

GigaScience

Chromosome-scale assembly comparison of the Korean Reference Genome KOREF from PromethION and PacBio with Hi-C mapping information --Manuscript Draft--

Manuscript Number:	GIGA-D-19-00240	
Full Title:	Chromosome-scale assembly comparison of the Korean Reference Genome KOREF from PromethION and PacBio with Hi-C mapping information	
Article Type:	Data Note	
Funding Information:	Ulsan National Institute of Science and Technology (1.190007.01)	Dr. Jong Bhak
	Ulsan National Institute of Science and Technology (1.190033.01)	Dr. Jong Bhak
Abstract:	<p>Background</p> <p>Long DNA reads produced by single molecule and pore-based sequencers are more suitable for assembly and structural variation discovery than short read DNA fragments. For de novo assembly, PacBio and Oxford Nanopore Technologies (ONT) are favorite options. However, PacBio's SMRT sequencing is expensive for a full human genome assembly and costs over 40,000 USD for 30x coverage as of 2019. ONT PromethION sequencing, on the other hand, is one-twelfth the price of PacBio for the same coverage. This study aimed to compare the cost-effectiveness of ONT PromethION and PacBio's SMRT sequencing in relation to the quality.</p> <p>Findings</p> <p>We performed whole genome de novo assemblies and comparison to construct an improved version of KOREF, the Korean reference genome, using sequencing data produced by PromethION and PacBio. With PromethION, an assembly using sequenced reads with 64x coverage (193 Gb, 3 flowcell sequencing) resulted in 3,725 contigs with N50s of 16.7 Mbp and a total genome length of 2.8 Gbp. It was comparable to a KOREF assembly constructed using PacBio at 62x coverage (188 Gbp, 2,695 contigs and N50s of 17.9 Mbp). When we applied Hi-C-derived long-range mapping data, an even higher quality assembly for the 64x coverage was achieved, resulting in 3,179 scaffolds with an N50 of 56.4 Mbp.</p> <p>Conclusion</p> <p>The pore-based PromethION approach provides a good quality chromosome-scale human genome assembly at a low cost with long maximum contig and scaffold lengths and is more cost-effective than PacBio at comparable quality measurements.</p>	
Corresponding Author:	Jong Hwa Bhak, Ph.D. UNIST Ulsan, Ulsan KOREA, REPUBLIC OF	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	UNIST	
Corresponding Author's Secondary Institution:		
First Author:	Hui-Su Kim, Ph.D.	
First Author Secondary Information:		
Order of Authors:	Hui-Su Kim, Ph.D.	
	Sungwon Jeon, B.S.	

	Changjae Kim, Ph.D.
	Yeon Kyung Kim, M.S.
	Yun Sung Cho, Ph.D.
	Jungeun Kim, Ph.D.
	Asta Blazyte, B.S.
	Andrea Manica, Ph.D.
	Semin Lee, Ph.D.
	Jong Bhak, Ph.D.
Order of Authors Secondary Information:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. 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>Resources</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 Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
Availability of data and materials	Yes

All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in [publicly available repositories](#) (where available and ethically appropriate), referencing such data using 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](#)?

[Click here to view linked References](#)

1 **Chromosome-scale assembly comparison of the Korean Reference Genome**
2 **KOREF from PromethION and PacBio with Hi-C mapping information**

3

4 Hui-Su Kim¹, Sungwon Jeon^{1,2}, Changjae Kim¹, Yeon Kyung Kim¹, Yun Sung Cho³, Jungeun
5 Kim⁵, Asta Blazyte¹, Andrea Manica⁴, Semin Lee^{1,2*}, Jong Bhak^{1,2,3,5*}

6

7 ¹KOGIC, Ulsan National Institute of Science and Technology (UNIST), Ulsan 44919, Republic
8 of Korea

9 ²Department of Biomedical Engineering, School of Life Sciences, UNIST, Ulsan 44919,
10 Republic of Korea

11 ³Clinomics Inc., Ulsan 44919, Republic of Korea

12 ⁴Department of Zoology, Cambridge, University, Cambridge, UK

13 ⁵Personal Genomics Institute, Genome Research Foundation, Cheongju 28160, Republic of
14 Korea

