

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,
8 Republic 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, PhD

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

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
5 options. However, PacBio's SMRT sequencing is expensive for a full human genome
6 assembly and costs over 40,000 USD for 30x coverage as of 2019. ONT PromethION
7 sequencing, on the other hand, is one-twelfth the price of PacBio for the same coverage. This
8 study aimed to compare the cost-effectiveness of ONT PromethION and PacBio's SMRT
9 sequencing in 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 64x
13 coverage (193 Gb, 3 flowcell sequencing) resulted in 3,725 contigs with N50s of 16.7 Mbp
14 and a total genome length of 2.8 Gbp. It was comparable to a KOREF assembly constructed
15 using PacBio at 62x coverage (188 Gbp, 2,695 contigs and N50s of 17.9 Mbp). When we
16 applied Hi-C-derived long-range mapping data, an even higher quality assembly for the 64x
17 coverage was 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
4 variation discovery using single molecule sequencing methods. Oxford Nanopore
5 Technologies (ONT) and PacBio platforms have advantages of a short run time and long read
6 lengths over short fragmented reads by Illumina [1, 2]. Unfortunately, both methods share
7 high base-calling error rates [3, 4]. However, bioinformatics pipelines for self-error
8 correction and/or polishing sequences with short reads have become an effective option, and
9 the overall accuracy of long read based assemblies is approaching what is required to be a
10 viable option for personal reference genome construction [5]. Despite its excellent
11 performance, PacBio's SMRT sequencing is expensive for the effective coverage required for
12 a full human genome assembly, costing over 40,000 USD for 30x coverage (with 15 SMRT
13 cells; from an estimated 6 Gbp raw reads production per SMRT cell) as of 2019 [6, 7, 8]. On
14 the other hand, the nanopore based single molecule, long read platform, PromethION from
15 ONT is highly cost-effective at one-twelfth the price of PacBio's for the same read amount,
16 with an advantage of even longer average and maximum read lengths [9]. Although the two
17 methods share some similarity, they are fundamentally different in that ONT uses a minimal
18 amount of reagents with small form factor devices, and can be a promising future technology
19 for a very broad scope of applications given its advantageous size and cost.

20 In this study, we performed benchmark tests of PromethION and PacBio with low
21 and high coverages of sequencing data and investigated the advantages of pairing these long
22 read 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 x 10⁶) 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
17 min at room temperature, and washed twice with 200 µL fresh 70% ethanol. The pellet was
18 allowed to dry for 30 s, and then DNA was eluted in 61 µL of nuclease-free water. An aliquot
19 of 1 µL 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-
21 LSK109 Ligation Sequencing Kit 1D, Oxford Nanopore Technologies (ONT)), 25 µL
22 Ligation Buffer (LNB, SQK-LSK109), and 10 µL NEBNext Quick T4 DNA Ligase (NEB,
23 E6056) to 60 µL bead cleaned-up DNA, followed by gentle mixing and incubation for 10 min
24 at room 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
6 new Eppendorf Lobind tube. A 1- μ L aliquot was quantified by Qubit to ensure \geq 500 ng
7 DNA was retained.

8 To load the library, 75 μ L of Sequencing Buffer (SQB, SQK-LSK109) was mixed
9 with 51 μ L of Loading Beads (LB, SQK-LSK109) and this mixture was added to 24 μ L DNA
10 library. This library was mixed by pipetting slowly and 150 μ L of sample was loaded through
11 the inlet 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
19 annealed to the SMRTbell template, DNA polymerase was bound to the complex (Sequel
20 Binding kit 2.0). We purified the complex using AMPure Purification to remove excess
21 primer and polymerase prior to sequencing. The SMRTbell library was sequenced using
22 SMRT cells (Pacific Biosciences) using Sequel Sequencing Kit v2.1 and 10 h movies were
23 captured for each SMRT Cell 1M v2 using the Sequel (Pacific Biosciences) sequencing

1 platform.

2

3 **Short read sequencing by Illumina HiSeq**

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

7

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

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

22

1 **Short and long sequence reads processing**

2 A total of 144 Gbp of short paired-end DNA raw reads were obtained from SRA2204706.
3 Adapter sequences were trimmed from sequenced raw reads using Trimmomatic v0.36 [11]
4 (ILLUMINACLIP:2:30:10 LEADING:5 TRAILING:5 SLIDINGWINDOW:4:20
5 HEADCROP:15 MINLEN:60), and screening for vectors and microbial contaminants were
6 performed using customized database from Refseq. After preprocessing, a total of 137 Gbp
7 cleaned reads were obtained.

