

## KOREF\_S1: the phased, parental Trio-binned Korean reference genome using long-reads and Hi-C sequencing methods

--Manuscript Draft--

|  |  |                   |
|--|--|-------------------|
| <b>Manuscript Number:</b>                            | GIGA-D-21-00219  |                   |
| <b>Full Title:</b>                                   | KOREF_S1: the phased, parental Trio-binned Korean reference genome using long-reads and Hi-C sequencing methods  |                   |
| <b>Article Type:</b>                                 | Data Note  |                   |
| <b>Funding Information:</b>                          | ministry of smes and startups (P0016193)   | Dr. Jong Hwa Bhak |
| <b>Abstract:</b>                                     | <p><b>Background</b></p> <p>KOREF is the Korean reference genome which was constructed with various sequencing technologies including long reads, short reads, and optical mapping methods. It is also the first East Asian multiomic reference genome accompanied by extensive clinical information, time series and multiomic data, and his parental sequencing data. However, it was still not a chromosome-scale reference. Here, we updated the previous KOREF assembly to a new chromosome-level haploid assembly of KOREF, KOREF_S1v2.1. ONT PromethION, PacBio Hifi-CCS, and Hi-C technology were used to build the most accurate East Asian reference assembled so far.</p> <p><b>Results</b></p> <p>We produced 705 Gb ONT reads and 114 Gb PacBio Hifi reads, and corrected ONT reads by PacBio reads. The corrected ultra-long reads reached higher accuracy of 1.4% base-errors than the previous KOREF_S1v1.0, which was mainly built with short reads. KOREF has parental genome information, and we successfully phased it using a trio-binning method acquiring a near-complete haploid-assembly. The final assembly resulted in total length of 2.9 Gb with an N50 of 150 Mb, and the longest scaffold covered 97.3% of GRCh38's chromosome 2. And the final assembly showed high base accuracy, less than 0.01% of base-errors.</p> <p><b>Conclusions</b></p> <p>KOREF_S1v2.1 is the first chromosome-scale haploid assembly of the Korean reference genome with high contiguity and accuracy. Our study provides useful resources of the Korean reference genome and demonstrates a new strategy of hybrid assembly which collaborates ONT's PromethION and PacBio's HiFi-CCS.</p> |                   |
| <b>Corresponding Author:</b>                         | Jong Hwa Bhak, Ph.D.<br>UNIST<br>Ulsan, Ulsan KOREA, REPUBLIC OF   |                   |
| <b>Corresponding Author Secondary Information:</b>   |  |                   |
| <b>Corresponding Author's Institution:</b>           | UNIST  |                   |
| <b>Corresponding Author's Secondary Institution:</b> |  |                   |
| <b>First Author:</b>                                 | Hui-su Kim, Ph.D.  |                   |
| <b>First Author Secondary Information:</b>           |  |                   |
| <b>Order of Authors:</b>                             | Hui-su Kim, Ph.D.  |                   |
|  | Sungwon Jeon   |                   |
|  | Yeonkyung Kim  |                   |
|  | Changjae Kim, Ph.D.  |                   |
|  |  |                   |

|   |                      |
|---|----------------------|
|   | Jihun Bhak           |
|   | Jong Hwa Bhak, Ph.D. |
| <b>Order of Authors Secondary Information:</b>  |                      |
| <b>Additional Information:</b>  |                      |
| <b>Question</b>   | <b>Response</b>      |
| Are you submitting this manuscript to a special series or article collection?   | No                   |
| <p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>  | Yes                  |
| <p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p> | Yes                  |
| <p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using</p>  | Yes                  |

a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

1 **KOREF\_S1: the phased, parental Trio-binned Korean**  
2 **reference genome using long-reads and Hi-C sequencing**  
3 **methods**

4

5 Hui-su Kim<sup>1</sup>, Sungwon Jeon<sup>1,2</sup>, Yeonkyung Kim<sup>1</sup>, Changjae Kim<sup>2</sup>, Jihun Bhak<sup>1,2</sup>, and Jong  
6 Bhak<sup>1,2,3,4\*</sup>

7 <sup>1</sup> Korean Genomics Center (KOGIC), Ulsan National Institute of Science and Technology  
8 (UNIST), Ulsan, 44919, Republic of Korea

9 <sup>2</sup> Department of Biomedical Engineering, College of Information and Biotechnology, Ulsan  
10 National Institute of Science and Technology (UNIST), Ulsan, 44919, Republic of Korea

11 <sup>3</sup> Clinomics LTD, Ulsan, 44919, Republic of Korea

12 <sup>4</sup> Personal Genomics Institute, Genome Research Foundation, Cheongju, 28160, Republic of  
13 Korea

14

15 \*Correspondence author:

16

17 Name: Jong Bhak, Ph.D.

18 Address: #110-303, Ulsan National Institute of Science and Technology, UNIST-gil 50, Eonyang-  
19 eup, Ulju-gu, Ulsan 44919, Republic of Korea