15

16

17 ***Correspondence author:**

18 Name: Jong Bhak, Ph.D.

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

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

2 Email: jongbhak@genomics.org, ORCID: 0000-0002-4228-1299

3 Name: Semin Lee, Ph.D.

4 Address: #110-302, Ulsan National Institute of Science and Technology, UNIST-gil 50,

5 Eonyang-eup, Ulju-gun, Ulsan 44919, Republic of Korea

6 Phone: (+82) 52-217-2663

7 Email: seminlee@gmail.com, ORCID: 0000-0002-9015-6046

8

1 **Abstract**

2 **Background:** Long DNA reads produced by single molecule and pore-based sequencers are
3 more suitable for assembly and structural variation discovery than short read DNA fragments.
4 For *de novo* assembly, PacBio and Oxford Nanopore Technologies (ONT) are favorite options.
5 However, PacBio's SMRT sequencing is expensive for a full human genome assembly and
6 costs over 40,000 USD for 30× coverage as of 2019. ONT PromethION sequencing, on the
7 other hand, is one-twelfth the price of PacBio for the same coverage. This study aimed to
8 compare the cost-effectiveness of ONT PromethION and PacBio's SMRT sequencing in
9 relation to the quality.

10 **Findings:** We performed whole genome *de novo* assemblies and comparison to construct an
11 improved version of KOREF, the Korean reference genome, using sequencing data produced
12 by PromethION and PacBio. With PromethION, an assembly using sequenced reads with 64×
13 coverage (193 Gb, 3 flowcell sequencing) resulted in 3,725 contigs with N50s of 16.7 Mbp and
14 a total genome length of 2.8 Gbp. It was comparable to a KOREF assembly constructed using
15 PacBio at 62× coverage (188 Gbp, 2,695 contigs and N50s of 17.9 Mbp). When we applied Hi-
16 C-derived long-range mapping data, an even higher quality assembly for the 64× coverage was
17 achieved, resulting in 3,179 scaffolds with an N50 of 56.4 Mbp.

18 **Conclusion:** The pore-based PromethION approach provides a good quality chromosome-
19 scale human genome assembly at a low cost with long maximum contig and scaffold lengths
20 and is more cost-effective than PacBio at comparable quality measurements.

21

22 **Keywords:** Korean reference genome; KOREF, PromethION; Hi-C; nanopore sequencing,
23 single molecule sequencing

1 **Data Description**

2 Next-generation sequencing (NGS) is a set of powerful sequencing technologies and a recent
3 trend in genomics is to use cost-effective long DNA reads for assembly and structural variation
4 discovery using single molecule sequencing methods. Oxford Nanopore Technologies (ONT)
5 and PacBio platforms have advantages of a short run time and long read lengths over short
6 fragmented reads by Illumina [1, 2]. Unfortunately, both methods share high base-calling error
7 rates [3, 4]. However, bioinformatics pipelines for self-error correction and/or polishing
8 sequences with short reads have become an effective option, and the overall accuracy of long
9 read based assemblies is approaching what is required to be a viable option for personal
10 reference genome construction [5]. Despite its excellent performance, PacBio's SMRT
11 sequencing is expensive for the effective coverage required for a full human genome assembly,
12 costing over 40,000 USD for 30× coverage (with 15 SMRT cells; from an estimated 6 Gbp raw
13 reads production per SMRT cell) as of 2019 [6, 7, 8]. On the other hand, the nanopore based
14 single molecule, long read platform, PromethION from ONT is highly cost-effective at one-
15 twelfth the price of PacBio's for the same read amount, with an advantage of even longer
16 average and maximum read lengths [9]. Although the two methods share some similarity, they
17 are fundamentally different in that ONT uses a minimal amount of reagents with small form
18 factor devices, and can be a promising future technology for a very broad scope of applications
19 given its advantageous size and cost.

20 In this study, we performed benchmark tests of PromethION and PacBio with low and
21 high coverages of sequencing data and investigated the advantages of pairing these long read
22 technologies with very long-range chromosome mapping information by Hi-C, using the
23 already existing high-quality Korean reference genome, KOREF, as a benchmark [10].

24

1 **Whole genome sequencing by ONT PromethION R9.4.1 platform**

2 Human KOREF cell lines (<http://koref.net>) were cultured at 37°C in 5% CO₂ in RPMI-1640
3 medium with 10% heat-inactivated fetal bovine serum. DNA was extracted from cells using
4 the DNeasy Blood & Tissue kit (Qiagen). The KOREF cells (5×10^6) were centrifuged at 300
5 g for 5 min; the pelleted cells were suspended in 200 µL of PBS and DNA was extracted
6 according to the manufacturer's instructions. To preserve large-sized DNA and purify DNA
7 fragments, we used Genomic DNA Clean & Concentrator kit (Zymo). The DNA quality and
8 size were assessed by running 1 µL of purified DNA on the Bioanalyzer system (Agilent).
9 Concentration of DNA was assessed using the dsDNA BR assay on a Qubit fluorometer
10 (Thermo Fisher).

