion, some reference bias could remain.

105

106 One promising approach to address these problems is to construct new reference genomes  
107 specific to ethnic populations of interest<sup>20</sup>. Although costly, highly contiguous *de novo*  
108 assembly—independent reconstruction—of the entire human genome is now feasible  
109 using, for example, Pacific Biosciences (PacBio) single molecule, real-time (SMRT) long  
110 reads (~10 kb in length) and Bionano optical mapping, which generates a high-resolution  
111 physical map<sup>18,21,22</sup>. Combining of these approaches is known as 'hybrid scaffolding,'  
112 which is carried out in three steps: 1) PacBio long reads are *de novo* assembled to yield  
113 primary contigs; 2) Bionano raw data are also *de novo* assembled (independent of the  
114 PacBio assembly) to yield optical maps; and 3) the PacBio-derived contigs are scaffolded  
115 by the Bionano optical maps. This strategy is analogous to the hierarchical shotgun  
116 sequencing strategy used in the Human Genome Project<sup>1</sup> with arrangements of long  
117 sequences from BAC/PAC on a physical map. Although assemblies generated in recent  
118 studies were highly contiguous and accurate, the assembled sequences were rarely  
119 anchored to a set of chromosomes, thus making their use as references for NGS analyses  
120 impractical. Moreover, a single haploid assembly from a single individual cannot be used  
121 to solve the rare reference allele problem. A notable exception is the KOREF genome  
122 sequence<sup>20</sup>, in which a Korean reference genome was constructed by *de novo* assembly  
123 of the genome sequence of a Korean individual, reconstructed as a set of chromosomes,  
124 and rare variants were substituted with short reads from 40 Korean individuals. However,  
125 the KOREF genome assembly was found to be less contiguous than long read-based  
126 assemblies because the primary sequencing platform was a short-read sequencer, and  
127 KOREF depended heavily on the reference genome because chromosome building was  
128 carried out by sequence-based alignment of scaffolds onto GRCh38.

129

130 Using a hybrid scaffolding strategy, in this study, we constructed a new reference genome,  
131 JG1, by integrating *de novo* assemblies of three Japanese individuals. After merging the  
132 three haploid assemblies constructed by hybrid scaffolding strategy, we defined major

133 variants among the three (i.e., majority decision) and adopted them as the reference allele.  
134 We also positioned the scaffolds along chromosomes with the aid of conventional genetic  
135 and RH maps. We then assessed the extent to which JG1 represents the major variants in  
136 the Japanese population in terms of SNVs and SVs. As an example potential application,  
137 we also demonstrated the utility of using JG1 as a reference genome in NGS analyses  
138 aimed at identifying the causal variants of several rare diseases.

139

## 140 **RESULTS**

### 141 ***Construction of JG1***

142 To construct a genome sequence with population reference-quality, the population  
143 background of the reference genome should not significantly diverge from the  
144 backgrounds of sample individuals in order to reduce unnecessary variant calls that  
145 merely reflect the difference in the population background. In the case of our study, the  
146 donor should therefore be chosen from the Japanese population originating from the main  
147 island of Japan. In addition, we built the Japanese reference genome independent of the  
148 GRC reference genome in order to eliminate known ethnic biases toward African and  
149 European backgrounds as well as any other (and possibly unknown) biases. We therefore,  
150 performed *de novo* assembly of the Japanese human genome. Majority-based decision-  
151 making regarding multiple *de novo* assemblies was implemented as an effective way to  
152 avoid inclusion of rare reference alleles. This majority-based decision-making strategy  
153 produced a haploid genome sequence amenable to analyses using currently available and  
154 standard bioinformatics tools for NGS data.

155

156 We recruited three male Japanese volunteers, and they were given the sample names jg1a,  
157 jg1b, and jg1c (jg1a is the same individual as JPN00001<sup>19</sup>). Principal component analysis  
158 (PCA) based on the genotypes inferred by whole-genome sequencing indicated that the  
159 subjects were scattered within the cluster of the Japanese population (Figure 1a). G-

160 banding analyses (Supplementary Fig. 1) indicated that all three individuals had a normal  
161 karyotype, although subject jg1a had a common pericentric inversion within chromosome  
162 9, inv(9)(p12q13). Because it was difficult to assemble the pericentric region of  
163 chromosome 9 equally for all three subjects, this variation does not appear to have  
164 affected the assembly results (Supplementary Fig. 2).

165

166 To construct a reference-quality haploid genome sequence, we integrated the three *de*  
167 *novo* assembled genomes (see Supplementary Fig. 3 for an overview; see Supplementary  
168 Tables 1-3 for materials). For each subject, we sequenced deeply (122× for jg1a, 123×  
169 for jg1b, and 128× for jg1c) using PacBio technology (Supplementary Fig. 4 and  
170 Supplementary Table 1) and then performed *de novo* assembly using Falcon software<sup>23</sup>.  
171 The *de novo* assemblies yielded 2,194, 2,227, and 2,120 primary contigs for jg1a, jg1b,  
172 and jg1c, respectively (Table 1). The contig N50 value was approximately 20 Mb for the  
173 three subjects (Table 1). Using ArrowGrid software<sup>24</sup>, the primary contigs were then  
174 error-corrected (polished) with the same long reads used for the initial *de novo* assembly.

175

176 We also obtained deep Bionano data for each subject (123× and 140× for two enzymes  
177 for jg1a; 160× and 175× for one enzyme for jg1b and jg1c, respectively; Supplementary  
178 Fig. 5 and Supplementary Table 2), and performed *de novo* assemblies of these data to  
179 generate optical maps (Supplementary Table 4). Each *de novo* assembly of the Bionano  
180 data was performed in two rounds (rough and full) to guarantee independence relative to  
181 the GRC reference genome (see Methods section). We then performed hybrid scaffolding  
182 between the PacBio-derived contigs and the Bionano-derived optical maps. The resulting  
183 hybrid scaffolds were then polished with 55×, 59×, and 57× Illumina short reads for  
184 subjects jg1a, jg1b, and jg1c, respectively (Supplementary Table 3). The number and N50  
185 value of the resulting hybrid scaffolds were 1,911, 1,893, and 1,797, and 86.28 Mb, 59.38  
186 Mb, and 58.20 Mb for subjects jg1a, jg1b, and jg1c, respectively. These and other

187 assembly statistics were better than or comparable to other published *de novo* assemblies  
188 (Table 1).

189

190 To enhance the quality of our genome assembly, we adopted a meta-assembly strategy in  
191 which multiple assemblies were merged to yield a single assembly. In meta-assembly  
192 strategies, individual assemblies are aligned, and one best assembly is selected for each  
193 aligned segment based on the absence of rare SVs, unresolved sequences, or possible mis-  
194 assembly inferred by other experimental evidence, such as mate-pair sequencing data<sup>25</sup>.  
195 For meta-assembly, Metassembler software<sup>25</sup> was applied to 37× mate-pair short reads  
196 from the three subjects in sum to infer discordance among the individual scaffolds  
197 (Supplementary Table 3). A total of 12 meta-assemblies, or sets of meta-scaffolds, were  
198 generated from the three sets of scaffolds, based on the order and combination of the  
199 processed sets of scaffolds (see Methods section). Among the 12 possible combinations,  
200 we found that one combination (jg1c + (jg1a + jg1b))—which merged the scaffolds of  
201 jg1c with the meta-scaffolds generated from that of jg1a and jg1b in this order—exhibited  
202 no apparent large chimeric mis-assembly in any autosomes. This combination was chosen  
203 for the downstream sophistications; the absence of chimeric mis-assembly was assessed  
204 using STS markers described later. This set of meta-scaffolds exhibited better contiguity  
205 and accuracy than the original set of scaffolds for subject jg1c (Table 1).

206

207 Although meta-scaffolds were more contiguous and accurate than individual sets of  
208 scaffolds, rare reference alleles should still be retained in the meta-scaffolds. To eliminate  
209 these rare reference alleles, we aligned the three individual sets of scaffolds against the  
210 meta-scaffolds, performed variant calling, defined the major allele among the three sets  
211 of scaffolds, and substituted the minor allele on the meta-scaffolds to the major allele both  
212 in terms of SNVs and SVs (Supplementary Fig. 6a). For tri-allelic sites, we chose one  
213 allele randomly among the three as a reference allele. We also found that two assemblies



214 among the three contained a 2.6-Mb inversion in the long arm of chromosome 9  
215 (Supplementary Fig. 2), and we confirmed that the meta-scaffolds also contained the  
216 inversion.

217

218 We next tried to anchor the majority-voted meta-scaffolds on each chromosome. To do  
219 so, we utilized a total of 85,386 distinct STS markers from three genetic maps and six RH  
220 maps pre-dated the reference genome: the Genethon<sup>26</sup>, deCODE<sup>27</sup>, and Marshfield<sup>28</sup>  
221 genetic maps and the Whitehead-RH<sup>29</sup>, GeneMap99-GB4<sup>30</sup>, GeneMap99-G3<sup>30</sup>, Stanford-  
222 G3<sup>31</sup>, NCBI\_RH<sup>32</sup>, and TNG<sup>33</sup> RH maps. We searched for STS marker amplifications by  
223 electronic PCR analysis of the meta-scaffolds and used ALLMAPS software<sup>34</sup> to order  
224 and orient the meta-scaffolds to build chromosomes. The co-linearities between the  
225 anchored meta-scaffolds and genetic and RH maps were  $0.999 \pm 0.004$  and  $0.986 \pm 0.021$ ,  
226 respectively (Pearson's correlation coefficient; mean  $\pm$  SD). However, we found that  
227 ALLMAPS using all nine above-mentioned maps did not assign any meta-scaffolds to  
228 the Y chromosome, probably because most of the maps did not include the Y chromosome.  
229 Nonetheless, we found that ALLMAPS using three of the nine maps (deCODE, TNG,  
230 and Stanford-G3) assigned some meta-scaffolds to the Y chromosome as well as  
231 autosomes and the X chromosome. Therefore, we adopted the ALLMAPS assignment  
232 with the nine maps for autosomal assignment and those with three maps to the sex  
233 chromosomes.

234

235 After anchoring these meta-scaffolds to chromosomes, we found a chimera in the sex  
236 chromosomes. A meta-scaffold harboring the *SRY* locus, a gene on the Y chromosome,  
237 was chimeric and anchored to the long arm of the X chromosome in the selected set of  
238 meta-scaffolds. We therefore chose a set of meta-scaffolds from another set of meta-  
239 scaffolds (jg1a + (jg1b + jg1c)) for the long arm of the X chromosome that had no  
240 apparent chimeric region.

