# GigaScience

## Comparative analysis of seven short-read sequencing platforms using the Korean Reference Genome: MGI and Illumina sequencing benchmark for whole-genome sequencing
### --Manuscript Draft--

| Manuscript Number: | GIGA-D-20-00072R1 |
|---|---|
| Full Title: | Comparative analysis of seven short-read sequencing platforms using the Korean Reference Genome: MGI and Illumina sequencing benchmark for whole-genome sequencing |
| Article Type: | Data Note |
| Funding Information: | Ulsan National Institute of Science and Technology (1.200047.01) — Dr. Jong Hwa Bhak |

| Abstract: | Background:  MGISEQ-T7 is a new whole-genome sequencer developed by Complete Genomics and MGI utilizing DNA nanoball and combinatorial probe anchor synthesis technologies to generate short reads at a very large scale – up to 60 human genomes per day. However, it has not been objectively and systematically compared against Illumina short-read sequencers.<br>Findings:  By using the same KOREF sample, the Korean Reference Genome, we have compared seven sequencing platforms including BGISEQ-500, MGISEQ-T7, HiSeq2000, HiSeq2500, HiSeq4000, HiSeqX10, and NovaSeq6000. We measured sequencing quality by comparing sequencing statistics (base quality, duplication rate, and random error rate), mapping statistics (mapping rate, depth distribution, and %GC coverage), and variant statistics (transition/transversion ratio, dbSNP annotation rate, and concordance rate with SNP genotyping chip) across the seven sequencing platforms. We found that MGI platforms showed a higher concordance rate for SNP genotyping than HiSeq2000 and HiSeq4000. The similarity matrix of variant calls confirmed that the two MGI platforms have the most similar characteristics to the HiSeq2500 platform.<br>Conclusions:  Overall, MGI and Illumina sequencing platforms showed comparable levels of sequencing quality, uniformity of coverage, %GC coverage, and variant accuracy, thus we conclude that the MGI platforms can be used for a wide range of genomics research fields at a lower cost than the Illumina platforms. |
|---|---|

| Corresponding Author: | Jong Hwa Bhak, Ph.D.<br>UNIST<br>Ulsan, Ulsan KOREA, REPUBLIC OF |
|---|---|
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | UNIST |
| Corresponding Author's Secondary Institution: | |
| First Author: | Hak-Min Kim, Ph.D. |
| First Author Secondary Information: | |
| Order of Authors: | Hak-Min Kim, Ph.D. |
| | Sungwon Jeon |
| | Oksung Chung |
| | Je Hoon Jun |
| | Hui-Su Kim, Ph.D. |
| | Asta Blazyte |
| | Hwang-Yeol Lee |