20 Phone: (+82) 10-4644-6754

21 Email: [jongbhak@genomics.org](mailto:jongbhak@genomics.org), ORCID: 0000-0002-4228-1299

## 22 **Abstract**

## 23 **Background**

24 KOREF is the Korean reference genome which was constructed with various sequencing  
25 technologies including long reads, short reads, and optical mapping methods. It is also the first  
26 East Asian multiomic reference genome accompanied by extensive clinical information, time  
27 series and multiomic data, and his parental sequencing data. However, it was still not a  
28 chromosome-scale reference. Here, we updated the previous KOREF assembly to a new  
29 chromosome-level haploid assembly of KOREF, KOREF\_S1v2.1. ONT PromethION, PacBio  
30 Hifi-CCS, and Hi-C technology were used to build the most accurate East Asian reference  
31 assembled so far.

## 32 **Results**

33 We produced 705 Gb ONT reads and 114 Gb PacBio Hifi reads, and corrected ONT reads by  
34 PacBio reads. The corrected ultra-long reads reached higher accuracy of 1.4% base-errors than the  
35 previous KOREF\_S1v1.0, which was mainly built with short reads. KOREF has parental genome  
36 information, and we successfully phased it using a trio-binning method acquiring a near-complete  
37 haploid-assembly. The final assembly resulted in total length of 2.9 Gb with an N50 of 150 Mb,  
38 and the longest scaffold covered 97.3% of GRCh38's chromosome 2. And the final assembly  
39 showed high base accuracy, less than 0.01% of base-errors.

## 40 **Conclusions**

41 KOREF\_S1v2.1 is the first chromosome-scale haploid assembly of the Korean reference genome  
42 with high contiguity and accuracy. Our study provides useful resources of the Korean reference

43 genome and demonstrates a new strategy of hybrid assembly which collaborates ONT's  
44 PromethION and PacBio's HiFi-CCS.

45 **Keywords:** Korean reference; KOREF\_S1; ONT PromethION; PacBio Hifi; Hi-C; hybrid  
46 assembly

47

48

## 49 **Introduction**

50 Since the human genome reference was released in 2003, it has been updated and recently was  
51 patched in 2019 (GRCh38.p13) by the Genome Reference Consortium (GRC) [1]. Despite high  
52 completeness of GRCh38 assembly, it derives from a single individual, mostly based on Caucasian  
53 and African ancestry [2]. It is the most precise and extensive among all human references  
54 constructed so far. Recently, due to recent cost-effective sequencing methods, especially long  
55 reads methods, one can construct human personal references fast and efficiently [3]. The first  
56 Korean reference, KOREF, has been constructed in two types [4]. The first is KOREF\_S1 which  
57 is a personal reference from an individual which is accompanied by parental *de novo* assemblies.  
58 The second one is KOREF\_C which is a consensus population reference that includes variome  
59 information of Koreans. KOREF was initiated by the Korean Ministry of Science and Technology  
60 in 2006 to generate a national genome and variome references and currently it is jointly developed  
61 by the Genome Research Foundation, National Standard Reference Research Center, and the  
62 Korean Genomics Center at UNIST (Ulsan National Institute of Science and Technology). The  
63 first version of KOREF\_S1, KOREF\_S1v1.0, had a clear limitation of short reads and long-  
64 distance mapping-based approaches that resulted in a relatively low-quality assembly compared to  
65 the current GRCh38. We used Oxford Nanopore Technologies (ONT) PromethION and PacBio  
66 HiFi sequencers to upgrade KOREF\_S1 by using a publicly available KOREF cell line.

67

## 68 **Materials and Methods**

### 69 **Sample preparation and genome sequencing**

70 Sample preparation steps were followed as the previous study [4]. Human KOREF cell lines  
71 (<http://koref.net>) were cultured at 37°C in 5% CO<sub>2</sub> in RPMI-1640 medium with 10% heat-  
72 inactivated fetal bovine serum. DNA was extracted from cells using the DNeasy Blood & Tissue  
73 kit (Qiagen) to the manufacturer's instructions. Sequencing libraries for the Oxford Nanopore  
74 Technologies PromethION were prepared using the 1D ligation sequencing kit (SQK-LSK109,  
75 Oxford Nanopore Technologies, UK) following the manufacturer's instructions. The products  
76 were quantified using the Bioanalyzer 2100 (Agilent, Santa Clara, CA, USA) and raw signals were  
77 generated by the PromethION R9.4.5 platform (Oxford Nanopore Technologies, UK). Base-  
78 calling the raw signals was performed using Guppy v4.0.11 with the Flip-flop hac model.