11 DNA repair (NEBNext FFPE DNA Repair Mix, NEB M6630) and end-prep
12 (NEBNext End Repair/dA-tailing, NEB E7546) were performed using 1 µg human genomic
13 DNA. The mixture of 1 µL DNA CS, 3.5 µL FFPE Repair Buffer, 2 µL FFPE DNA Repair
14 Mix, 3.5 µL Ultra II End-prep reaction buffer, and 3 µL Ultra II End-prep enzyme mix was
15 added to 47 µL DNA sample. The final mixture was incubated at 20°C for 5 min and then at
16 65°C for 5 min, cleaned up using 60 µL AMPure XP beads, incubated on Hula mixer for 5 min
17 at room temperature, and washed twice with 200 µL fresh 70% ethanol. The pellet was allowed
18 to dry for 30 s, and then DNA was eluted in 61 µL of nuclease-free water. An aliquot of 1 µL
19 was quantified by Qubit to ensure ≥ 1 µg DNA was retained.

20 Adaptor ligation was performed by adding 5 µL of Adaptor Mix (AMX, SQK-LSK109
21 Ligation Sequencing Kit 1D, Oxford Nanopore Technologies (ONT)), 25 µL Ligation Buffer
22 (LNB, SQK-LSK109), and 10 µL NEBNext Quick T4 DNA Ligase (NEB, E6056) to 60 µL
23 bead cleaned-up DNA, followed by gentle mixing and incubation for 10 min at room
24 temperature.

1 The adaptor-ligated DNA was cleaned up by adding 40 μ L of AMPure XP beads,
2 incubating for 5 min at room temperature and re-suspending the pellet twice in 250 μ L L
3 Fragment Buffer (LFB, SQK-LSK109). The purified ligated DNA was re-suspended in 25 μ L
4 of Elution Buffer (ELB, SQK-LSK109), incubated for 10 min at room temperature, followed
5 by pelleting the beads, and transferring the supernatant (pre-sequencing mix or PSM) to a new
6 Eppendorf Lobind tube. A 1- μ L aliquot was quantified by Qubit to ensure \geq 500 ng DNA was
7 retained.

8 To load the library, 75 μ L of Sequencing Buffer (SQB, SQK-LSK109) was mixed with
9 51 μ L of Loading Beads (LB, SQK-LSK109) and this mixture was added to 24 μ L DNA library.
10 This library was mixed by pipetting slowly and 150 μ L of sample was loaded through the inlet
11 port.

12

13 **Whole genome sequencing by PacBio Sequel platform**

14 Genomic DNA was extracted from human KOREF blood samples using QIAGEN Blood &
15 Cell Culture DNA Kit (cat no 13323). A total of 5 μ g of each sample was used as input for
16 library preparation. The SMRTbell library was constructed using SMRTbell® Express
17 Template Preparation Kit (101-357-000). Using the BluePippin Size selection system we
18 removed the small fragments for large-insert library. After sequencing primer v4 was annealed
19 to the SMRTbell template, DNA polymerase was bound to the complex (Sequel Binding kit
20 2.0). We purified the complex using AMPure Purification to remove excess primer and
21 polymerase prior to sequencing. The SMRTbell library was sequenced using SMRT cells
22 (Pacific Biosciences) using Sequel Sequencing Kit v2.1 and 10 h movies were captured for
23 each SMRT Cell 1M v2 using the Sequel (Pacific Biosciences) sequencing platform.