| | Youngseok Yu |
| :--- | :--- |
| | Yun Sung Cho, Ph.D. |
| | Dan M. Bolser, Ph.D. |
| | Jong Hwa Bhak, Ph.D. |
| **Order of Authors Secondary Information:** | |
| **Response to Reviewers:** | Reviewer reports:<br><br>Reviewer #1: In this manuscript, Kim et al. compared seven sequencing platforms, including 2 MGI platforms (BGISEQ-500 and MGISEQ-T7) and 5 Illumina platforms (HiSeq2000, HiSeq2500, HiSeq4000, HiSeqX10, and NovaSeq6000), by using one human genome. The sequencing quality of different sequencing platform was assessed by basic sequencing statistics, mapping statistic and variant statistic. Overall the manuscript is suitable to be published on Giga Science after a major revision. There are several major issues with the work presented in the manuscript, as listed below:<br>=> Thank you for precise and critical feedback. We have modified the text and added further analysis to accommodate the reviewer's suggestions. (See below for point-by-point responses).<br><br>1. This work only contains samples from one human individual. It's really hard to reach a confident conclusion based on such a small sample size.<br>=> It is a generally correct point. However, both platforms produce massive amounts of sequences and the sample number would not affect the conclusion much as our study rely on how the two sets of platforms are similar or dissimilar in terms of variant calling.<br><br>This work still needs more samples and even replicates (both Cross-platform replicates and intra-platform replicates) to do further analysis, and provide confident evidence.<br>=> We think this is a practically important point. Unfortunately, we have not generated replicates for each sequencer. First, this study is based on years of sequencing history with one reference sample and each sequencing batch can contain multiple replicates or not. It is because each platform has a different amount of sequence output per run, it is impossible to produce a controlled amount of sequences in a certain common replicate number. We stated these limitations in the discussion part of the manuscript. The purpose of this benchmarking work was to compare two major platforms (MGI and Illumina).<br><br>2. The samples for sequencing were extracted on different points of time from the individual, that we wonder if the differences between mutation sets of seven sequencing platforms were caused by different sampling time and the bias of sampling process.<br>=> There must be some problems caused by the different sampling time and the sampling process mentioned by the reviewer. We used a Korean male sample and the difference between the first and the last sampling time is about 7 years. It is known that the human germline mutation rate is approximately $0.5 \times 10^{-9}$ per base pair per year (Scally A, 2016. [10.1016/j.gde.2016.07.008]), which means that 10.5 germline mutations can be accumulated in 7 years. In this respect, although the mutation rate of DNA of leukocyte, a somatic cell, is expected to be higher than that of a germline cell, the number of mutations accumulated over the 7 years would be much lower than the difference between platforms. Therefore, we think that the different sampling time had no significant effect on the results.<br>For the case of sampling process bias, we stated in the discussion part of the manuscript that there is a clear limitation in the sampling process. Although there are some limitations as the reviewer mentioned, we think our study is still meaningful in that it provides the data generated by the short read-based whole genome sequencing platform, which is the most used in the field. We compared the long existing common Illumina platforms with the relatively new MGISEQ-T7 platform using one human whole genome sequence (WGS) data which has not been done before.<br><br>3. This manuscript needs to show more detail about the sequencing process, such as the number of the flow cell and sequencing cycle, the run time of the sequencing process, the amount of DNA each sequencing platform needs. |

=> We added the detailed methods for DNA extraction, library preparation, and sequencing process in the Materials and Methods section.

4. In order to compare, the sequencing data of seven sequencing platforms need to have the same genome coverage.
=> Very good point. As pointed out by the reviewer, we set the same genome coverage of the seven platforms and updated all subsequent analyses after analyzing the whole data. Please see Figure S5 and Table S4.

5. The results of the manuscript let me worry about the quality of the sequencing data generated from Hiseq2000 and Hiseq4000. More samples or replicates were needed to prove these results that the author found were normal.
=> HiSeq2000 and HiSeq4000 platforms are old, and their quality is not good compared to other platforms in our case. Currently it is not possible to have more replicates as these machines are often not available in sequencing centers and, also, it is quite expensive to run them now. Still, to compare with MGI platforms, we decided to add as many Illumina platforms as possible.

6. According to the official information, MGI platforms have low duplicate rate than any sequencing platform which needs PCR. But this work showed MGISEQ T7 had highest duplicate rate, I suggest the authors prove their finding by using other samples or individuals.
=> The official information showed a duplicate rate of less than 3% when using a PCR free library kit. However, we used the FS library kit that included the PCR process. Therefore, it seems that the duplicate rate is higher than the manufacturer's official information. We provide the table presenting the mapping rates and duplicate rates of other human samples produced simultaneously with the KOREF sample. We found that the duplicate rates of the other human samples that were sequenced simultaneously with the KOREF sample were also high (see link below).

https://github.com/howmany2/SequencingPlatformComparison/raw/master/Mapping%20and%20duplicate%20rate%20of%20samples%20using%20PE100%20protocol%20and%20MGISEQ-T7.xlsx

An FS library kit containing PCR steps was used for MGISEQ-T7 sequencing of the KOREF sample. Furthermore, according to the sequencing vendor, the PE100 (Paired-end 100 bp) protocol has a high duplication rate, and the new PE150 (Paired-end 150 bp) protocol has a duplication rate less than 3%. We used the PE100 protocol for the KOREF sample and it can be a reason for why relatively many duplicated reads were found from the reads generated by the MGISEQ-T7 platform. However, we think the duplicate rate does not affect variant results much because it was analyzed after removing duplicate reads and matching to the same genome coverage for the seven sequencing platforms.