79 Genomic DNA from KOREF blood samples was extracted using QIAGEN Blood & Cell Culture  
80 DNA Kit (cat no 13323). A total of 5 µg of each sample was used as input for library preparation.  
81 The SMRTbell library was constructed using the SMRTbell® Express Template Preparation Kit  
82 (101-357-000). Using the BluePippin Size selection system we removed the small fragments for a  
83 large-insert library. After sequencing primer v4 was annealed to the SMRTbell template, DNA  
84 polymerase was bound to the complex (Sequel Binding kit 2.0). We purified the complex using  
85 AMPure Purification to remove excess primer and polymerase prior to sequencing. The SMRTbell  
86 library was sequenced using SMRT cells (Pacific Biosciences) using Sequel Sequencing Kit v2.1  
87 and 10 hr movies were captured for each SMRT Cell 1M v2 using the Sequel II (Pacific  
88 Biosciences) sequencing platform.

89 Hi-C libraries were generated using the Arima-Hic kit (A160105v01, San Diego, CA, USA).  
90 KOREF cell lines and blood samples were prepared for the construction of Hi-C libraries. Briefly,  
91 chromatin from cross-linked cells was solubilized and then digested using restriction enzymes  
92 MboI or Arima's multiple enzymes (GATC and GANTC). The digested ends were labeled using



93 a biotinylated nucleotide, and ends were ligated to create ligation products. Ligation products were  
94 purified, fragmented, and selected by size using AMPure XP Beads. Illumina-compatible  
95 sequencing libraries were constructed on end repair, dA-tailing, and adaptor ligation using a  
96 modified workflow of the Hyper Prep kit (KAPA Biosystems, Inc.). The bead-bound libraries were  
97 amplified and purified using AMPure XP beads and sequenced using Illumina NovaSeq platform  
98 with a read-length of 150 bp by Novogene (Beijing, China).

99 Short paired-end raw reads using Illumina HiSeq 2000 platform were acquired from a previous  
100 study, accession no. SRR2204706 (<ftp://ftp.sra.ebi.ac.uk/vol1/srr/SRR220/006/SRR2204706>).

101 For generating parental sequencing reads, we prepared samples from both of KOREF\_S1's parents.  
102 DNA was extracted from the donor's blood using DNAeasy Blood & Tissue Kit from QIAGEN  
103 according to the manufacturer's instruction. The quality and concentration of the extracted DNA  
104 were evaluated using NanoDrop™ One/OneC UV-Vis Spectrophotometer (Thermo Scientific™).  
105 Library construction and whole-genome sequencing were performed by Illumina HiSeq platform  
106 (Illumina, USA) with a 100 bp paired-end sequencing.

107

## 108 **Preprocessing of sequenced reads**

109 The sequenced long- and short-reads data were performed preprocessing steps as adapter trimming,  
110 quality trimming, and error correction. For the long reads, adapter trimming was performed using  
111 Porechop v0.2.4 (<https://github.com/rrwick/Porechop>) (Porechop, RRID:SCR\_016967) and  
112 removing reads with below quality-score 7 was performed using Guppy. For the short reads,  
113 adapter- and quality trimming were performed using Trimmomatic v0.39 [5] (Trimmomatic,  
114 RRID:SCR\_011848), and an error correction was performed using the tadpole.sh program of

115 BBtools suite v38.26 (<https://jgi.doe.gov/data-and-tools/bbtools>) (Bestus Bioinformaticus Tools,  
116 RRID:SCR\_016968).

117

## 118 **Trio-binning and read correction**

119 To obtain more accurate and longer haplotype-resolved reads from ONT PromethION sequencing,  
120 we applied a trio-binning with KOREF's parental sequencing data and an error-correction with  
121 PacBio Hifi sequencing data. The whole procedure is described in figure 1. To obtain haplotype-  
122 resolved reads from ONT PromethION and PacBio Hifi sequencing, we performed a trio-binning  
123 using TrioCanu v2.1 [6] (Canu, RRID:SCR\_015880) with the parental short-reads. In this step,  
124 reads from eleven PromethION flow-cells and six PacBio Hifi cells were participated. We merged  
125 unclassified reads to the classified paternal-reads and classified maternal-reads each. To correct  
126 base-errors on the PromethION reads, we corrected the errors with the haplotype-resolved reads  
127 from PacBio Hifi sequencing using Racon v1.4.3 [7] (Racon, RRID:SCR\_017642). We acquired  
128 KOREF's parental sequencing data from the KOREF homepage  
129 ([http://koref.net/KOREF\\_Data\\_Download](http://koref.net/KOREF_Data_Download)).

130

## 131 ***De novo* assembly of KOREF\_S1 genome**

132 Contig assembly was processed using wtdbg2 v2.5 [8] (WTDBG, RRID:SCR\_017225) and Flye  
133 assembler v2.8.1 [9] (Flye, RRID:SCR\_017016). For a wtdbg2 assembly, parameters were set as  
134 '-x corrected -g 3g -L 5000 -X 70.0'. An error correction of the assembled contigs was conducted  
135 using Racon with a single iteration. The Flye assembly was performed with parameters of '--

136 pacbio-hifi --hifi-error 0.008 --genome-size 3g'. For error correction, we carried out the same  
137 procedure as the wtdbg2 assembly.