241

242 We also manually modified the length of unresolved regions in the telomeric, centromeric,  
243 and constitutive heterochromatic regions represented as a stretch of Ns (see Methods  
244 section). We then masked a pseudo-autosomal region (PAR) in the Y chromosome to  
245 guarantee that the resulting sets of sequences represented a haploid. In addition, we  
246 shifted the start position of the mitochondrial meta-scaffold to match the revised  
247 Cambridge Reference Sequence (rCRS) coordinates<sup>35</sup>, which provides the reference  
248 coordinate system for the mitochondrial genome.

249

250 The procedure described above yielded a set of chromosome-level sequences for 22  
251 autosomes, 2 sex chromosomes, and 1 mitochondrial chromosome, along with 599  
252 unplaced scaffolds, and we designated this set of sequences JG1 (Figure 1b). The total  
253 length of JG1 was approximately 3.1 Gb (Table 1). Notably, in the JG1 genome assembly,  
254 19 chromosomal arms were successfully represented as single scaffolds (Figure 1b). After  
255 constructing these chromosome-level sequences, we then aligned them to reference  
256 genome GRCh38 using minimap2 software<sup>36</sup> and found an overall high similarity  
257 between the two genomes at the sequence level (Figure 1c). Because JG1 was built  
258 independently from the reference genome GRCh38, this overall high similarity provided  
259 strong support for our approach for building JG1 described above.

260

### 261 ***Representativeness of the JG1 haplotype in terms of SNVs***

262 To assess whether JG1 is a representative reflection of the SNV composition of the  
263 Japanese population, we performed PCA using JG1 and the reference hg19, along with  
264 2,022 haplotypes constructed from 11 HapMap3 populations (see Methods section). The  
265 PCA plot shows three major clusters representing African, Asian, and European  
266 populations (Figure 2a). The JG1 haplotype localized near the cluster of Asian  
267 populations, whereas the hg19 haplotype localized between the African and European

268 populations, as expected based on the donors' ancestries (Figure 2a). Notably, the JG1  
269 haplotype did not localize within the Asian cluster; instead it was the most distant site  
270 both from the European and African populations, suggesting an “Asianness” when  
271 compared with the other two populations. We also performed PCA with JG1 and 505  
272 haplotypes constructed from three Asian populations: Japanese (JPT), Han Chinese  
273 (CHB) and Chinese in Denver (CHD) (Figure 2b). The PCA plot included two distinct  
274 clusters (namely, Japanese and others), with the JG1 haplotype associated with the  
275 Japanese cluster.

276

277 To assess whether JG1 harbors the major allele among the Japanese population across  
278 SNV sites, we first aligned JG1 against the reference genome hs37d5 and detected SNVs.  
279 The genome-by-genome alignment and comparison by minimap2 and paf tools software<sup>36</sup>  
280 called 2,501,575 SNVs between hs37d5 and JG1 in the autosomes and X chromosome.  
281 We then extracted the frequency of the allele employed in JG1 from the allele frequency  
282 (AF) panel of 3,552 Japanese individuals (namely, 3.5KJPNv2 AF panel<sup>37</sup>) to create a site  
283 frequency spectrum, in which the horizontal axis indicates the non-hg19-type allele and  
284 the vertical axis indicates the number of such SNV sites (Figure 2c). From these data, we  
285 found 241,500 SNV sites with an AF = 1.0, indicating that all of the Japanese  
286 chromosomes in the AF panel carried the JG1-type allele at the 241,500 sites. This  
287 corresponds to 97.99% of all such SNV sites that had an AF = 1.0 (246,464) in the  
288 3.5KJPNv2 AF panel. Similarly, we identified 367,271 and 626,254 SNV sites with an  
289  $AF \geq 0.99$  or  $\geq 0.90$ , respectively, corresponding to 97.11% and 96.24% of such SNV  
290 sites in the 3.5KJPNv2 AF panel, respectively (378,211 and 650,718). A peak observed  
291 at an AF of ~0.22 was associated with the SNPs clustered in the XTR region—a region  
292 known to harbor complex duplications—within 88.8 to 92.4 Mb on the X chromosome.  
293 A peak at an AF of approximately zero could most likely be attributed to artificial SNVs  
294 called at the edges of alignments. We also assessed the effectiveness of the majority

295 decision approach. Of the 2,501,575 SNVs, 1,176,922 (47%) and 1,204,762 (48%) were  
296 detected in three and two of the three JG1-donor individuals, respectively (Supplementary  
297 Fig. 6b).

298

### 299 ***Representativeness of the JG1 haplotype in terms of SVs***

300 To investigate differences between JG1 and GRCh38 in terms of SVs, we aligned JG1  
301 against the reference GRCh38 and detected SVs (insertions and deletions) using the  
302 minimap2 and paf tools software programs<sup>36</sup>. A genome-by-genome comparison detected  
303 8,689 insertions and 6,177 deletions >50 bp but <10,000 bp in length. The length  
304 distribution of the SVs exhibited two peaks, at approximately 300 bp and 6 kb (Figure  
305 3a). We confirmed that the 300-bp and 6-kb peaks were associated with *Alu* and LINE1,  
306 respectively. Most of the SV-associated *Alu* and LINE1 were classified as *AluY* and L1HS,  
307 respectively, both of which constitute the currently-active subclass of these transposable  
308 elements (the length distributions of detected transposable elements are shown in  
309 Supplementary Fig. 7). In addition, the detected SVs were often observed in the sub-  
310 telomere–telomere regions (Figure 3b), consistent with a previous report<sup>12</sup>.

311

312 To investigate the extent to which JG1 represents a Japanese population in terms of SVs,  
313 we mapped short reads of 200 Japanese individuals to JG1 and GRCh38 to compare the  
314 average read depth among the 200 individuals around the SVs in JG1 and GRCh38. As  
315 shown in Figure 3c, insertions were typically associated with a 'piling-up' of the average  
316 depth, whereas deletions were typically associated with a depression of the average depth  
317 in GRCh38. Neither pattern was clearly evident in the corresponding region in JG1  
318 (Figure 3c), suggesting that most of the Japanese samples shared the SVs. To determine  
319 whether this pattern is common among SVs throughout the genome, we compared the  
320 maximum difference in average depth between the SV region and its adjacent upstream  
321 region of the same length and found that the difference in the average depth was smaller

322 in JG1 than GRCh38 (Figure 3d;  $P = 7.6 \times 10^{-11}$  for  $n = 3,950$  pair of insertions;  $P < 2.2$   
323  $\times 10^{-16}$  for  $n = 2,763$  pair of deletions; Wilcoxon signed rank tests).

324

### 325 ***Utility of JG1 as a reference for rare disease exome analyses***

326 To evaluate whether JG1 is a suitable reference genome for clinical NGS analyses, we  
327 examined exomes of Japanese families with rare diseases<sup>38</sup>. The sample cohort consisted  
328 of 22 individuals from six trio families and one quartet family. All of the families had one  
329 child affected with diplegia and eight causal variants were identified in previous analyses  
330 using the reference genome hg19 (Table 2). The diseases exhibited autosomal recessive,  
331 compound heterozygous, and autosomal dominant modes of inheritance, including *de*  
332 *nov*o mutations, and the causal variants included both single nucleotide and deletion  
333 variants.

334

335 To facilitate exome re-analysis with JG1, we lifted over the resource bundles of Genome  
336 Analysis ToolKit (GATK) software (which are used for accurate variant calling) and  
337 GENCODE gene annotation information<sup>39</sup> to predict variant effects (see Methods section)  
338 and performed exome analyses according to GATK best practices. The JG1-based exome  
339 analyses correctly identified all (8/8) of the previously reported causal variants. In  
340 addition, the total number of called variants was lower in JG1 than the reference hs37d5  
341 (Figure 4a). This comparison was done in the 225,888 exome regions with one-to-one  
342 correspondence between JG1 and hs37d5 (87,971,409 bp and 87,997,786 bp for JG1 and  
343 hs37d5, respectively). Moreover, the number of both high- and moderate-impact variants  
344 (which are the primary causal variant candidates) was also lower in JG1 than hs37d5  
345 (Figure 4b;  $473 \pm 16$  vs  $671 \pm 13$  high-impact and  $8,774 \pm 97$  vs  $10,599 \pm 89$  moderate-  
346 impact variants for JG1 and hs37d5, respectively; mean  $\pm$  SD). These findings suggest  
347 that JG1 produces fewer false-positives while successfully detecting disease-causing  
348 variants in whole-exome analyses. In addition, we compared the variants detected with

349 JG1 and hs37d5 by lifting over the JG1-detected variants to hs37d5 and found ~15,000,  
350 ~29,000, and ~52,000 specific to JG1, hs37d5, and both references, respectively (Figure  
351 4c). Moreover, we extracted the non-GRC-type AF in the JG1-specific, hs37d5-specific,  
352 and shared variant sites from the 3.5KJPNv2 AF panel and found that most of the hs37d5-  
353 specific variants were major alleles among the Japanese population, whereas the shared  
354 and JG1-specific variants tended to be biased toward the minor AFs (Figure 4d).

355

## 356 **DISCUSSION**

357 Here, we report the first construction of a Japanese haploid genome sequence, JG1, by  
358 integrating three highly contiguous *de novo* hybrid assemblies from three Japanese donor  
359 individuals to build a population-specific (i.e. ethnicity-matched) reference genome.  
360 Employing a meta-assembly approach produced a more-contiguous and accurate  
361 assembly, and relying on majority decision among the three genomes substituted most of  
362 the rare reference alleles. The results of both SNV and SV analyses suggested that the  
363 JG1 haplotype represents major variation among the Japanese population. Moreover, we  
364 demonstrated that JG1 exhibits several advantages as an ethnicity-matched reference for  
365 NGS analyses, at least within the clinical context of whole exomes of Japanese samples.  
366 Using JG1 could thus facilitate detecting the proverbial needle in a haystack, by reducing  
367 the size of the haystack in NGS analyses of the Japanese population.

368

369 Integration and majority decision regarding multiple assemblies to yield a single haploid  
370 genome can produce a highly contiguous and accurate assembly, thus effectively  
371 eliminating most rare reference variations. Haploid representation of the genome is  
372 beneficial because it is compatible with many conventional bioinformatics tools  
373 developed to date for mapping, variant calling, predicting variant effects, and subsequent  
374 interpretations. Although we appreciate that the development of a pan-human genome  
375 graph could be the next milestone reached in comprehensively assessing human genetic

376 variations among diverse populations, we expect that population-specific reference  
377 genome such as JG1 will prove to be practical and beneficial options for genome analyses  
378 of individuals originated from the population.