7. The methods for identifying the platform-specific covered region are unreasonable as different sequencing platforms had different coverage.
=> We agree with the reviewer's comment. We set the same genome coverage of the seven platforms and updated the result. As a result, the number of platform-specific covered regions of MGI platform decreased from 1,516 to 1,436, and in the case of Illumina, increased from 2,264 to 2,881. However, it was confirmed that the %GC ratio of the platform-specific covered region is the same as before meaning that the MGI platform covers a higher GC area (see Figure S10).

8. The Comparison of variants detected among seven platforms needs further analysis. Authors need a standard SNP and indel list of the Korean reference genome, which is verified by Sanger sequencing or other methods, to replace the dbSNP and SNP genotype chip as a compare object. What the relationship of FP, FN and the sequencing errors?
=> We agree with the reviewer's comment that it is a powerful tool to compare the variants to the gold standard variant set. However, to our knowledge, there is no gold standard variant set for the KOREF, which can give FP, FN, and sequencing error information, and, for this reason, we could not make a design for this study to conduct more precise and accurate comparison among the NGS platforms. As an alternative, we examined how much difference exists among the sequences generated by different

NGS platforms which are generally used methods for genome sequencing.

9. The introduction of this manuscript is too simple.
=> We added several sequencing platform comparative studies to the introduction section.

Minor revisions:
1. The coverages of BGISEQ-500 and HiseqX10 were not mentioned in the first section.
=> We added the coverages of BGISEQ-500 and HiSeqX10 in the first section.

2. Using the ratio of singletons may help you to bring out your findings more clearly.
=> We agree with the reviewer's comment. We examined the concordance rate of the singleton variants with SNP genotyping data to determine the accuracy of the singleton variants (see link below). However, it was difficult to obtain statistically significant results because there were very few overlapping positions between the singleton variants and the SNP chip data.

https://github.com/howmany2/SequencingPlatformComparison/raw/master/Compariso n%20between%20singleton%20variant%20and%20SNP%20genotyping%20chip.xlsx

Reviewer #2: The submitted study has characterized sequencing quality, uniformity of coverage, %GC coverage, and variant accuracy of seven sequencing platforms. They found that MGI platforms showed a higher concordance rate of SNP genotyping than HiSeq series. The study is of interest to genomics and sequencing technologies areas. Two concerns must be addressed prior to acceptance.
=> Thank you for the feedback. We have modified the text and added further analysis to accommodate the reviewer's suggestion. (See below point-by-point responses).

1)The author defined low-quality reads as those that had more than 30% of bases with a sequencing quality score lower than 20. I am wondering whether the results is stable once the definition changed?
=> As a supplementary analysis, we conducted an analysis without the filtering step to see how much the read filtering step affects in the result of this study. The supplementary analysis was conducted by matching the number of unfiltered reads with that of clean reads of prior analysis. The two tables below are the results of comparing the read mapping and variant statistics between the cases using clean (filtered) and unfiltered sequences (see link below).

https://github.com/howmany2/SequencingPlatformComparison/raw/master/Mapping%2 0rate%20and%20Variant%20statistics%20between%20clean%20reads%20and%20un filtered%20reads.xlsx

As a result of using the unfiltered sequences, there was no notable difference in mapping and duplicate rates. The number of SNVs increased by 0.8% on average, and as the number of heterozygous SNVs increased, the hetero/homo ratio increased by 0.02 on average. Interestingly, the differences in total SNVs between clean and unfiltered reads in the two MGI platforms were less than that of the Illumina platforms. In the case of the Illumina platforms, on average, 44,000 additional SNVs were discovered when unfiltered reads were used compared to the case of the clean reads, while the increment in MGI platform was 800 SNVs on average when using unfiltered reads.

2) It looks the author ignored a highest duplicate ratio was found in MGISEQ-T7. More discussion and analysis should be performed to make this clear. The author claimed that duplicates and adapter contamination may be more affected by the process of sample preparation than by the sequencing instrument. However, again, no evidence was provided.
=> We agree with the reviewer's concerns about the high duplicate ratio. We provide the table presenting the mapping rates and duplicate rates of other human samples produced simultaneously with the KOREF sample. We found that the duplicate rates of

other human samples that were sequenced simultaneously with the KOREF sample were also high (see Table below).