138 To construct scaffolds with a chromosome-scale, we conduct scaffolding using PromethION reads  
139 and Hi-C data. To scaffold contigs using PromethION reads, LINKS v1.8.7 [10] was used with a  
140 single flow-cell of PromethION reads. To construct chromosome-scale scaffolds using Hi-C data,  
141 3D-DNA pipeline v180922 [11] with Juicer v1.6.2 program [12] (Juicer, RRID:SCR\_017226) was  
142 performed with the scaffolds by LINKS. Hi-C raw reads were mapped against the extended contigs  
143 using Juicer, and the 3D-DNA pipeline was initiated to correct mis-joined contigs and construct  
144 scaffolds. To correct misassemblies on the scaffolds, a manual curation was performed using JBAT  
145 (JuiceBox Assembly Tool) v1.11.08 program (<https://github.com/aidenlab/Juicebox>) (Juicebox,  
146 RRID:SCR\_021172). To polish base-errors and small indels, we performed Pilon v1.23 program  
147 [13] (Pilon, RRID:SCR\_014731) with KOREF's short read data and parameters of '--fix snps and  
148 indels' were used.

149

## 150 **Constructing high-confident regions, and the assessment of base-errors on** 151 **long-reads and genome assemblies**

152 For an assessment of base-errors, we constructed high-confident regions of KOREF\_S1 v1 against  
153 chromosome sequences of the GRCh38.p13. The procedure was referred to Heng Li's study [14].  
154 We aligned the KOREF\_S1v1.0 assembly to GRCh38 using the Minimap2 program v2.17-r941  
155 [15] (Minimap2, RRID:SCR\_018550). Alignments with mapping quality >5 and aligned segments  
156 shorter than 50 kb were discarded. The filtered alignments were converted to the BED format and  
157 sorted.

158 To assess base-errors of long-reads and genome assemblies, we compared them to the  
159 KOREF\_S1v1.0 assembly using the assembly\_assess program from Pomoxis v0.3.4  
160 (<https://github.com/nanoporetech/pomoxis>). And the Merqury v1.0 [16] program was performed  
161 to assess assemblies using k-mers.

162

## 163 **Genome annotation**

164 To identify protein coding genes on KOREF\_S1v2.1 genome, we performed a reference-guided  
165 transcriptome assembly with liftover with a gene annotation from GENCODE 38. The liftover was  
166 processed using Liftoff v1.6.1 program [17], and the reference-guided assembly was performed  
167 using Stringtie v2.1.5 program [18] (Stringtie, RRID:SCR\_016323) and TransDecoder v5.5.0  
168 program (<https://github.com/TransDecoder/TransDecoder>) (TransDecoder, RRID:SCR\_017647)  
169 with KOREF's RNASeq data. Mapping RNASeq data was conducted using HISAT2 v2.2.1  
170 program [19] (HISAT2, RRID:SCR\_015530) for short reads and GMAP v2020-06-01 program  
171 [20] (GMAP, RRID:SCR\_008992) for PacBio ISOSeq reads. Other genes including lncRNAs and  
172 pseudo-genes were annotated by the liftover. The result of genome annotation was stored in the  
173 KOREF genome browser, built by the JBrowse v1.16.9 [21] (JBrowse, RRID:SCR\_001004). All  
174 RNASeq data were collected from the KOREF homepage  
175 ([http://koref.net/KOREF\\_Data\\_Download](http://koref.net/KOREF_Data_Download)).

176

## 177 **Results**

### 178 **KOREF\_S1v2.1 assembly**

179 We obtained 235× coverage (705 Gb) of long-reads from twelve ONT PromethION flow-cells and  
180 38× coverage (114 Gb) of long reads from six PacBio HiFi cells (Table S1). We also acquired 274  
181 Gb corrected paternal haplotype-resolved reads and 265 Gb corrected maternal haplotype-resolved  
182 reads after trio-binning and read-correction. The corrected reads were identified about 1.4% base-  
183 errors (Table S2). Contigs from both haplotypes were assembled using wtdbg and Flye. The Flye  
184 assembly showed better results of higher N50 values (19.47 Mb for a paternal and 25.86 Mb for a  
185 maternal assembly) and longer length of the longest contig (70.97 Mb for a paternal and 109.79  
186 Mb for a maternal assembly) (Table 1).