379

380 Several limitations of the current version of JG1 should be noted: (1) sequence  
381 incompleteness and gaps/un-localized fragments remaining, which could result in  
382 erroneous mapping and variant calling; (2) few original annotations on the JG1  
383 coordinates; and (3) incomplete representation of the major variations in the Japanese  
384 population. The incompleteness of the genome sequence could be largely overcome by  
385 applying other genome sequencing technologies, including ultra-long reads of Oxford  
386 nanopore technologies in combination with targeted cloning from whole-genome BAC  
387 libraries. Chromosome-scale scaffolding using Hi-C<sup>40</sup> could also contribute to the  
388 generation of more-contiguous assemblies. The genome of a Japanese complete  
389 hydatidiform mole, characterized by a duplicated haploid genome, could also contribute  
390 to gap-filling due to ease of assembly<sup>41</sup>. The limitation of few original annotations could  
391 be overcome by constructing an AF panel with JG1 as the reference and by *de novo*  
392 prediction or experimental inference regarding gene regions. More comprehensive  
393 lifting-over of many annotations would also be practically important. The  
394 representativeness of the major alleles would be improved by adding more assemblies.  
395 One approach that could be used for addition is the phased diploid assembly<sup>23</sup>, which  
396 provides a pair of haplotype (i.e., diplotype) assemblies from a single individual. Because  
397 the two haploid genomes can be regarded as a random sample from a panmictic  
398 population, assembling two haploid genomes per individual can increase the  
399 representativeness of variations in the population. Despite its limitations above, the  
400 current version of JG1 represents a useful tool for efficient causal variant detection.

401

402 Additional samples, for example hundreds of samples from a single population, would be

403 beneficial for constructing population-specific reference genomes in the future, not only  
404 with respect to SNVs but also SVs, although less is known regarding the entire repertoire  
405 of SVs present in a population than that of SNVs. Both integrative haploid reference  
406 genomes such as JG1 and collective genome reference developed in the future such as  
407 genome graphs—both of which can be constructed from hundreds of *de novo*  
408 assemblies—should advance the accuracy of genome analyses and facilitate development  
409 of personalized medicine approaches.

410

## 411 **METHODS**

### 412 **Selection and analysis of donor individuals**

413 ***Donor selection:*** Three male Japanese volunteers were recruited and participated in this  
414 study with written, informed consent.

415 ***G-banding analysis (Supplementary Fig. 1a–c):*** G-banding analyses for the three  
416 volunteers were performed using phytohemagglutinin-stimulated lymphocytes at the  
417 laboratory of SRL Inc. (Tokyo, Japan).

418 ***PCA of donors with Japanese samples (Fig. 1a):*** Paired-end reads with length of 162 bp  
419 from the three donors (jg1a, jg1b, and jg1c) were individually mapped to hs37d5.fa, and  
420 variant calling was performed according to previously described methods<sup>37</sup>, following  
421 GATK Best practices. The resulting VCF file was subjected to PCA using EIGENSOFT  
422 software (ver. 4.2). We chose 310 Japanese samples from the 3.5KJPNv2 cohort<sup>37</sup>; 100  
423 from Miyagi Prefecture in northern Japan; 29 from Nagahama City, in western Japan; and  
424 181 from Nagasaki Prefecture, in southern Japan. Variants shared among the 313 samples  
425 were selected and filtered using plink software (ver. 1.9) with the '--geno 0.05 --maf 0.05  
426 --hwe 0.05', and '--indep-pairwise 1500 150 0.03' options. The resulting dataset consisted  
427 of 18,658 variants.



## 428 **Genome analyses**

429 **Long-read SMRT sequencing:** Long-read SMRT sequencing was performed as  
430 previously described<sup>19</sup>. Briefly, genomic DNA from nucleated blood cells was sheared to  
431 ~20 kb and used for library preparation with a DNA template prep kit 2.0 (Pacific  
432 Biosciences; Menlo Park, CA). Size selection was carried out using the Blue Pippin  
433 system (Sage Science; Beverly, MA), targeting 18 kb (10-15 kb for some libraries of  
434 jg1a). The libraries were sequenced on a PacBio RSII instrument using P6-C4 chemistry.

435 **Optical mapping:** Optical mapping was performed using Irys system or Saphyr system,  
436 according to the manufacturer's protocol (Bionano Genomics; San Diego, CA). For  
437 sample jg1a, high-molecular-weight genomic DNA from nucleated blood cells was  
438 nicked using the endonucleases Nt.BspQI or Nb.BssSI and then labeled with fluorophore-  
439 tagged nucleotides. The labeled DNA was imaged on the Irys system. For samples jg1b  
440 and jg1c, high-molecular-weight genomic DNA from nucleated blood cells was labeled  
441 using direct labeling and staining (DLS) chemistry. The labeled DNA was imaged on the  
442 Saphyr system.

443 **Short-read paired-end sequencing:** Short-read paired-end sequencing was performed as  
444 previously described<sup>42</sup>. Briefly, genomic DNA from buffy coat samples was fragmented  
445 to an average target size of 550 bp, and then subjected to library construction using a  
446 TruSeq DNA PCR-Free HT sample prep kit (Illumina; San Diego, CA). The libraries  
447 were sequenced on a HiSeq 2500 system (Illumina) with a TruSeq Rapid PE Cluster kit  
448 (Illumina) and TruSeq Rapid SBS kit (Illumina) to obtain 162- or 259-bp paired-end reads.

449 **Mate-pair sequencing:** Genomic DNA from nucleated blood cells was used for library  
450 construction with a Nextera Mate Pair Library Preparation kit, gel-free protocol  
451 (Illumina). The libraries were size-selected to an average of 500 bp using AMPure XP  
452 beads (Beckman Coulter; Indianapolis, IN) and sequenced on a HiSeq 2500 system

453 (Illumina) with a TruSeq Rapid PE Cluster kit (Illumina), and TruSeq Rapid SBS kit  
454 (Illumina) to obtain 201-bp paired-end reads.

#### 455 **Overview of the computational methods for JG1 construction**

456 A diagram showing an overview of the construction of JG1 is provided in Supplementary  
457 Fig. 3. JG1 was constructed according to the following steps, which are also described in  
458 the download page for the JG1 sequence file from the jMorp website  
459 ([https://jmorp.megabank.tohoku.ac.jp/dj1/datasets/tommo-jg1.0.0.beta-](https://jmorp.megabank.tohoku.ac.jp/dj1/datasets/tommo-jg1.0.0.beta-20190424/files/tech-notes-for-computation.pdf)  
460 [20190424/files/tech-notes-for-computation.pdf](https://jmorp.megabank.tohoku.ac.jp/dj1/datasets/tommo-jg1.0.0.beta-20190424/files/tech-notes-for-computation.pdf)). The computation was carried out by  
461 using the Tohoku Medical Megabank Organization (ToMMO) Super Computer  
462 (<https://sc.megabank.tohoku.ac.jp/en/outline>).

463 **De novo assembly of PacBio subreads:** PacBio subreads were assembled using Falcon  
464 software<sup>23</sup> (build ver. falcon-2017.11.02-16.04-py2.7-ucs2.tar.gz) with the following  
465 configurations: reads shorter than 9 kb were used for error-correction of the longer reads  
466 ('length\_cutoff = 9000'), and error-corrected reads longer than 15 kb were used for  
467 assembly ('length\_cutoff\_pr = 15000'). Detailed settings are provided below:

```
468 length_cutoff = 9000  length_cutoff_pr = 15000  genome_size = 3200000000  
469 pa_HPCdaligner_option = -v -dal128 -t16 -e.75 -M16 -l4800 -k18 -h480 -w8 -  
470 s100 -T1  
471 ovlp_HPCdaligner_option = -v -dal128 -M24 -k24 -h1024 -e.96 -l2500 -s100 -T1  
472 pa_DBsplit_option = -x500 -s400  ovlp_DBsplit_option = -s400  
473 falcon_sense_option = --ouput_multi --min_idt 0.70 --min_cov = 4 --max_n_read  
474 200 --n_core 1 overlap_filtering_setting = --max_diff 60 --max_cov 60 --  
475 min_cov 0 --n_core 12
```

476 The contigs were then polished with the PacBio subreads using ArrowGrid software<sup>24</sup>  
477 (ver. 81b03f1; GitHub commit tag), with slight modifications to accommodate our  
478 number of data files and UGE settings.

479 **De novo assembly of Bionano optical maps:** We obtained two sets of Bionano data using  
480 two different enzymes, Nt.BspQI and Nb.BssSI, for subject jg1a, and one set of Bionano  
481 data was obtained with DLE-1 for jg1b and jg1c. In both cases, the Bionano data were  
482 assembled in two steps—a rough assembly step and a full assembly step—to perform *de*  
483 *novo* assembly as independently as possible from the reference. For the rough assembly  
484 step for jg1a, we ran pipelineCL.py software using the following settings:

```
485 -T 128 -j 8 -f 0.2 -i 0 -b ${data}/Molecules.bnx -l ${work} -V 0 -A -z -u -m  
486 -t ${bin}/Solve3.1_08232017/RefAligner/6700.6920rel/avx/ -a  
487 ${bin}/Solve3.1_08232017/RefAligner/6700.6920rel/optArguments_nonhaplotype_ir  
488 ys.xml -C ${work}/clusterArguments_${ver}.xml  
489
```

490 For the full assembly step for subject jg1a, we ran the software using the following  
491 settings:

```
492 -T 128 -j 8 -f 0.2 -i 5 -b ${data}/Molecules.bnx -l ${work} -V 0 -y -m  
493 -t ${bin}/Solve3.1_08232017/RefAligner/6700.6920rel/avx/ -a  
494 ${bin}/Solve3.1_08232017/RefAligner/6700.6920rel/optArguments_nonhaplotype_ir  
495 ys.xml -C ${work}/clusterArguments_${ver}.xml -r  
496 ${rough_assembly_output}/exp_mrg0/EXP_MRG0A.cmap  
497
```

498 For the rough assembly step for subjects jg1b and jg1c, we ran the software using the  
499 following settings:

```
500 -f 0 -i 5 -b ${data}/all.bnx -l ${work} -V 0 -N 4 -R  
501 -t ${bin}/Solve3.2.1_04122018/RefAligner/7437.7523rel/avx/ -a  
502 ${bin}/Solve3.2.1_04122018/RefAligner/7437.7523rel/avx/optArguments_nonhaplot  
503 ype_DLE1_saphyr_human.xml  
504 -C ${work}/clusterArgumentsBG_saphyr_phi_${ver}.xml
```

505

506 For the full assembly step of subjects jg1b and jg1c, we ran the software using the  
507 following settings:

```
508 -f 0 -i 5 -b ${data}/all.bnx -l ${work} -V 0 -N 4 -R -y  
509 -t ${bin}/Solve3.2.1_04122018/RefAligner/7437.7523rel/avx/ -a  
510 ${bin}/Solve3.2.1_04122018/RefAligner/7437.7523rel/avx/optArguments_nonhaplot  
511 ype_DLE1_saphyr_human.xml  
512 -C ${work}/clusterArgumentsBG_saphyr_phi_${ver}.xml  
513 -r ${rough_assembly_output}/exp_mrg0/EXP_MRG0A.cmap  
514
```

515 The '-T' and '-j' options were varied for computational efficiency. The BionanoSolve  
516 software suite was used for the above computation. We used BionanoSolve (ver. 3.1) for  
517 the assembly for subject jg1a, and ver.3.2 for the assembly for subjects jg1b and jg1c.