1

2 **Short read sequencing by Illumina HiSeq**

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

5

6 **Hi-C chromosome conformation captured reads sequencing**

7 Long distance Hi-C chromosome conformation capture data were generated using the Arima-
8 HiC kit (A160105 v01), and double restriction enzymes were used for chromatin digestion. To
9 prepare KOREF cell line samples for Hi-C analysis, cells were harvested and cross-linked as
10 instructed by the manufacturer. One million cross-linked cells were used as input in the Hi-C
11 protocol. Briefly, chromatin from cross-linked cells or nuclei was solubilized, and then digested
12 using restriction enzymes A1 and A2. The digested ends were then labeled using a biotinylated
13 nucleotide, and ends were ligated to create ligation products. Ligation products were purified,
14 fragmented, and selected by size using AMPure XP Beads. Biotinylated fragments were then
15 enriched using Enrichment beads, and Illumina-compatible sequencing libraries were
16 constructed on End Repair, dA-tailing, and Adaptor Ligation using a modified workflow of the
17 Hyper Prep kit (KAPA Biosystems, Inc.). The bead-bound library was then amplified, and
18 amplicons were purified using AMPure XP beads and subjected to deep sequencing.

19

20 **Short and long sequence reads processing**

21 A total of 144 Gbp of short paired-end DNA raw reads were obtained from SRA2204706.
22 Adapter sequences were trimmed from sequenced raw reads using Trimmomatic v0.36 [11]

1 (ILLUMINACLIP:2:30:10 LEADING:5 TRAILING:5 SLIDINGWINDOW:4:20
2 HEADCROP:15 MINLEN:60) (Trimmomatic, RRID:SCR_011848), and screening for vectors
3 and microbial contaminants were performed using customized database from Refseq. After
4 preprocessing, a total of 137 Gbp cleaned reads were obtained.

5 A total of 80.7 Gbp and 193 Gbp raw reads (27× and 64× coverage) were obtained as
6 a result of PromethION nanopore sequencing using one and three flowcells. Removing adapter
7 sequences from the raw reads was performed using Porechop v0.2.4 (Porechop,
8 RRID:SCR_016967) [12]. We also acquired 92.2 Gbp and 187.9 Gbp raw reads from PacBio
9 Sequel sequencing resulting in 30× and 62× coverage (Table 1).

10

11 **Long-read sequence based *de novo* genome assemblies**

12 *De novo* assemblies for the 27× and 64× PromethION raw reads were performed using wtdbg2
13 v2.3 (WTDBG, RRID:SCR_017225) [13]. To compare the accuracy, two sets of raw reads with
14 30× and 62× coverage of PacBio Sequel were also used employing the same assembler.
15 Parameters for the assembler were set optimally for each sequencing platform with multiple
16 trials (https://github.com/macarima/KOREF_PromethION_paper). For self-error correction
17 with long reads, we generated consensus sequences using Racon v1.3.2 [14]. To improve the
18 accuracy of assemblies, polishing consensus sequences with 48.2× coverage short reads was
19 performed using Pilon v1.23 (Pilon, RRID:SCR_014731) [15]. To assess the completeness of
20 the long-read genome assemblies, BUSCO v3.0.2 (BUSCO, RRID:SCR_015008) [16] with the
21 default AUGUSTUS model for human was used to locate the presence and absence of 4,104
22 single copy orthologous genes from mammalian OrthoDB v9.

1 For constructing chromosome-scale assemblies for the PromethION long-reads data,
2 map assembly with Hi-C reads was performed using SALSA2 v2.2 [17]. Duplicated Hi-C reads
3 were removed using clumpify.sh program from BBTools suite v38.32 (Bestus Bioinformaticus
4 Tools, RRID:SCR_016968) [18]. Mapping Hi-C reads to the assembled genome was conducted
5 using the pipeline provided by Arima-Genomics
6 (https://github.com/ArimaGenomics/mapping_pipeline).

7 Long read assemblies from 27× and 64× PromethION sequencing yielded total
8 assembly sizes of 2,757 Mbp and 2,827 Mbp, with scaffold N50s of 7.6 Mbp and 16.7 Mbp,
9 respectively (Table 2). Assemblies from PacBio sequencing at 30× and 62× coverage yielded
10 the total assembly sizes of 2,800 Mbp and 2,815 Mbp, with scaffold N50s of 11.1 Mbp and
11 17.9 Mbp, respectively. Adding Hi-C reads to assemblies led to 3.4- to 4.3-fold increase in the
12 scaffold N50 lengths of PromethION (32.7 Mbp for 27× coverage and 56.4 Mbp for 64×
13 coverage). For the PacBio assemblies, 2.2- to 3.3-fold increase was achieved for the scaffold
14 N50 lengths (38.1 Mbp for 30× coverage and 59.3 Mbp for 62× coverage). The longest scaffold
15 from both PromethION and PacBio assemblies with Hi-C was two times the length of the
16 assemblies without Hi-C.