187 We extended the contigs to chromosome-scale scaffolds using 76.5 Gb of PromethION reads  
188 (Flow-cell no.2) and 884 Gb of Hi-C data (294× sequencing-depth). Scaffolds from a  
189 mitochondrial genome were excluded using the KOREF's mtDNA sequence from the previous  
190 study [4]. As a result, we acquired the paternal assembly of 2.82 Gb length with 2,230 scaffolds  
191 and an N50 of 141.04 Mb (Table 1). The maternal assembly resulted in 2,616 scaffolds with an  
192 N50 of 150.05 Mb, and its total length was 2.88 Gb. For generating the final assembly of  
193 KOREF\_S1v2.1, we substituted sequences of autosomal chromosomes and a Y chromosome from  
194 the paternal assembly, and a X chromosome from the maternal assembly. As a result, the  
195 KOREF\_S1v2.1 was acquired a total length of 2.9 Gb with an N50 of 150.05 Mb.

196

## 197 **Genome annotation**

198 We annotated genes in KOREF\_S1v2.1 by integrating a liftover of gene annotations from the  
199 GENCODE release 38 (<https://www.gencodegenes.org/human/>) and homology information of  
200 RNASeq data. The genes included 20,378 protein-coding genes with 166,570 transcripts, 46,973

201 lncRNAs and 17,535 pseudogenes (Table 3). 1,391 genes from the Gencode38 annotation were  
202 not transferred to the KOREF by liftover, and a list of these genes can be found in the  
203 supplementary table 4.

204

## 205 **Assessment of KOREF and comparison with other human genome assemblies**

206 Using the Merqury program for a quality assessment, we estimated QV scores of Q43.88 for the  
207 paternal assembly and Q44.49 for the maternal assembly. The final assembly showed QV score of  
208 Q43.88, indicating >99.99% accuracy (Table S5), and it is higher than KOREF\_S1v1.0's (Q33.58)  
209 and KOREF\_S1v2.0 (Q39.52) which were assembled with the PromethION data. We  
210 compared KOREF\_S1v2.1 and other human reference genome assemblies (AK1\_v2, JG2.0.0 Beta,  
211 HuRef, CHM13\_v1.1, and GRCh38.p13) [22-25]. The results showed that KOREF\_S1v2.1 is  
212 more contiguous than AK1 and HuRef, and comparable to JG2.0.0 Beta and CHM13\_v1.1 (Table  
213 2). Among six genome assemblies, KOREF\_S1v2.1 and CHM13 were a haplotype-resolved  
214 assembly with a chromosome-scale. AK1 was haplotype-resolved using a read-based phasing  
215 method but could not reach a chromosome-scale without a guidance of the reference genome.

216

## 217 **Discussion**

218 In previous version of KOREF\_S1, we generated a chromosome-level genome assembly with a  
219 guidance of GRCh38. A new version of KOREF assembly, KOREF\_S1v2.1, was assembled with  
220 high accurate (less than 0.01% of base error) and contiguity from multiple sequencing technologies  
221 including ONT, PacBio, Illumina, and Hi-C. Furthermore, the new KOREF assembly was phased

222 with parental sequencing data. To generate ultra-long and high accurate reads, we corrected ONT  
223 reads using PacBio HiFi reads. Most genomic regions were covered by the corrected reads, but  
224 some highly competitive regions including telomere and centromere were not covered. They were  
225 remained as gaps with unknown length. Especially on a chromosome Y, we found more gaps and  
226 less contiguity than other chromosomes. The genomic sequences of a chromosome X and Y have  
227 high similar regions and they probably make difficulties to phase genomic sequences on sex  
228 chromosomes.

229 Recently, new *de novo* assembly pipelines, such as the Hifiasm [26] (Hifiasm, RRID:SCR\_021069)  
230 and HiCanu [27], have been developed for PacBio's HiFi-CCS. Especially, Hifiasm supports a  
231 trio-binning from parental sequencing and Hi-C, and the assembly resulted in high base-accuracy  
232 and contiguity (Table S3). Despite these advantages, scaffolding contigs from Hifiasm has  
233 difficulties for using Hi-C data. From a pilot study, an error-correction module of the 3D-DNA  
234 pipeline seemed to split long repetitive regions complicatedly, and it made difficult to construct  
235 scaffolds or curate misassemblies (Fig. S1). However, the high-quality contigs from Hifiasm can  
236 be helpful to remove gaps and resolve highly repetitive regions. Also, a recent study of the T2T  
237 consortium shared a complete structure of centromeric regions [25], and it will be a useful resource  
238 to complete the KOREF\_S1 genome.

239 In conclusion, we upgraded a high-quality Korean reference genome, KOREF. Our study provides  
240 useful resources of the Korean reference genome and demonstrates a new strategy of hybrid  
241 assembly which collaborates ONT's PromethION and PacBio's HiFi-CCS.

242

243 **Data availability**

244 The Korean reference genome project has been deposited at DDBJ/ENA/GenBank under the  
245 accession PRJNA735947. The version described in this paper is version JHRJT000000000. Raw  
246 DNA and RNA sequence reads for KOREF and KPGP have been submitted to the NCBI Sequence  
247 Read Archive database (from SRR14759111 to SRR14759134). The immortalized cell line of  
248 KOREF was deposited in the Korean Cell Line Bank (KCLB, #60211). KOREF\_S1 data is found  
249 from <http://koreanreference.org>

250

## 251 **Competing financial interests**

252 The authors declare no competing financial interests.

253

## 254 **Funding**

255 This work was supported by the Promotion of Innovative Businesses for Regulation-Free Special  
256 Zones funded by the Ministry of SMEs and Startups (MSS, Korea)(P0016193).