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

13

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

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

1 locate the presence and absence of 4,104 single copy orthologous genes from mammalian
2 OrthoDB v9.

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

8 Long read assemblies from 27x and 64x PromethION sequencing yielded total
9 assembly sizes of 2,757 Mbp and 2,827 Mbp, with scaffold N50s of 7.6 Mbp and 16.7 Mbp,
10 respectively (Table 2). Assemblies from PacBio sequencing at 30x and 60x coverage yielded
11 the total assembly sizes of 2,800 Mbp and 2,815 Mbp, with scaffold N50s of 11.1 Mbp and
12 17.9 Mbp, respectively. Adding Hi-C reads to assemblies led to 3.4- to 4.3-fold increase in
13 the scaffold N50 lengths of PromethION (32.7 Mbp for 27x coverage and 56.4 Mbp for 64x
14 coverage). For the PacBio assemblies, 2.2- to 3.4-fold increase was achieved for the scaffold
15 N50 lengths (38.1 Mbp for 30x coverage and 40.9 Mbp for 62x coverage). The longest
16 scaffold from both PromethION and PacBio assemblies with Hi-C was two times the length
17 of the assemblies without Hi-C.

18

19 **Comparison between PromethION and PacBio assemblies**

20 The comparison between PromethION and PacBio assemblies without Hi-C mapping
21 information using sequenced reads at 64x coverage showed comparable quality. In terms of
22 N50, the PromethION assembly at 64x coverage yielded 1.45-fold and 0.93-fold longer N50s
23 compared with the PacBio assemblies at 30x and 62x coverage, respectively (Figure 1a).

1 When we compared the longest contigs, the PromethION assembly at 64x coverage yielded
2 1.7-fold and 1.1-fold length increase compared with the PacBio assemblies at 30x and 62x
3 coverage, respectively (Figure 1b). Comparing the number of scaffolds, PacBio assembly at
4 30x coverage showed the fewest (2,443) compared with that of PromethION assembly at 64x
5 coverage (3,725).

6 When Hi-C mapping information was added to the assembly construction, the
7 PromethION assembly at 64x coverage showed the best statistics as N50s of 56.4 Mbp and
8 the longest scaffold length of 175.2 Mbp. The PromethION assembly at 27x coverage with
9 Hi-C mapping information yielded 32.7 Mbp for N50s, which was comparable to both 30x
10 and 62x coverage PacBio assemblies with Hi-C; 0.85-fold and 0.79-fold for N50s,
11 respectively (Table 2).

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

18

19 **Conclusions**

20 We generated high-quality assemblies of the Korean reference genome, KOREF, using
21 ONT's PromethION long-reads accompanied with Hi-C mapping information and compared
22 them against PacBio sequencing and assemblies of the same sample. Comparing the results
23 from the PromethION 64x sequencing to the PacBio 62x sequencing, we found that the

1 former provided high contiguity and completeness at one-twelfth the cost of PacBio. Results
2 from just 27x PromethION sequencing combined with Hi-C mapping information were also
3 comparable to the 30x and even 62x coverage PacBio sequencing data. Therefore, to generate
4 a chromosome-scale assembly with a long-read technology, at present, the ONT's
5 PromethION sequencing is a good alternative to PacBio's, owing to its quality and cost-
6 effectiveness. Simple pore-based long read sequencing has potential to dramatically improve
7 sequencing and subsequent bioinformatics analysis for personal genome projects and cancer
8 genome analyses where *de novo* assemblies are necessary for structural and copy number
9 variations that cannot be detected easily by conventional short read only methods.

10

11 **Availability of supporting data**

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

15

16 **Abbreviations**

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

19

20 **Competing interests**

21

1 Y.S.C. is an employee, and J.B. is the CEO of Clinomics Inc. J.B. and Y.S.C. have an equity
2 interest in the company. All other coauthors have no conflicts of interest to declare.