17

18 **Comparison between PromethION and PacBio assemblies**

19 The comparison between PromethION and PacBio assemblies without Hi-C mapping
20 information using sequenced reads at 64× coverage showed comparable quality. In terms of
21 N50, the PromethION assembly at 64× coverage yielded 1.5-fold and 0.93-fold longer N50s
22 compared with the PacBio assemblies at 30× and 62× coverage, respectively (Figure 1a). When
23 we compared the longest contigs, the PromethION assembly at 64× coverage yielded 1.7-fold
24 and 1.1-fold length increase compared with the PacBio assemblies at 30× and 62× coverage,

1 respectively (Figure 1b). Comparing the number of scaffolds, PacBio assembly at 30× coverage
2 showed the fewest (2,443) compared with that of PromethION assembly at 64× coverage
3 (3,725).

4 When Hi-C mapping information was added to the assembly construction, the
5 PromethION assembly at 64× coverage showed the best statistics as N50s of 56.4 Mbp and
6 the longest scaffold length of 175.2 Mbp. The PromethION assembly at 27× coverage with
7 Hi-C mapping information yielded 32.7 Mbp for N50s, which was comparable to both 30×
8 and 62× coverage PacBio assemblies with Hi-C; 0.85-fold and 0.55-fold for N50s,
9 respectively (Table 2).

10 When we compared assessment results from BUSCO, all the assemblies that had been
11 polished with short reads showed good quality; around 92% completed orthologous genes
12 with less than 1.1% completed and duplicated orthologous genes. Comparing the accuracy of
13 the assemblies to the single assembly of KOREF (KOREF_S), which is the current standard,
14 both showed around 99.8% accuracy (Table 3). The accuracy comparison was performed
15 using assess_assembly program from Pomoxis [19].

16

17 **Conclusions**

18 We generated high-quality assemblies of the Korean reference genome, KOREF, using ONT's
19 PromethION long-reads accompanied with Hi-C mapping information and compared them
20 against PacBio sequencing and assemblies of the same sample. Comparing the results from the
21 PromethION 64× sequencing to the PacBio 62× sequencing, we found that the former provided
22 high contiguity and completeness at one-twelfth the cost of PacBio. Results from just 27×
23 PromethION sequencing combined with Hi-C mapping information were also comparable to

1 the 30× coverage PacBio sequencing data. Therefore, to generate a chromosome-scale
2 assembly with a long-read technology, at present, the ONT's PromethION sequencing is a good
3 alternative to PacBio's, owing to its quality and cost-effectiveness. Simple pore-based long
4 read sequencing has potential to dramatically improve sequencing and subsequent
5 bioinformatics analysis for personal genome projects and cancer genome analyses where *de*
6 *novo* assemblies are necessary for structural and copy number variations that cannot be detected
7 easily by conventional short read only methods.

8

9 **Availability of supporting data**

10 Raw long-read sequencing data from PromethION and PacBio is available at NCBI genbank
11 under the project accession number PRJNA549351. All genome assemblies of KOREF are
12 available at KOREF website (<http://koref.net>).

13

14 **Abbreviations**

15 BUSCO: Benchmarking Universal Single-Copy Orthologs; PacBio: Pacific Biosciences;
16 SMRT: single-molecule real-time

17

18 **Competing interests**

19

20 Y.S.C. is an employee, and J.B. is the CEO of Clinomics Inc. J.B. and Y.S.C. have an equity

1 interest in the company. All other coauthors have no conflicts of interest to declare.

2

3 **Funding**

4 This work was supported by U-K BRAND Research Fund (1.190007.01) of UNIST; Research
5 Project Funded by Ulsan City Research Fund (1.190033.01) of UNIST and Clinomics internal
6 funding for KOREF sequencing using PromethION machine.