257

## 258 **Author contributions**

259 J.B. supervised and coordinated the national Korean reference genome project and Personal  
260 Genome Project Korea. J.B. conceived and designed the reference genome project. H.K.  
261 performed the analyses and assembly. H.K. and J.B. wrote the manuscript.

262

## 263 **Acknowledgements**



264 This work was supported by the Ministry of Trade, Industry & Energy (MOTIE, Korea) under  
265 Industrial Technology Innovation Programs ('Pilot study of building of Korean Reference  
266 Standard Genome map', No.10046043; 'Developing Korean Reference Genome', No.10050164;  
267 and 'National Center for Standard Reference Data', No.10063239) and Industrial Strategic  
268 Technology Development Program ('Bioinformatics platform development for next generation  
269 bioinformation analysis', No.10040231). Korea Institute of Science and Technology Information  
270 (KISTI) provided us with Korea Research Environment Open NETwork (KREONET) which is  
271 the internet connection service for efficient information and data transfer. We thank Jaesu Bhak  
272 for editing the manuscript.

273

## 274 **References**

- 275 1. Schneider VA, Graves-Lindsay T, Howe K, Bouk N, Chen HC, Kitts PA, et al. Evaluation of  
276 GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of  
277 the reference assembly. *Genome Res.* 2017;27 5:849-64. doi:10.1101/gr.213611.116.
- 278 2. Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, et al. A Draft Sequence of  
279 the Neandertal Genome. *Science.* 2010;328 5979:710-22. doi:10.1126/science.1188021.
- 280 3. Logsdon GA, Vollger MR and Eichler EE. Long-read human genome sequencing and its  
281 applications. *Nat Rev Genet.* 2020;21 10:597-614. doi:10.1038/s41576-020-0236-x.
- 282 4. Cho YS, Kim H, Kim HM, Jho S, Jun J, Lee YJ, et al. An ethnically relevant consensus Korean  
283 reference genome is a step towards personal reference genomes. *Nat Commun.* 2016;7  
284 doi:ARTN 13637 10.1038/ncomms13637.
- 285 5. Bolger AM, Lohse M and Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence  
286 data. *Bioinformatics.* 2014;30 15:2114-20. doi:10.1093/bioinformatics/btu170.
- 287 6. Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, et al. De novo assembly  
288 of haplotype-resolved genomes with trio binning. *Nat Biotechnol.* 2018;  
289 doi:10.1038/nbt.4277.
- 290 7. Vaser R, Sovic I, Nagarajan N and Sikic M. Fast and accurate de novo genome assembly  
291 from long uncorrected reads. *Genome Res.* 2017;27 5:737-46.  
292 doi:10.1101/gr.214270.116.
- 293 8. Ruan J and Li H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods.*  
294 2020;17 2:155-+. doi:10.1038/s41592-019-0669-3.