518 **Hybrid scaffolding:** Hybrid scaffolding was performed using BionanoSolve software  
519 (ver. 3.2). Hybrid scaffolding for subject jg1a was performed in the two-enzyme hybrid  
520 scaffolding mode using the runTGH.R script with the following options:

```
521 -N ${jg1a}-p_ctg.arrow.fa -e1 BSPQI -e2 BSSSI  
522 -b1 ${BspQI_work}/contigs/exp_refineFinal1/EXP_REFINEFINAL1.cmap  
523 -b2 ${BssSI_work}/contigs/exp_refineFinal1/EXP_REFINEFINAL1.cmap  
524 -O ${jg1a_hybscf}/${prefix}  
525 -R ${bin}/Solve3.2.1_04122018/RefAligner/7437.7523rel/avx/RefAligner  
526 ${bin}/Solve3.2.1_04122018/HybridScaffold/04122018/TGH/hybridScaffold_two_enz  
527 ymes.xml  
528
```

529 Hybrid scaffolding for subjects jg1b and jg1c was performed in the single-enzyme hybrid  
530 scaffolding mode, using the hybridScaffold.pl script with the following options:

```
531 -n ${arrow_work}/${individual}-p_ctg.arrow.fa
532 -b ${bionano_work}/contigs/exp_refineFinal1/EXP_REFINEFINAL1.cmap
533 -c ${hybscf_work}/hybridScaffold_DLE1_config.tmem.xml
534 -r ${bin}/Solve3.2.1_04122018/RefAligner/7437.7523rel/avx/RefAligner
535 -o ${work} -B 2 -N 2 -f
536 -e ${bionano_work}/contigs/auto_noise/autoNoise1.errbin
537
```

538 **Error correction with short reads:** Two sets of Illumina paired-end short reads with  
539 lengths of 162 bp and 259 bp were mapped to the hybrid scaffolds using BWA-MEM  
540 software<sup>4</sup> (ver. 0.7.17) with the option '-t 22 -K 1000000'. The alignment file was  
541 coordinate-sorted and compressed using the Picard tools (ver. 2.18.4) SortSam command.  
542 The resulting BAM files for the 162- and 259-bp paired-end reads were merged using the  
543 Picard tools MergeSamFiles command. The merged BAM files were then split to each  
544 scaffold using SAMtools<sup>43</sup> (ver. 1.8) view command, and then each scaffold was polished  
545 using Pilon software<sup>44</sup> (ver. 1.22, modified to correct the issue reported at  
546 <https://github.com/broadinstitute/pilon/issues/48>) with the option '--threads 22 --diploid -  
547 -changes --vcf --tracks'. The FASTA files for each polished scaffold were then merged  
548 into a single multi-FASTA format file.

549 **Meta-assembly:** The three sets of polished scaffolds were then meta-assembled using  
550 Metassembler software<sup>25</sup> (ver. 1.5; with modification of the type of 'totalBases' variable  
551 in the CEstat.hh from int to long to accommodate large genomes). There were 12 possible  
552 combinations to meta-assemble the three sets: (a + (b + c)), (a + (c + b)), ((a + b) + c), ((a  
553 + c) + b), (b + (a + c)), (b + (c + a)), ((b + a) + c), ((b + c) + a), (c + (a + b)), (c + (b + a)),  
554 ((c + a) + b), and ((c + b) + a), where x + y indicates meta-assemble x and y in this order.  
555 For each round of meta-assembly, we aligned the two assemblies using the NUCmer  
556 command of MUMmer software<sup>45</sup> (ver. 4.0.0beta2) with the option '--maxmatch -c 50 -l  
557 300'. The resulting DELTA file was filtered using delta-filter software with the option '-

558 1' to extract one-to-one correspondence. Next, the DELTA file was converted to  
559 COORDS format using the show-coords command with '-clrTH' option. Short mate-pair  
560 reads were classified into four categories (mp, pe, se, and unknown) using NxTrim  
561 software<sup>46</sup> (ver. 0.4.3), and the resulting set of reads with the correct mate-pair orientation  
562 (mp) were mapped using Bowtie2 software<sup>47</sup> (ver. 2.3.4.1) with the '--minins 1000 --  
563 maxins 16000 --rf' options. The output SAM file was then processed using the mateAn  
564 command with '-A 2000 -B 15000' option, indicating that the range of insert length was  
565 2 to 15 kb. The NUCmer alignment information and the mate-pair mapping information  
566 were integrated using the asseMerge command with '-i 5 -c 6' option. Finally, the resulting  
567 METASSEM file was converted to FASTA format using the meta2fasta command.

568 **Major allele substitution:** The three sets of polished hybrid scaffolds were aligned to the  
569 12 sets of meta-scaffolds using minimap2 (ver. 2.12), and variants were called using the  
570 paftools call command. After normalizing the manner of variant representation using the  
571 BCFtools norm command (ver. 1.8), SNVs and SVs shared by two of the three genomes  
572 were detected using the BCFtools isec command, and these were regarded as the major  
573 alleles and employed in JG1 using the BCFtools consensus command. For multi-allelic  
574 sites, one allele was chosen randomly.

575 **Detection of STS marker amplification by electronic PCR:** We detected amplification  
576 of the STS markers in the three genetic and six RH maps (Genethon, Marshfield, and  
577 deCODE genetic maps; GeneMap-G3, GeneMap99-GB4, TNG, NCBI\_RH, Stanford-G3,  
578 and Whitehead-RH maps) from the meta-scaffolds using the in-house electronic PCR  
579 software gPCR (ver. 2.6a) with the '-S -D' option ('-S' to show amplicon sequence, '-D':  
580 to show direction of markers). The STS markers were obtained from the UniSTS database  
581 ([ftp://ftp.ncbi.nih.gov/pub/ProbeDB/legacy\\_unists/](ftp://ftp.ncbi.nih.gov/pub/ProbeDB/legacy_unists/)). The results were used to infer the  
582 presence of chimeric scaffolds. One set of meta-scaffolds (jg1c + (jg1a + jg1b)) was

583 selected for the primary downstream analysis. In addition, to build the X and Y  
584 chromosomes, an additional set of meta-scaffolds (jg1a + (jg1b + jg1c)) was selected.

585 **Anchoring scaffolds to chromosomes:** The electronic PCR results were converted to  
586 BED format files, and the coordinates of some RH maps were scaled to approximately  
587 2,000 to fit those for the genetic maps; this was done to make it easier to understand the  
588 visualization results of the ALLMAPS software<sup>34</sup> (ver. 0.8.12) but did not affect the  
589 anchoring results. These maps were merged using the ALLMAPS mergebed command,  
590 and then processed using the ALLMAPS path command with the option '--  
591 gapsize=10000' to anchor the meta-scaffolds to the chromosomes. The weights of each  
592 of the three genetic maps was set to 5, and that of each of the RH maps was set to 1 in the  
593 weights.txt file. To anchor the sex chromosomes, three maps (deCODE, TNG, and  
594 Stanford-G3) that could anchor some scaffolds to the Y chromosome were used.

595 **Manual modification:** The physical lengths of the short arms of acrocentric  
596 chromosomes 13, 14, 15, 21, and 22 were obtained from Table 4 of Morton (1991)<sup>48</sup>. The  
597 relative length estimates of constitutive heterochromatin regions in the chromosome 1, 9,  
598 16 were obtained from ref. 49 and ref. 50. The relative length estimate of heterochromatin  
599 segment of the Y chromosome was obtained from ref. 51–53. These relative lengths were  
600 converted to the base-pair length (Mb) by using the chromosomal arms shown in Table 4  
601 of Morton (1991)<sup>48</sup>. The length of consecutive Ns for each chromosome is provided in  
602 Supplementary Table 5. For all chromosomes except 8 and 11, 3-Mb consecutive Ns were  
603 inserted instead of 10-kb Ns inserted by the ALLMAPS software, between the two  
604 scaffolds flanking the centromere. For chromosomes 8 and 11, in which the centromere-  
605 specific sequence repeats were identified in the midst of a scaffold by aligning the  
606 LinearCen1.1 sequences<sup>54</sup> using minimap2 software, no centromeric Ns were inserted.  
607 The position of the centromere was inferred from the Whitehead-RH and GeneMap99-  
608 GB4 maps, in which the centromeric or constitutive heterochromatin region could be

609 inferred from the region sparsely covered by STS markers, possibly due to the radiation  
610 conditions.

611 ***Building the X and Y chromosomes:*** We noted that one set of meta-scaffolds, (jg1c +  
612 (jg1a + jg1b)), contained a chimeric scaffold between the long arm of the X chromosome  
613 and the *SRY* locus of the Y chromosome. To reduce the chimeric meta-scaffolds, we chose  
614 apparently non-chimeric scaffolds anchored to the long arm of the X chromosome from  
615 another set of meta-scaffolds, (jg1a + (jg1b + jg1c)), and linked them to the scaffold of  
616 the short arm of the X chromosome.

617 ***Masking the pseudo-autosomal region:*** To locate the pseudo-autosomal regions, we  
618 aligned both the X and Y chromosomes from JG1 using minimap2 with the option '-cx  
619 asm5', and vice versa. The alignment started from the terminal region of the Y  
620 chromosome and ended at 2.26 Mb. This region was regarded as the putative PAR1 region.  
621 Other regions such as PAR2 and XTR were probably unresolved for unknown reasons.  
622 The putative PAR1 region was masked using the BEDTools software<sup>55</sup> (ver. 2.27.1)  
623 maskfasta command.

624 ***Mitochondrial chromosome:*** We aligned the set of meta-scaffolds to GRCh38, the  
625 mitochondrial sequence of which was obtained from the rCRS using minimap2 with the  
626 option '-cx asm5' to identify a scaffold that corresponds to the mitochondrial genome. We  
627 found a scaffold of 16,568 bp in length corresponding to the mitochondrial sequence in  
628 another set of meta-scaffolds (jg1a + (jg1b + jg1c)). The start site of the scaffold and the  
629 rCRS sequence differed; therefore, we shifted the start site of the scaffold to match that  
630 of the rCRS sequence.

631

632 **Idiogram drawing (Fig. 1b)**



633 Idiograms were depicted using JG1 BED files scaled to 90% of the original length so that  
634 the drawing of JG1 chromosomes longer than that of GRCh38 would be successful using  
635 the NCBI Genome Decoration Page (<https://www.ncbi.nlm.nih.gov/genome/tools/gdp>).  
636 The length of the chromosomal arms and the centromeric regions of the idiograms were  
637 manually modified to fit the scaffold length of JG1.