7

1 **Figure Legends**

2

3 **Figure 1.** Comparison of N50s and the longest contig/scaffold lengths for PromethION and
4 PacBio assemblies of KOREF

5

1 **References**

- 2 1. Mccarthy A. Third generation DNA sequencing: Pacific biosciences' single molecule
3 real time technology. Chem Biol [Internet]. Elsevier Ltd; 2010;17:675–6.
- 4 2. Laver T, Harrison J, O'Neill PA, Moore K, Farbos A, Paszkiewicz K, et al. Assessing
5 the performance of the Oxford Nanopore Technologies MinION. Biomol Detect
6 Quantif [Internet]. Elsevier GmbH; 2015;3:1–8.
- 7 3. Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Nonhybrid,
8 finished microbial genome assemblies from long-read SMRT sequencing data. Nat
9 Methods [Internet]. 2013;10:563–9.
- 10 4. Laver T, Harrison J, O'Neill PA, Moore K, Farbos A, Paszkiewicz K, et al. Assessing
11 the performance of the Oxford Nanopore Technologies MinION. Biomol Detect
12 Quantif [Internet]. Elsevier GmbH; 2015;3:1–8.
- 13 5. Fu S, Wang A, Au KF. A comparative evaluation of hybrid error correction methods
14 for error-prone long reads. Genome Biol [Internet]. Genome Biology; 2019;20:1–17.
- 15 6. Desk R. Review article. Toenail onychomycosis: an important global disease burden.
16 [Internet]. N. Engl. J. Med. 2005. p. 1–4.
- 17 7. University of Washington PacBio Sequencing Services. Available from:
18 <https://pacbio.gs.washington.edu/>
- 19 8. UC Davis Genome Center. Available from:
20 <https://dnatech.genomecenter.ucdavis.edu/prices/>
- 21 9. Nanopore tech. Available from: <https://nanoporetech.com/products/comparison/>
- 22 10. Cho YS, Kim H, Kim HM, Jho S, Jun JH, Lee YJ, et al. Corrigendum: An ethnically
23 relevant consensus Korean reference genome is a step towards personal reference
24 genomes. Nat Commun [Internet]. 2017;8:16168.

- 1 11. Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina
2 sequence data. *Bioinformatics* [Internet]. 2014;30:2114–20
- 3 12. Porechop, adapter trimmer for Oxford Nanopore reads.
4 <https://github.com/rrwick/Porechop>
- 5 13. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *bioRxiv* [Internet].
6 2019;530972.
- 7 14. Racon, ultrafast consensus module for raw de novo genome assembly of long
8 uncorrected reads. <https://github.com/isovic/racon>
- 9 15. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An
10 integrated tool for comprehensive microbial variant detection and genome assembly
11 improvement. *PLoS One* [Internet]. 2014;9.
- 12 16. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM. BUSCO:
13 Assessing genome assembly and annotation completeness with single-copy orthologs.
14 *Bioinformatics* [Internet]. 2015;31:3210–2.
- 15 17. Ghurye J, Pop M, Koren S, Bickhart D, Chin CS. Scaffolding of long read assemblies
16 using long range contact information. *BMC Genomics* [Internet]. *BMC Genomics*;
17 2017;18:1–11.
- 18 18. BBMap, short read aligner and other bioinformatics tools.
19 <https://sourceforge.net/projects/bbmap/>
- 20 19. Pomoxis, bioinformatics tools for nanopore research.
21 <https://nanoporetech.github.io/pomoxis/>

22

1 **Table 1. Statistics of raw sequenced reads**

	ONT PromethION R9.4.1		PacBio Sequel		Short read
	27×	64×	30×	62×	Illumina HiSeq 2000
Number of reads	15,004,723	47,591,997	11,195,434	20,683,965	1,433,779,680
Total length of reads	80,770,821,288	193,027,803,978	92,229,416,062	187,914,740,184	144,811,747,680
N50	12,736	9,190	13,426	14,568	101
Max contig length	774,322	1,160,324	65,865	169,910	101

2

3

4

5

6

7

8

9

10

11

12

1 **Table 2.** Statistics of KOREF genome assemblies using ONT PromethION and PacBio Sequel sequencing

	ONT PromethION R9.4.1				PacBio Sequel			
	27× assembly*	64× assembly**	27× assembly with Hi-C	64× assembly with Hi-C	30× assembly	62× assembly	30× assembly with Hi-C	62× assembly with Hi-C
Contigs / Scaffolds No.	3,262	3,725	2,313	3,179	2,443	2,695	1,476	2,139
Total length	2,757,297,803	2,827,624,042	2,757,776,303	2,827,900,542	2,800,962,512	2,815,311,932	2,801,450,512	2,815,594,432
Scaffold N50	7,655,153	16,706,773	32,758,624	56,457,651	11,137,362	17,931,968	38,113,117	59,361,327
Max contig / scaffold length	60,569,695	88,903,341	120,666,262	175,227,974	50,101,007	77,816,513	126,818,544	174,360,016
Gap	0.00%	0.00%	0.02%	0.01%	0.00%	0.00%	0.02%	0.01%
GC content	40.82%	40.81%	40.82%	40.81%	40.90%	40.92%	40.90%	40.90%