- 295 9. Kolmogorov M, Yuan J, Lin Y and Pevzner PA. Assembly of long, error-prone reads using  
296 repeat graphs. *Nat Biotechnol.* 2019;37 5:540-+. doi:10.1038/s41587-019-0072-8.
- 297 10. Warren RL, Yang C, Vandervalk BP, Behsaz B, Lagman A, Jones SJM, et al. LINKS: Scalable,  
298 alignment-free scaffolding of draft genomes with long reads. *Gigascience.* 2015;4  
299 doi:ARTN 35 10.1186/s13742-015-0076-3.
- 300 11. Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, et al. De novo  
301 assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds.  
302 *Science.* 2017;356 6333:92-5. doi:10.1126/science.aal3327.
- 303 12. Durand NC, Shamim MS, Machol I, Rao SS, Huntley MH, Lander ES, et al. Juicer Provides a  
304 One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Syst.* 2016;3 1:95-  
305 8. doi:10.1016/j.cels.2016.07.002.
- 306 13. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An  
307 Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly  
308 Improvement. *Plos One.* 2014;9 11 doi:ARTN e112963 10.1371/journal.pone.0112963.
- 309 14. Li H, Bloom JM, Farjoun Y, Fleharty M, Gauthier L, Neale B, et al. A synthetic-diploid  
310 benchmark for accurate variant-calling evaluation. *Nat Methods.* 2018;15 8:595-7.  
311 doi:10.1038/s41592-018-0054-7.
- 312 15. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018;34  
313 18:3094-100. doi:10.1093/bioinformatics/bty191.
- 314 16. Shumate A and Salzberg SL. Liftoff: accurate mapping of gene annotations.  
315 *Bioinformatics.* 2020; doi:10.1093/bioinformatics/btaa1016.
- 316 17. Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT and Salzberg SL. StringTie  
317 enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat*  
318 *Biotechnol.* 2015;33 3:290-5. doi:10.1038/nbt.3122.
- 319 18. Kim D, Paggi JM, Park C, Bennett C and Salzberg SL. Graph-based genome alignment and  
320 genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol.* 2019;37 8:907-15.  
321 doi:10.1038/s41587-019-0201-4.
- 322 19. Wu TD and Watanabe CK. GMAP: a genomic mapping and alignment program for mRNA  
323 and EST sequences. *Bioinformatics.* 2005;21 9:1859-75.  
324 doi:10.1093/bioinformatics/bti310.
- 325 20. Buels R, Yao E, Diesh CM, Hayes RD, Munoz-Torres M, Helt G, et al. JBrowse: a dynamic  
326 web platform for genome visualization and analysis. *Genome Biol.* 2016;17:66.  
327 doi:10.1186/s13059-016-0924-1.
- 328 21. Rhie A, Walenz BP, Koren S and Phillippy AM. Merqury: reference-free quality,  
329 completeness, and phasing assessment for genome assemblies. *Genome Biol.* 2020;21  
330 1:245. doi:10.1186/s13059-020-02134-9.
- 331 22. Seo JS, Rhie A, Kim J, Lee S, Sohn MH, Kim CU, et al. De novo assembly and phasing of a  
332 Korean human genome. *Nature.* 2016;538 7624:243-+. doi:10.1038/nature20098.
- 333 23. Takayama J, Tadaka S, Yano K, Katsuoka F, Gocho C, Funayama T, et al. Construction and  
334 integration of three de novo Japanese human genome assemblies toward a population-  
335 specific reference. *Nat Commun.* 2021;12 1:226. doi:10.1038/s41467-020-20146-8.
- 336 24. Levy S, Sutton G, Ng PC, Feuk L, Halpern AL, Walenz BP, et al. The diploid genome  
337 sequence of an individual human. *PLoS Biol.* 2007;5 10:e254.  
338 doi:10.1371/journal.pbio.0050254.

- 339 25. Cheng H, Concepcion GT, Feng X, Zhang H and Li H. Haplotype-resolved de novo  
340 assembly using phased assembly graphs with hifiasm. *Nat Methods*. 2021;18 2:170-5.  
341 doi:10.1038/s41592-020-01056-5.
- 342 26. Nurk S, Walenz BP, Rhie A, Vollger MR, Logsdon GA, Grothe R, et al. HiCanu: accurate  
343 assembly of segmental duplications, satellites, and allelic variants from high-fidelity long  
344 reads. *Genome Res*. 2020;30 9:1291-305. doi:10.1101/gr.263566.120.
- 345 27. Nurk S, Koren S, Rhie A, Rautiainen M, Bzikadze AV, Mikheenko A, et al. The complete  
346 sequence of a human genome. *bioRxiv*. 2021.  
347

348

349

350 **Figures**

351

352 **Figure 1. The flowchart of KOREF genome assembly**

353

354 **Tables**

355

356 **Table 1. The statistics of KOREF\_S1v2.1 assembly**

|                   | Contig          |               |                 |               | Scaffold      |               |
|-------------------|-----------------|---------------|-----------------|---------------|---------------|---------------|
|                   | Wtdbg2_paternal | Flye_paternal | Wtdbg2_maternal | Flye_maternal | Paternal      | Maternal      |
| Sequence no.      | 3,059           | 4,463         | 2,426           | 2,475         | 2,230         | 2,616         |
| Total length (bp) | 2,652,350,533   | 2,820,210,305 | 2,691,371,348   | 2,885,670,065 | 2,821,407,033 | 2,886,600,011 |
| N50 (bp)          | 15,085,508      | 19,472,363    | 15,312,743      | 25,861,606    | 141,044,433   | 150,051,441   |
| Longest (bp)      | 70,969,653      | 87,371,841    | 70,444,093      | 109,786,075   | 235,665,501   | 234,237,609   |
| Gaps              | 0.000%          | 0.000%        | 0.000%          | 0.000%        | 0.048%        | 0.037%        |
| GC contents       | 40.90%          | 40.92%        | 40.84%          | 40.86%        | 40.92%        | 40.88%        |

357

358 **Table 2. Comparison between KOREF and other human genomes**

|                    | KOREF_S1v2.1   | AK1_v2         | JG2.0.0 Beta   | HuRef            | CHM13             | GRCh38.p13     |
|--------------------|----------------|----------------|----------------|------------------|-------------------|----------------|
| Scaffold no.       | 2,230          | 2,832          | 1,173          | 4,530            | 24                | 472            |
| Total length (bp)  | 2,901,828,151  | 2,904,207,228  | 3,059,652,438  | 2,844,000,504    | 3,054,832,041     | 3,272,089,205  |
| Scaffold N50 (bp)  | 150,051,441    | 44,846,623     | 152,668,378    | 143,733,266      | 154,259,566       | 67,794,783     |
| Phasing approach   | <i>De novo</i> | <i>De novo</i> | <i>De novo</i> | Reference-guided | <i>De novo</i>    | <i>De novo</i> |
| Assembly level     | Chromosome     | Scaffold       | Chromosome     | Chromosome       | Chromosome        | Chromosome     |
| Haplotype-resolved | Trio-binning   | Read-based     | No             | No               | Haploid cell line | No             |