3

4 **Funding**

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

8

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.
8 Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing
9 data. Nat 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. BMAP, 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
	27x	64x	30x	62x	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			
	27x assembly	64x assembly	27x assembly with Hi-C	64x assembly with Hi-C	30x assembly	62x assembly	30x assembly with Hi-C	62x assembly with Hi-C
Contigs / Scaffolds No.	3,262	3,725	2,313	3,179	2,443	2,695	1,476	2,338
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	2,818,181,112
Scaffold N50	7,655,153	16,706,773	32,758,624	56,457,651	11,137,362	17,931,968	38,113,117	40,953,073
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	131,630,574
Gap	0.00%	0.00%	0.02%	0.01%	0.00%	0.00%	0.02%	0.02%
GC content	40.82%	40.81%	40.82%	40.81%	40.90%	40.92%	40.90%	40.92%

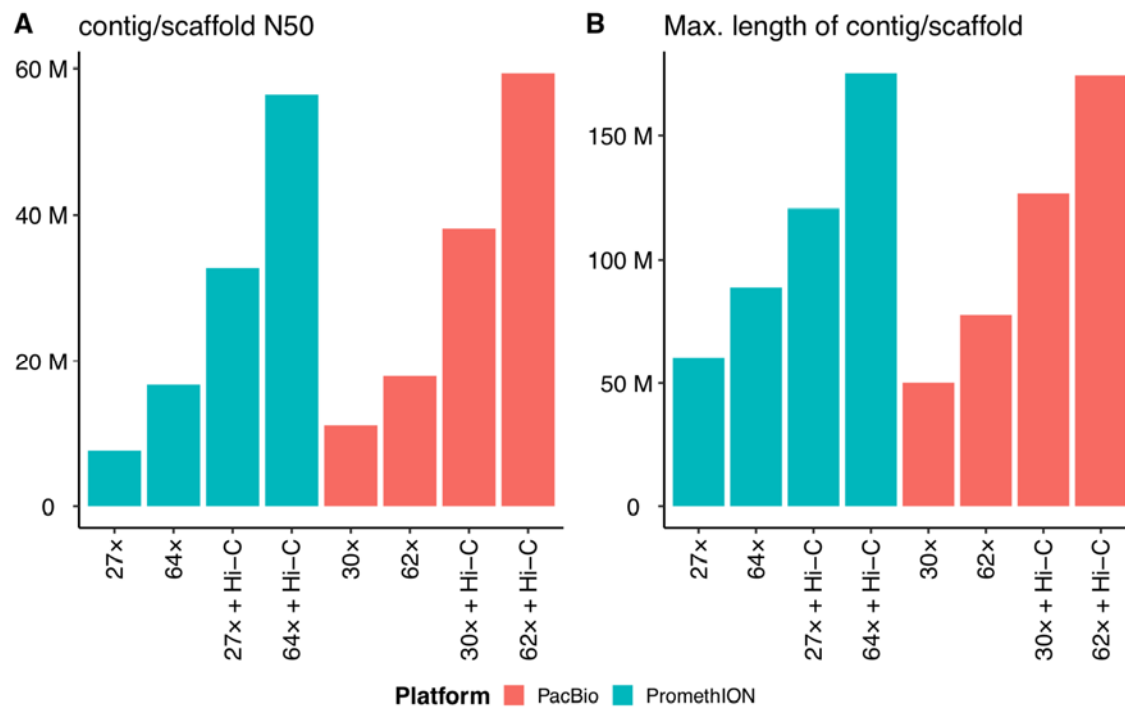
2
3
4
5
6
7
8

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

BUSCO assessment	ONT PromethION R9.4.1				PacBio Sequel			
	27x assembly	64x assembly	27x assembly with Hi-C	64x assembly with Hi-C	30x assembly	62x assembly	30x assembly with Hi-C	62x assembly with Hi-C
Complete	92.50%	92.70%	92.60%	94.00%	93.80%	93.10%	93.80%	93.50%
Complete and single-copy	91.80%	91.60%	91.90%	93.20%	93.00%	92.00%	93.00%	92.70%
Complete and duplicated	0.70%	1.10%	0.70%	0.80%	0.80%	1.10%	0.80%	0.80%
Fragmented	3.10%	3.70%	3.20%	3.10%	3.00%	3.20%	3.00%	3.10%
Missing	4.40%	3.60%	4.20%	2.90%	3.20%	3.70%	3.20%	3.40%
Accuracy comparison*	99.81%	99.79%	99.78%	99.73%	99.86%	99.79%	99.84%	99.80%

* Compared with KOREF_S, the single assembly of KOREF

1 **Figure 1**



2

3