638

### 639 **Possible shared large inversion (Supplementary Fig. 2)**

640 Two large scaffolds corresponding to chromosome 9 were extracted from each assembly  
641 using the faSomeRecords command. Orientation was carried out by using the seqtk  
642 software (ver. 1.3) 'seq -r' command. Next, the chromosome 9 sequence from GRCh38  
643 and the two large scaffolds from each subject were aligned using minimap2 (ver. 2.12)  
644 with the '-t 12 -x asm5 --cs' option. Harr plots were drawn using the minidot command  
645 (bundled with miniasm software<sup>56</sup> [ver. 0.2]) with the '-L' option. The idiogram of  
646 chromosome 9 was drawn using the NCBI Genome Decoration Page.

647

### 648 **PCA of the JG1 and hg19 haplotypes with HapMap3 haplotypes (Fig. 2a, b)**

649 Two haplotypes were constructed for each individual of the HapMap3 variant information.  
650 JG1 haplotypes were constructed by aligning JG1 to hs37d5 using minimap2 software  
651 and identifying the allele at the marker sites. Variants shared among the 2,022 HapMap3  
652 haplotypes and JG1 and hg19 haplotypes were filtered the using plink software with '--  
653 geno 0.05 --maf 0.05' option. The resulting dataset consisted of 178,047 variants. PCA  
654 was performed using EIGENSOFT software. For PCA of JG1 and three Asian  
655 populations, 505 haplotypes from the JPT, CHB, and CHD populations were chosen. Four

656 CHD samples were omitted due to apparent inconsistency inferred from a pre-analysis of  
657 the PCA plots including these samples.

658

### 659 **SV analysis (related to Fig. 3)**

660 **SV detection:** The GRCh38 sequence was downloaded from illumina iGenome website  
661 ([ftp://ussd-](ftp://ussd-ftp.illumina.com/Homo_sapiens/NCBI/GRCh38/Homo_sapiens_NCBI_GRCh38.tar.gz)

662 [ftp.illumina.com/Homo\\_sapiens/NCBI/GRCh38/Homo\\_sapiens\\_NCBI\\_GRCh38.tar.gz](ftp.illumina.com/Homo_sapiens/NCBI/GRCh38/Homo_sapiens_NCBI_GRCh38.tar.gz)).

663 The JG1 sequence was aligned to GRCh38 using minimap2 (ver. 2.12) with the '-t 24 -cx  
664 asm5 --cs=long' option. The resulting PAF file was subjected to variant calling using the  
665 paftools call command. The VCF file was normalized using the BCFtools (ver. 1.8) norm  
666 command with the '--threads 4 --remove-duplicates' option. SVs  $\geq 51$  bp and  $< 10$  kb were  
667 subjected to the downstream analyses.

668 **Average depth analysis:** Two hundred Japanese individuals (100 males and 100 females)  
669 were selected from the 3,552 samples<sup>37</sup>. The 162-bp paired-end reads were mapped using  
670 BWA-MEM software as described<sup>37</sup>. Next, accessible regions were defined as the regions  
671 where the average depth among the 200 individuals is  $\geq 5$  and  $\leq \text{mean} + 2\text{SD}$ ; mean and  
672 SD were computed for each chromosome. SVs detected within the accessible regions and  
673 detected by comparing same autosomes of GRCh38 and JG1 were considered. The mean  
674 value of the average depth within the adjacent upstream region of SV was regarded as the  
675 reference value, and the  $\Delta$  average depth was defined as the difference between the  
676 reference value and the value of the position showing the maximum absolute difference  
677 within the SV region.

### 678 **Rare disease exome analysis (related to Fig. 4)**

679 Exome analyses were carried out following the GATK best practices for germline variant  
680 detection. Short reads were mapped using BWA-MEM software, and the resulting  
681 alignment files were sorted and duplication-marked using SAMtools<sup>43</sup> software. Variants  
682 of the disease cohort families were called using GATK software (ver. 4.0 to 4.1), and the  
683 joint calling process was carried out with samples from other Japanese subjects with  
684 various rare diseases. The BED files describing the exome capture regions (SureSelect  
685 Human All Exon V5, Agilent) were lifted over using the paftools liftover command. The  
686 GATK resource bundles were lifted over to the JG1 coordinates using the Picard tools  
687 LiftoverVcf command. GENCODE (ver. 29) annotations were lifted over to the JG1  
688 coordinates using an in-house script. The chain files, which were required for lifting over,  
689 were generated from the results of minimap2 with an in-house script. The SnpEff  
690 database<sup>58</sup> was constructed using the lifted-over GENCODE annotation file and used for  
691 variant effect predictions. Variants called against JG1 were lifted over to the hs37d5  
692 coordinates by using the Picard tools LiftoverVcf command. Overlap relationships  
693 between the variants were assessed using the BCFtools (ver. 1.9) isec command.

#### 694 **Statistical tests and graph drawing**

695 Statistical tests were performed using R software (ver. 3.5.1). Histograms were drawn  
696 using R software (ver. 3.5.1) and ggplot2 library (ver. 3.0.0).

#### 697 **Data availability**

698 JG1 sequence, chain files and GENCODE annotation files are available from the jMorp  
699 website (<https://jmorp.megabank.tohoku.ac.jp/201911/downloads#sequence>).

700

701

702 **ACKNOWLEDGEMENTS**

703 This work was supported in part by the Tohoku Medical Megabank (TMM) Project from  
704 the Ministry of Education, Culture, Sports, Science and Technology (MEXT) and the  
705 Reconstruction Agency; the Japan Agency for Medical Research and Development  
706 (AMED; Grant Numbers JP19km0105001 and JP19km0105002) for Tohoku University.  
707 All computational resources were provided by the ToMMo supercomputer system  
708 (<http://sc.megabank.tohoku.ac.jp/en>), which is supported by Facilitation of R&D  
709 Platform for AMED Genome Medicine Support conducted by AMED (Grant Number  
710 JP19km0405001). We appreciate all the volunteers who participated in the TMM project.

711

712

713 **AUTHOR CONTRIBUTIONS**

714 J.T., S.T, K.Y., C.G., T.F., S.M., and Y.O. performed computational analyses. J.T., A.K.,  
715 S.K., and G.T. interpreted the results of rare disease re-analyses. F.K., J.K., A.O., and J.Y.  
716 designed and conducted experiments. J.T. and G.T. wrote the manuscript with the  
717 assistance of the others. K.K., M.Y., and G.T. conceived and supervised the project. All  
718 authors read and approved the final manuscript.

719

720

721 **COMPETING FINANCIAL INTERESTS**

722 The authors declare no competing financial interests.

723

724 **REFERENCES**

- 725 1. Lander E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature*  
726 **409**, 860-921 (2001).
- 727 2. Venter, J. C. *et al.* The Sequence of the Human Genome. *Science* **291**, 1304-1351  
728 (2001).
- 729 3. Metzker, M. L. Sequencing technologies — the next generation. *Nat. Rev. Genet.*  
730 **11**, 31-46 (2010).
- 731 4. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-  
732 MEM. *arXiv [q-bio.GN]* 1303.3997 (2013).
- 733 5. Church, D.M. *et al.* Modernizing reference genome assemblies. *PLoS Biol.* **9**,  
734 e1001091 (2011).
- 735 6. Schneider, V. A. *et al.* Evaluation of GRCh38 and *de novo* haploid genome  
736 assemblies demonstrates the enduring quality of the reference assembly. *Genome*  
737 *Res.* **27**, 849-864 (2017).
- 738 7. Magi, A. *et al.* Characterization and identification of hidden rare variants in the  
739 human genome. *BMC Genomics* **16**, 340 (2015).
- 740 8. Koko, M., Abdallah, M. O. E., Amin, M. & Ibrahim, M. Challenges imposed by  
741 minor reference alleles on the identification and reporting of clinical variants from  
742 exome data. *BMC Genomics* **19**, 46 (2018).
- 743 9. Degner, J. F. *et al.* Effect of read-mapping biases on detecting allele-specific  
744 expression from RNA-sequencing data. *Bioinformatics* **25**, 3207-3212 (2009).
- 745 10. Green, R. E. *et al.* A draft sequence of the Neandertal genome. **328**, 710-722  
746 (2010).
- 747 11. Sherman, R. M. *et al.* Assembly of a pan-genome from deep sequencing of 910  
748 humans of African descent. *Nat. Genet.* **51**, 30-35 (2019).
- 749 12. Audano, P. A. *et al.* Characterizing the major structural variant alleles of the  
750 human genome. *Cell* **176**, 663-675 (2019).

- 751 13. Dewey, F. E. *et al.* Phased whole-genome genetic risk in a family quartet using a  
752 major allele reference sequence. *PLoS Genet.* **7**, e1002280 (2011).
- 753 14. Paten, B. *et al.* Cactus: Algorithms for genome multiple sequence alignment.  
754 *Genome Res.* **21**, 1512-28 (2011).
- 755 15. Garrison, E. *et al.* Variation graph toolkit improves read mapping by representing  
756 genetic variation in the reference. *Nat. Biotechnol.* **36**, 875-879 (2018).
- 757 16. Rakocevic, G. *et al.* Fast and accurate genomic analyses using genome graphs.  
758 *Nat. Genet.* **51**, 354-362 (2019).
- 759 17. Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome  
760 alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.*  
761 **37**, 907-915 (2019).
- 762 18. Ameur, A. *et al.* *De novo* assembly of two Swedish genomes reveals missing  
763 segments from the human GRCh38 reference and improves variant calling of  
764 population-scale sequencing data. *Genes* **9**, 486 (2018).
- 765 19. Nagasaki, M. *et al.* Construction of JRG (Japanese reference genome) with single-  
766 molecule real-time sequencing. *Hum. Genome Var.* **6**, 27 (2019).
- 767 20. Cho, Y.S. *et al.* An ethnically relevant consensus Korean reference genome is a  
768 step towards personal reference genomes. *Nat. Commun.* **7**, 13637 (2016).
- 769 21. Seo, J.-S. *et al.* *De novo* assembly and phasing of a Korean human genome. *Nature*  
770 **538**, 243-247 (2016).
- 771 22. Shi, L. *et al.* Long-read sequencing and *de novo* assembly of a Chinese genome.  
772 *Nat. Commun.* **7**, 12065 (2016).
- 773 23. Chin, C.-S. *et al.* Phased diploid genome assembly with single-molecule real-time  
774 sequencing. *Nat. Methods* **13**, 1050-1054 (2016).
- 775 24. Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-  
776 mer weighting and repeat separation. *Genome Res.* **27**, 722-736 (2017).
- 777 25. Wences, A. H. & Schatz, M. C. Metassembler: merging and optimizing *de novo*