359

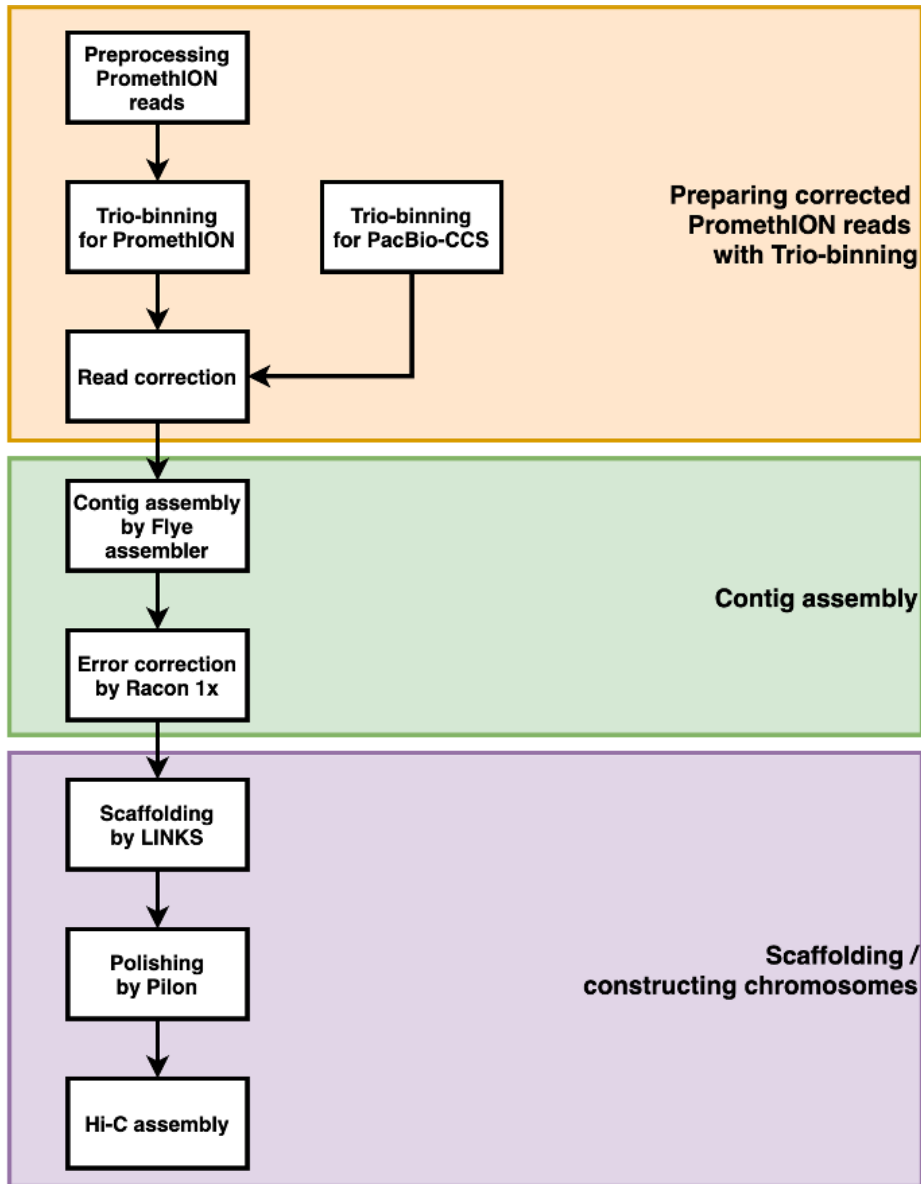
360

361 **Table 3. The statistics of KOREF genome annotation**

|                                  | <b>KOREF_S1v2.1 gene</b> |
|----------------------------------|--------------------------|
| Genes no.                        | 20,378                   |
| Transcripts no.                  | 166,570                  |
| Total length of transcripts (bp) | 216,532,041              |
| N50 (bp)                         | 1,851                    |
| Length of longest transcripts    | 107,976                  |
| GC contents                      | 52.22%                   |
| lncRNAs no.                      | 46,973                   |
| Pseudogenes no.                  | 17,535                   |

362

363 **Figure 1**



364

365



[Click here to access/download](#)

**Supplementary Material**

[KOREF\\_S1v2.1\\_supplementary\\_tables.20210726.xlsx](#)