2

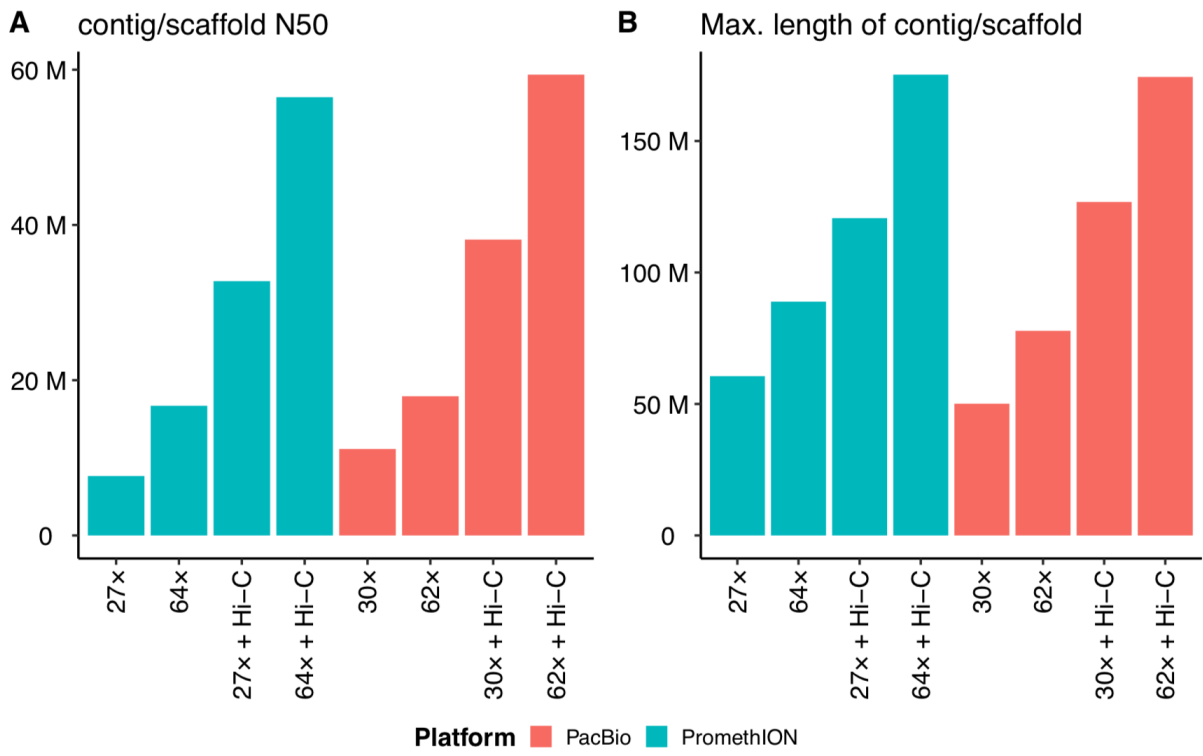
1 **Table 3.** Statistics of KOREF genome assembly assessment using BUSCO and accuracy comparison

BUSCO assessment	ONT PromethION R9.4.1				PacBio Sequel			
	27× assembly	64× assembly	27× assembly with Hi-C	64× assembly with Hi-C	30× assembly	62× assembly	30×assembly with Hi-C	62× assembly with Hi-C
Complete	92.5%	92.7%	92.6%	94.0%	93.8%	93.9%	93.8%	93.5%
Complete and single-copy	91.8%	91.6%	91.9%	93.2%	93.0%	93.1%	93.0%	92.7%
Complete and duplicated	0.7%	1.1%	0.7%	0.8%	0.8%	0.8%	0.8%	0.8%
Fragmented	3.1%	3.7%	3.2%	3.1%	3.0%	2.9%	3.0%	3.1%
Missing	4.4%	3.6%	4.2%	2.9%	3.2%	3.2%	3.2%	3.4%
Accuracy comparison*	99.78%	99.73%	99.78%	99.73%	99.83%	99.79%	99.86%	99.80%

* Compared with KOREF_S, the single assembly of KOREF

2

1 **Figure 1**



2

3