- 778 genome assemblies. *Genome Biol.* **16**, 207 (2015).
- 779 26. Dib, C. *et al.* A comprehensive genetic map of the human genome based on 5,264  
780 microsatellites. *Nature* **380**, 152-154 (1996).
- 781 27. Kong, A. *et al.* A high-resolution recombination map of the human genome. *Nat.*  
782 *Genet.* **31**, 241-247 (2002).
- 783 28. Broman, K. W. *et al.* Comprehensive human genetic maps: individual and sex-  
784 specific variation in recombination. *Am. J. Hum. Genet.* **63**, 861-869 (1998).
- 785 29. Hudson, T. J. *et al.* An STS-based map of the human genome. *Science* **270**, 1945-  
786 1954 (1995).
- 787 30. Stewart, E. A. *et al.* An STS-based radiation hybrid map of the human genome.  
788 *Genome Res.* **7**, 422-433 (1997).
- 789 31. Deloukas, P. *et al.* A physical map of 30,000 human genes. *Science* **282**, 744-746  
790 (1998).
- 791 32. Agarwala, R. *et al.* A fast and scalable radiation hybrid map construction and  
792 integration strategy. *Genome Res.* **10**, 350-364 (2000).
- 793 33. Olivier, M. *et al.* A high-resolution radiation hybrid map of the human genome  
794 draft sequence. *Science* **291**, 1298-1302 (2001).
- 795 34. Tang, H. *et al.* ALLMAPS: robust scaffold ordering based on multiple maps.  
796 *Genome Biol.* **16**, 3 (2015).
- 797 35. Andrews, R. M. *et al.* Reanalysis and revision of the Cambridge reference  
798 sequence for human mitochondrial DNA. *Nat. Genet.* **23**, 147-147 (1999).
- 799 36. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**,  
800 3094-3100 (2018).
- 801 37. Tadaka, S. *et al.* 3.5KJPNv2: an allele frequency panel of 3552 Japanese  
802 individuals including the X chromosome. *Hum. Genome Var.* **6**, 28 (2019).
- 803 38. Takezawa, Y. *et al.* Genomic analysis identifies masqueraders of full-term  
804 cerebral palsy. *Ann. Clin. Transl. Neurol.* **5**, 538-551 (2018).

- 805 39. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse  
806 genomes. *Nucleic Acids Res.* **8**, D766-D773 (2019).
- 807 40. Dudchenko, O. *et al.* De novo assembly of the *Aedes aegypti* genome using Hi-C  
808 yields chromosome-length scaffolds. *Science* **356**, 92-95 (2017).
- 809 41. Chaisson, M. J., *et al.* Resolving the complexity of the human genome using  
810 single-molecule sequencing. *Nature* **517**, 608-611 (2015).
- 811 42. Nagasaki, M. *et al.* Rare variant discovery by deep whole-genome sequencing of  
812 1,070 Japanese individuals. *Nat. Commun.* **6**, 8018 (2015).
- 813 43. Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics*  
814 **25**, 2078-2079 (2009).
- 815 44. Walker, B. J. *et al.* Pilon: An integrated tool for comprehensive microbial variant  
816 detection and genome assembly improvement. *PLoS One* **9**, e112963 (2014).
- 817 45. Marçais, G. *et al.* MUMmer4: A fast and versatile genome alignment system.  
818 *PLoS Comput. Biol.* **14**, 1-14 (2018).
- 819 46. O'Connell J. *et al.* NxTrim: Optimized trimming of Illumina mate pair reads.  
820 *Bioinformatics* **31**, 2035-2037 (2015).
- 821 47. Langmead, B. & Salzberg. S. L. Fast gapped-read alignment with Bowtie2. *Nat.*  
822 *Methods* **9**,357-359 (2012).
- 823 48. Morton, N. E. Parameters of the human genome. *Proc. Natl. Acad. Sci. USA* **88**,  
824 7474-7476 (1991).
- 825 49. Ludgren, R. *et al.* Constitutive heterochromatin C-band polymorphism in  
826 prostatic cancer. *Cancer Genet. Cytogenet.* **51**, 57-62 (1991).
- 827 50. Podugolnikova, O. A. & Blumina, M. G. Heterochromatic regions on  
828 chromosomes 1, 9, 16, and Y in children with some disturbances occurring during  
829 embryo development. *Hum. Genet.* **63**, 183-188 (1983).
- 830 51. Erçal, M. D. & Brøndum-Nielsen, K. Length polymorphism of heterochromatic  
831 segment of the Y chromosome in boys with acute leukemia. *Acta Paediatr. Jpn.*



- 832           **37**, 614-616 (1995).
- 833   52.    Petković, I. Variability of euchromatic and heterochromatic segment of the Y  
834           chromosome in men with malignant tumors and in a control group. *Cancer Genet.*  
835           *Cytogenet.* **13**, 29-36 (1984).
- 836   53.    Petković, I. *et al.* Heterochromatic segment length of Y chromosome in 55 boys  
837           with malignant diseases. *Cancer Genet. Cytogenet.* **25**, 351-353 (1987).
- 838   54.    Miga, K. H. *et al.* Centromere reference models for human chromosomes X and  
839           Y satellite arrays. *Genome Res.* **24**, 697-707 (2014).
- 840   55.    Quinlan, A. R. and Hall, I. M. BEDTools: A flexible suite of utilities for  
841           comparing genomic features. *Bioinformatics* **26**, 841-842 (2010).
- 842   56.    Li, H. Minimap and miniasm: fast mapping and de novo assembly for noisy long  
843           sequences. *Bioinformatics* **32**, 2103-2110 (2016).
- 844   57.    Altshuler, D. *et al.* Integrating common and rare genetic variation in diverse  
845           human populations. *Nature* **467**, 52–58 (2010).
- 846   58.    Cingolani, P. *et al.* A program for annotating and predicting the effects of single  
847           nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila*  
848           *melanogaster* strain *w<sup>1118</sup>*; *iso-2*; *iso-3*. *Fly (Austin)* **6**, 80-92 (2012).
- 849   59.    Neph, S. *et al.* BEDOPS: high-performance genomic feature operations.  
850           *Bioinformatics* **28**, 1919-1920 (2012).
- 851   60.    Mikheenko, A., Prjibelski, A., Saveliev, V., Antipov, D., & Gurevich, A. Versatile  
852           genome assembly evaluation with QUAST-LG. *Bioinformatics* **34**, i142-i150  
853           (2018).
- 854
- 855
- 856

857 **FIGURE LEGENDS**

858 **Figure 1.** Construction of JG1. **(a)** PCA plot showing that the three sample donors are  
859 within the Japanese population cluster. **(b)** Idiogram showing the regions sequenced for  
860 each chromosome in JG1. Red and blue boxes indicate scaffolds; the red box spans an  
861 entire chromosomal arm. Dark gray boxes denote E-gaps, which represent links  
862 connected by genetic and RH maps or gaps inserted according to other evidence. Pink  
863 boxes denote N-gaps, which are unresolved regions linked by Bionano optical maps, or  
864 the putative PAR1 region in the Y chromosome. **(c)** Harr plot representing the co-linearity  
865 between the reference genome GRCh38 and JG1.

866  
867

868 **Figure 2.** SNV characteristics of JG1. **(a)** PCA plot of the haplotype SNP composition of  
869 JG1, the reference hg19, and HapMap3 samples. **(b)** PCA plot of the haplotype SNP  
870 composition of JG1 and Asian samples from HapMap3. **(c)** Unfolded site frequency  
871 spectrum representing the frequencies of alleles employed in the JG1 sequence in the  
872 Japanese population of 3.5KJPNv2.

873

874

875 **Figure 3.** Analysis of JG1 SV. **(a)** Length histogram of small ( $\leq 500$  bp) and large ( $> 500$   
876 bp) insertions and deletions detected by comparing JG1 and GRCh38. **(b)** Distribution of  
877 insertions and deletions among the chromosomes of GRCh38. **(c)** JBrowse snapshots of  
878 one insertion/deletion example. Tracks are GENCODE gene annotations, detected SVs,  
879 and average depth of short reads from 200 Japanese samples. **(d)** Difference in average  
880 depth between the SV and upstream regions of same length as the SV for GRCh38 and  
881 JG1.

882

883

884 **Figure 4.** Comparison of variants called in exome analyses employing JG1 or hs37d5 as  
885 a reference genome. **(a)** Number of total variants, SNVs, and short indels called per  
886 individual. **(b)** Number of high- and moderate-impact variants. **(c)** Venn diagram showing  
887 overlap relationships between variants detected in JG1 (lifted over to the hs37d5  
888 coordinates) and those detected in hs37d5. Shown are results for a representative  
889 individual. **(d)** Unfolded site frequency spectra representing the frequency of non-GRC-  
890 type alleles in the variant sites detected specifically in JG1, in both genomes, and  
891 specifically in hs37d5, respectively. Shown are results for the same individual as in **(c)**.  
892  
893

894 **SUPPLEMENTARY FIGURE LEGENDS**

895 **Supplementary Fig. 1.** Karyotypes of the three subjects: **(a)** jg1a, **(b)** jg1b, and **(c)** jg1c.

896 The arrow in panel **(a)** indicates the normal variation inv(9)(p12q13).

897

898 **Supplementary Fig. 2.** Harr plot of the alignment between chromosome 9 of GRCh38

899 and two largest scaffolds aligned to chromosome 9 from the **(a)** jg1a, **(b)** jg1b, and **(c)**

900 jg1c assemblies, indicating that the two individual genomes harbor a possible shared

901 inversion. 'Super-scaffold' is the default prefix designated by BionanoSolve software.

902

903 **Supplementary Fig. 3.** Workflow of the construction of JG1. **(a)** Workflow of the

904 construction of each draft assembly. **(b)** Workflow of the integration of three draft

905 assemblies. Rectangles indicate substrates such as reads, contigs, and scaffolds.

906 Rectangles with rounded corners indicate software or processes.

907

908 **Supplementary Fig. 4:** Histogram of PacBio subreads length for **(a)** jg1a, **(b)** jg1b, and

909 **(c)** jg1c. The length of each subread was calculated using the SAMtools (ver. 1.8) faidx

910 command.

911

912 **Supplementary Fig. 5:** Histogram of Bionano data for **(a)** Nt.BspQI of jg1a, **(b)**

913 Nb.BssSI of jg1a, **(c)** jg1b, and **(d)** jg1c. The length of each molecule was extracted from

914 the BNX file.