An FS library kit containing PCR steps was used for MGISEQ-T7 sequencing of the KOREF sample. Furthermore, according to the sequencing vendor, the PE100 (Paired-end 100 bp) protocol has a high duplication rate, and the new PE150 (Paired-end 150 bp) protocol reduces the duplication rate to less than 3%. We used the PE100 protocol for the KOREF sample sequencing and it can be a reason why relatively many duplicated reads were found from the reads generated by the MGISEQ-T7 platform. However, we think the duplicate rate does not affect variant calling results because it was analyzed after removing the duplicate reads and matching to the same genome coverage for the seven sequencing platforms (see link below).

https://github.com/howmany2/SequencingPlatformComparison/raw/master/Mapping%20and%20duplicate%20rate%20of%20samples%20using%20PE100%20protocol%20and%20MGISEQ-T7.xlsx

There are three main causes of duplicate reads generated by NGS technology.
1. Natural duplication
2. PCR duplicates (occur in library preparation step)
3. Optical duplicates (occur in sequencing step)
Natural duplications are not discussed in this section because it is difficult to distinguish them from PCR duplicates and optical duplicates. The following table showed the ratio of PCR duplication and optical duplication of the seven platforms (see link below).

https://github.com/howmany2/SequencingPlatformComparison/raw/master/Statistics%20of%20PCR%20duplicate%20and%20optical%20duplicate%20in%20seven%20sequencing%20platforms.xlsx

This result showed that PCR duplication occurs at least 2 times more than the optical duplication. (Unfortunately, the two MGI platforms were unable to calculate optical duplication.) This means that most duplication occurs during the library preparation rather than the sequencing steps.

The adapter contamination is caused by the sequencing of short DNA fragments that are shorter than the read length (Turner FS, 2014. 10.3389/fgene.2014.00005). For this reason, it can be expected that adapter contamination is mainly affected by the library preparation step, because size selection of DNA fragments is a part of the library preparation step; improper operation of size selection can introduce the shorter DNA fragments into the DNA library for sequencing.

Reviewer #3: The authors compare various short-insert, short-read whole-genome sequencing platforms used by academic researchers and clinical scientists.

My minor comments and suggestions are:

• As stated by the authors, Illumina platforms are indeed now considered 'historical.' However, many Illumina sequencers are still heavily used - in particular in pathology labs. This manuscript may prove very useful when arguing for an instrument upgrade in such a setting.

• You may like to comment on single tube long fragment read (stLFR), which enables the sequencing of long transcripts by sequencing bar-coded reads on the BGISEQ-500 platform [and, thus, probably also MGISEQ-T7) (10.1101/gr.245126.118). This technology is relatively cheap and is likely to decrease in cost - another argument for the adaption of MGI platforms in the laboratory.

• You may want to comment on Illumina library kits. It is possible that revisions [in the five-six years since the data in your study were generated] to these kits could improve the sequencing results (e.g., see 10.1371/journal.pone.0113501). I realize the effect may be minor, but it may nevertheless be useful to remind the reader about the potential for *slightly* better raw read statistics.
=> Thank you for your positive feedback and the suggestions. We added the idea suggested in your comments to the discussion part of the manuscript. (See Discussion

| | section lines 209-210) |
|---|---|
| **Additional Information:** | |
| **Question** | **Response** |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>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.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>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.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials**<br><br>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. | Yes |

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

1 **Comparative analysis of seven short-read sequencing platforms**

2 **using the Korean Reference Genome: MGI and Illumina**

3 **sequencing benchmark for whole-genome sequencing**

4 **Hak-Min Kim[1], Sungwon Jeon[2,3], Oksung Chung[1], Je Hoon Jun[1], Hui-Su Kim[2], Asta**

5 **Blazyte[2,3], Hwang-Yeol Lee[1], Youngseok Yu[1], Yun Sung Cho[1], Dan M. Bolser[4]\*, and Jong**

6 **Bhak[1,2,3,4]\***

7 [1]Clinomics, Ulsan National Institute of Science and Technology, Ulsan 44919, Republic of

8 Korea

9 [2]Korean Genomics Center (KOGIC), Ulsan National Institute of Science and Technology,

10 Ulsan 44919, Republic of Korea

11 [3]Department of Biomedical Engineering, School of Life Sciences, Ulsan National Institute of