915

916 **Supplementary Fig. 6:** Majority decision. **(a)** Schematic representation of majority

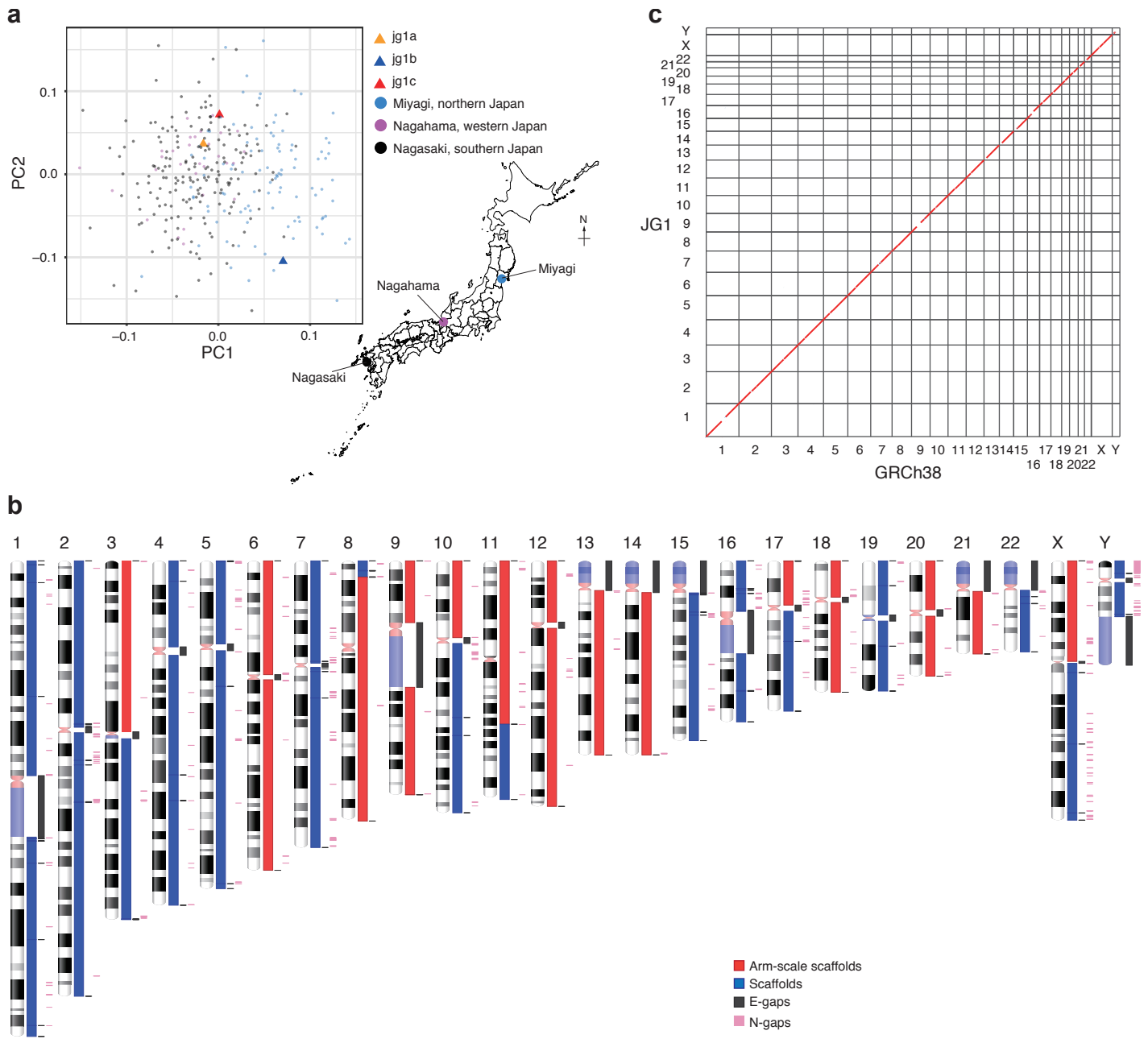
917 decision approach. **(b)** Venn diagram of SNVs detected in JG1, jg1a, jg1b, and jg1c by

918 comparison with hs37d5. The intersection relationship was inferred using the BCFtools

919 (ver. 1.8) isec command.

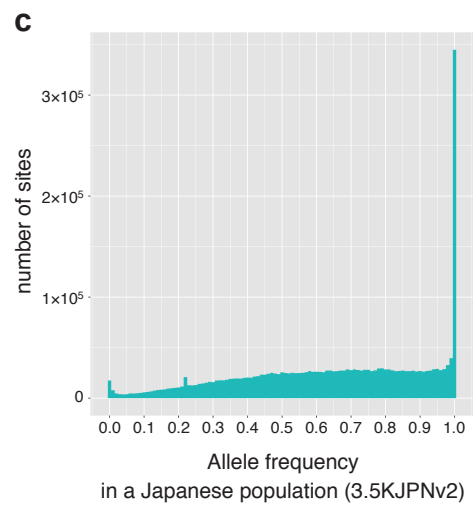
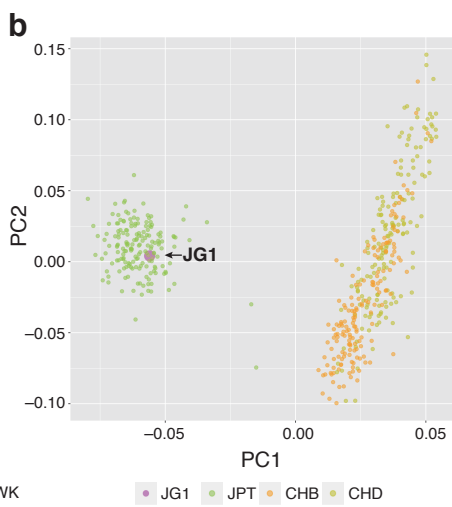
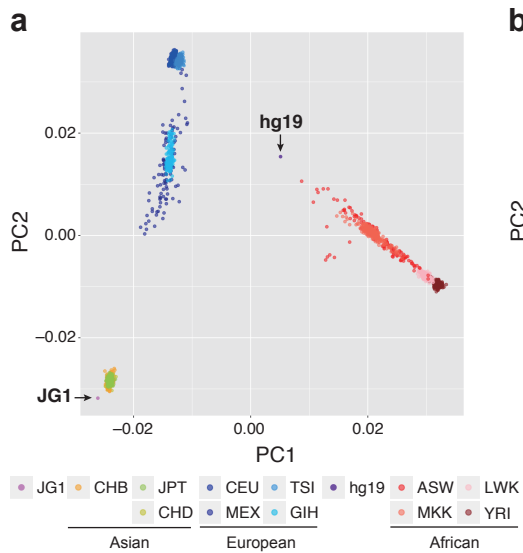
920

921 **Supplementary Fig. 7:** Length distributions of detected transposable elements in the  
922 GRCh38 and JG1 genomes. Shown are *Alu*, SVA, and LINE1. Transposable elements and  
923 their subclasses were identified using RepeatMasker software (ver. 4.0.7) with the '-  
924 species human' option. The resulting OUT format files were converted to BED format  
925 using the `rmsk2bed` command of BEDOPS software<sup>59</sup> (ver. 2.4.35). Transposable  
926 elements disrupted by other elements were counted as distinct.  
927