12 Science and Technology, Ulsan 44919, Republic of Korea

13 [4]Geromics Ltd, 23 King Street, Cambridge, CB1 1AH, UK

14 \*Corresponding authors

15

16 Email address:

17 H.M.K.: howmany2@gmail.com

18 S.J.: jsw0061@gmail.com

19 O.C.: okokookk219@gmail.com

20    J.H.J.: junjh0701@gmail.com

21    H.S.K.: hskim3824@gmail.com

22    A.B.: blazyte.asta@gmail.com

23    H.Y.L.: hyeol911@gmail.com

24    Y.Y.: yung7449@gmail.com

25    Y.S.C.: joys0406@gmail.com

26    D.M.B.: dan@geromics.co.uk

27    J.B.: jongbhak@genomics.org

28

29

## Abstract

**Background:** MGISEQ-T7 is a new whole-genome sequencer developed by Complete Genomics and MGI utilizing DNA nanoball and combinatorial probe anchor synthesis technologies to generate short reads at a very large scale – up to 60 human genomes per day. However, it has not been objectively and systematically compared against Illumina short-read sequencers. **Findings:** By using the same KOREF sample, the Korean Reference Genome, we have compared seven sequencing platforms including BGISEQ-500, MGISEQ-T7, HiSeq2000, HiSeq2500, HiSeq4000, HiSeqX10, and NovaSeq6000. We measured sequencing quality by comparing sequencing statistics (base quality, duplication rate, and random error rate), mapping statistics (mapping rate, depth distribution, and %GC coverage), and variant statistics

40 (transition/transversion ratio, dbSNP annotation rate, and concordance rate with SNP

41 genotyping chip) across the seven sequencing platforms. We found that MGI platforms showed

42 a higher concordance rate for SNP genotyping than HiSeq2000 and HiSeq4000. The similarity

43 matrix of variant calls confirmed that the two MGI platforms have the most similar

44 characteristics to the HiSeq2500 platform. **Conclusions:** Overall, MGI and Illumina

45 sequencing platforms showed comparable levels of sequencing quality, uniformity of

46 coverage, %GC coverage, and variant accuracy, thus we conclude that the MGI platforms can

47 be used for a wide range of genomics research fields at a lower cost than the Illumina platforms.

48 *Keywords*: MGISEQ-T7; whole-genome sequencing; sequencing platform comparison;

49

# Introduction

51 Recently, due to the rapid technological advancement, the second- and third-generation

52 sequencing platforms can produce a large amount of short- or long-read data at relatively low

53 cost [1]. Depending on the application, these sequencers offer several distinct advantages.

54 Short-read based second-generation sequencing can be used to efficiently and accurately

55 identify genomic variations. Long-read based third-generation sequencing can be used to

56 identify structural variations and build high quality *de novo* genome assemblies [2]. Short-read

57 sequencing technologies are routinely used in large-scale population analyses and molecular

58 diagnostic applications because of the low cost and high accuracy [3]. The recent platforms

59 from Illumina are the HiSeqX10 and NovaSeq6000 short-read sequencers. A competing

60 sequencer developed by Complete Genomics and MGI Tech is the MGISEQ-T7 (also known

61 as DNBSEQ-T7). MGISEQ-T7 is a new sequencing platform after BGISEQ-500 that uses

62 DNA nanoball and combinatorial probe anchor synthesis to generate short reads at a very large

63    scale [4].

64        Recently, a paper was published showing similar accuracy of SNP detection for the

65    BGISEQ-500 platform compared to the HiSeq2500 [5]. The quality of the data generated by

66    BGISEQ-500 was shown to be of high quality. However, some of its characteristics showed

67    lower quality compared to Illumina HiSeq2500. In addition, the comparison results for DNA,

68    RNA, and metagenome sequencing of the Illumina and the MGI platforms have been reported

69    [6-8]. Also, coronavirus analysis studies using an MGI platform have been reported in 2020 [9,

70    10]. Still, no study has compared Illumina platforms with MGISEQ-T7 for whole-genome

71    sequencing (WGS). In the present study, we compared seven short-read based sequencers; two

72    MGI platforms (BGISEQ-500 and MGISEQ-T7) and five Illumina platforms (HiSeq2000,

73    HiSeq2500, HiSeq4000, HiSeqX10, and NovaSeq6000) (Table 1). We focused on how similar

74    the two sets of platforms are rather than the accuracy of each sequencer, by comparing variants,

75    platform-specific covered regions as well as the concordance rate to SNP genotyping chip.

76

77    **Results**

78    **Sequencing data summary**

79    We analyzed and benchmarked the whole-genome sequencing data quality generated by the

80    seven sequencers using the KOREF (the Korean Reference Genome) [11] DNA. Due to the

81    sequential release and distribution of the sequencers, KOREF sequencing has been carried out

82    in nine years since 2010. Therefore, the blood samples, library construction, and sequencing

83    conditions were not the same, although all the sample were from one individual. The Illumina

84    platform data used here were from 2012 to 2019, while the MGI platform data were from 2017

85    and 2019. Also, the read length differs depending on the platform. The Illumina HiSeq2000

86    had the shortest read length of 90 bp paired-end (PE) and the HiSeq4000, HiSeqX10, and

87    NovaSeq6000 had 151 bp PE. The read length of the HiSeq2500 is 101 bp PE and that of the

88    BGISEQ-500 and MGISEQ-T7 is 100 bp PE. Also, there is a difference in the amount of data

89    as well. Thus, we randomly selected 35× coverage sequencing data for HiSeq2500 and

90    NovaSeq6000 which have that much sequencing data matching to that of BGISEQ-500 and

91    HiSeqX10. HiSeq2000, HiSeq4000, and MGISEQ-T7 had roughly 30× coverage.

92

93    **Assessment of base quality and sequencing error in raw reads**

94           Base quality is an important factor in evaluating the performance of sequencing

95    platforms. We analyzed the sequencing quality by identifying low-quality reads. First, we

96    investigated the base quality distribution of raw reads with the FastQC (FastQC,

97    RRID:SCR_014583) [12]. All seven sequencing platforms showed that the quality of each

98    nucleotide gradually decreased towards the end of a read (Fig. S1). The quality value of the

99    HiSeq4000 and HiSeqX10 reads showed a tendency to decrease rapidly towards the end of the

100   read. We defined low-quality reads as those that had more than 30% of bases with a sequencing

101   quality score lower than 20. The fraction of low-quality reads ranged from 2.8% to 18.3%

102   across the seven sequencing platforms (Fig. S2 and Table S1). Based on the filtering criteria,

103   the newest platforms, NovaSeq6000 and MGISEQ-T7, showed the lowest percentage of low-

104   quality reads (2.8% and 4.2%, respectively).

105          We analyzed the frequency of random sequencing errors (ambiguous base, N), which

106   is also an important factor to evaluate the quality of the sequencing platform. We found that

107   the HiSeq2000, HiSeq4000, and HiSeqX10 showed a high random error ratio in certain

108   sequencing cycles (Fig. S3 and Table S2). Furthermore, in the case of HiSeq2000, the random

109   error tended to increase gradually after each sequencing cycle. We also investigated the

110   sequencing error using *K*-mer analysis. Most erroneous *K*-mers caused by sequencing error

111   appeared at very low frequency and form a sharp left-side peak [13, 14]. Distribution of *K*-mer

112   frequencies showed similar distributions between the platforms (Fig. 1). However, there was a

113   difference in the proportion of low-frequency *K*-mer ($\leq 3$ *K*-mer depth), which was considered

114   as putative sequencing errors (Table S3). The NovaSeq6000 showed the lowest amount of

115   erroneous *K*-mer (3.91%), while the HiSeq4000 contained the highest amount of erroneous K-

116   mer (13.91%) among the seven sequencing platforms. The BGISEQ-500 and MGISEQ-T7

117   showed a moderate level of erroneous *K*-mer (7.72% and 6.39%, respectively).

118          We examined the duplication rate and adapter contamination in the seven sequencing

119   platforms (Table S2). We examined the exact duplicates, which are identical sequence copies,

120   from raw sequence data. The HiSeq2000 and MGISEQ-T7 showed the highest duplicate ratio

121   (8.71% in HiSeq2000 and 3.04% in MGISEQ-T7). The HiSeq4000, HiSeqX10 and

122   NovaSeq6000 showed higher adapter contamination rates than other platforms, probably due

123   to longer sequence length (151 bp). However, duplicates and adapter contamination may be

124   more affected by the process of sample preparation than by the sequencing instrument.