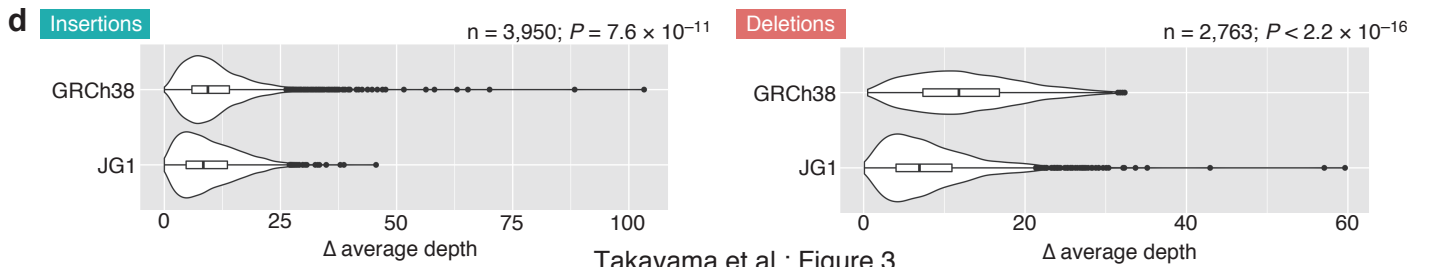
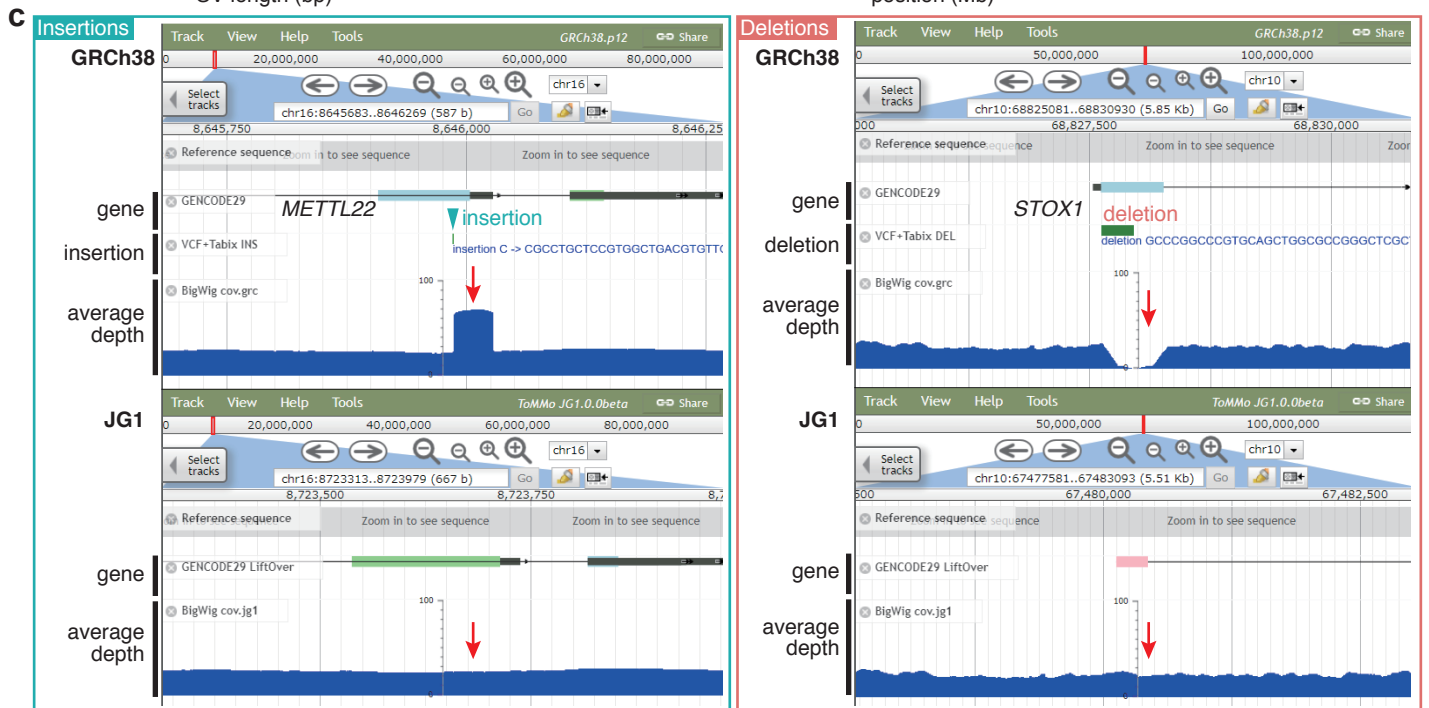
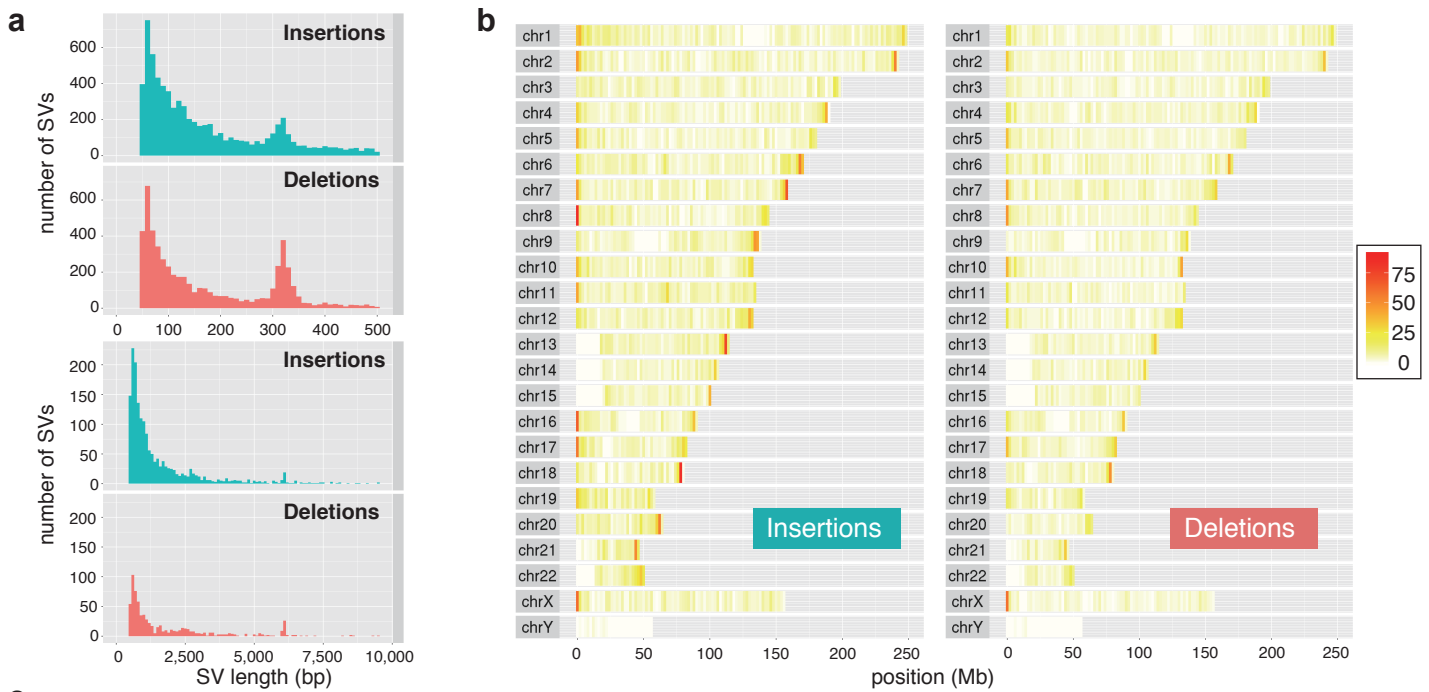


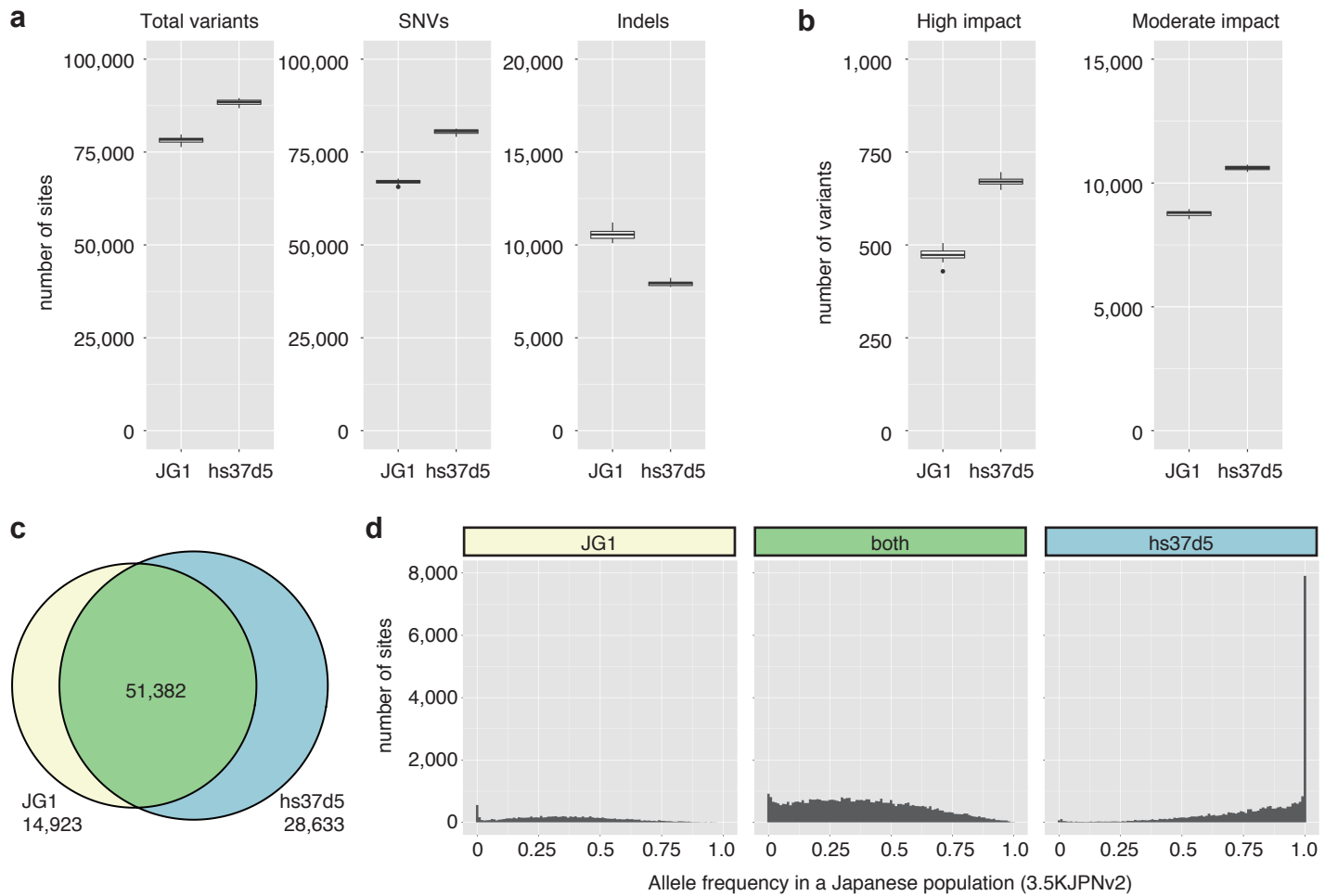
Takayama et al.; Figure 1





Takayama et al.; Figure 2





Takayama et al.; Figure 4

**Table 1. Assembly results.**

	Total length (bp)	Number of fragments	N50 (bp)	NG50 (bp)	Number of misassemblies	Reference
jg1a primary contigs	2,855,392,439	2,194	20,631,146	19,227,791	1,912	this study
jg1b primary contigs	2,852,624,381	2,227	21,603,629	19,007,577	1,673	this study
jg1c primary contigs	2,851,554,649	2,120	19,616,169	17,539,317	1,673	this study
jg1a scaffolds	2,889,327,167	1,911	86,280,884	59,417,266	2,071	this study
jg1b scaffolds	2,880,572,022	1,893	59,380,744	57,047,228	1,762	this study
jg1c scaffolds	2,875,657,275	1,797	58,198,703	58,048,754	1,867	this study
meta-scaffolds (jg1c + (jg1a + jg1b))	2,858,691,982	708	66,367,161	58,207,422	1,581	this study
JG1	3,085,782,898	624	141,953,703	141,953,703	1,654	this study
AK1 scaffolds	2,904,207,228	2,832	44,846,623	39,609,866	2,138	21
HX1 scaffolds	2,934,082,568	5,323	21,979,250	20,700,129	2,688	22
Swe1 scaffolds*	3,127,010,000	NA	49,799,000	NA	NA	18
Swe2 scaffolds*	3,103,497,000	NA	45,443,000	NA	NA	18

\* Results of Swe1 and Swe2 scaffolds were obtained from Aneur et al<sup>18</sup>. All other results were calculated by using Quast-LG software<sup>60</sup> with the reference GRCh38 as the truth set.

**Table 2. Diplegia cohort.**

<b>ID*</b>	<b>family</b>	<b>type</b>	<b>locus</b>	<b>variant(s)</b>	<b>variant effect prediction</b>	<b>mode of inheritance</b>
3	FD-05	trio	<i>CTNNB1</i>	c.1683+2T>C	HIGH (splice donor)	<i>de novo</i> SNV
5	FD-07	trio	<i>CYP2U1</i>	c.651delC	HIGH (frame shift)	autosomal recessive deletion
6	FD-08	trio	<i>SPAST</i>	c.1276C>T	MODERATE (missense)	<i>de novo</i> SNV
7	FD-09	quartet	<i>GNAO1</i>	c.736G>A	MODERATE (missense)	<i>de novo</i> SNV
9	FD-11	trio	<i>CACNA1A</i>	c.653C>T	MODERATE (missense)	<i>de novo</i> SNV
10	FD-12	trio	<i>SPAST</i>	c.1496G>A	MODERATE (missense)	<i>de novo</i> SNV
11	FD-13	trio	<i>AMPD2</i>	c.515+1G>A c.1724C>T	HIGH (splice donor) MODERATE (missense)	compound heterozygous SNV

\* Case ID from Table 2 in Takezawa et al<sup>38</sup>.