125

126   **Genome coverage and sequencing uniformity**

127          In order to assess genomic coverage and sequencing uniformity, we aligned quality-

128   filtered reads to the human reference genome (GRCh38). All seven sequencing platforms

129   showed a mapping rate of more than 99.98% and genome coverage of more than 99.6% ($\geq 1\times$;

130   Table 2). We observed a higher duplicate mapping rate in the HiSeq2000 (15.35%) and

131   MGISEQ-T7 (8.77%) than the other platforms and the same pattern as the duplication rates of

132   raw reads (see Table S2). The insert-size for paired-end libraries corresponds to the targeted

133    fragment size for each platform (Fig. S4). It has been reported that the depth of coverage is

134    often far from evenly distributed across the sequenced genome [15]. To assess the sequencing

135    uniformity, we analyzed the distribution of mapping depth for all chromosomes (Fig. S5). All

136    seven platforms showed a similar pattern of depth distribution, but interestingly, we found that

137    the depth near the centromere regions was lower exclusively in the HiSeq4000 (Figs. S6-S9).

138    We speculate that this may have been due to a bias in the library preparation step on the

139    HiSeq4000 platform.

140        In order to examine the platform-specific covered region of MGI and Illumina

141    platforms, we defined a platform-specific covered region that had significantly different depths

142    (five times difference with an average depth between MGI and Illumina platforms) based on

143    the 100 bp non-overlapping windows. Prior to examining the platform-specific covered regions,

144    mapped reads were down-sampled for all platforms to 24x coverage, which is the minimum

145    coverage among the platforms, for a fair comparison. (Table S4). We found 144 Kb and 288

146    Kb of the platform-specific covered regions from MGI and Illumina platforms, respectively

147    (Table S5). A total of 172 and 854 genes were overlapped in MGI and Illumina specific covered

148    regions, respectively, and most of them were intronic. Interestingly, however, the platform-

149    specific covered regions showed a significantly different distribution of GC ratios between the

150    MGI and Illumina platforms (Fig. S10). The MGI platforms tend to cover regions relatively

151    high in GC content (Wilcoxon rank-sum test, $P = 7.92 \times 10^{-187}$). Nevertheless, it is obvious that

152    platform-specific covered regions for Illumina platforms are slightly longer than those of the

153    MGI platforms, and these regions were not sufficiently covered by the MGI platforms.

154        Biases in PCR amplification create uneven genomic representation in classical

155    Illumina libraries [16, 17] as PCR is sensitive to extreme GC-content variation [18]. Thus, we

156    analyzed the GC biases for seven sequencing platforms. We examined the distribution of GC

157 content in sequencing reads and found that raw reads of all the seven sequencing platforms

158 showed a similar GC content distribution to the human reference genome (Fig. S11). To better

159 understand what parts of the genome were not covered properly, we generated GC-bias plots,

160 showing relative coverage at each GC level. Unbiased sequencing would not be affected by

161 GC composition, resulting in a flat line along with relative coverage = 1. We found that all the

162 seven sequencing platforms provided nearly even coverage in the moderate-GC range 20% to

163 60%, which represents approximately 95% of the human genome (Fig. 2). On the other hand,

164 the relative coverage of the HiSeq2000 platform dropped fast above 60% GC than other

165 platforms, while the NovaSeq6000 covered well above 60% GC, unlike the other platforms.

166

167 **Comparison of variants detected among seven sequencing platforms**

168 To investigate the performance of variant calling for the seven sequencing sequencers, we

169 adopted the widely used pipeline BWA-MEM (BWA, RRID:SCR_010910) [19] and GATK

170 (GATK, RRID:SCR_001876) [20-22]. We identified an average of 4.14 million single

171 nucleotide variants (SNVs), and 0.61 million indels (insertion and deletion) on each of the

172 seven sequencing platforms (Table 3). The statistics of SNVs were similar across all the seven

173 in terms of the dbSNP annotation rate (dbSNP153) and the transition/transversion (Ti/Tv) ratio,

174 which indirectly reflects SNV calling accuracy. About 3.7 million SNV loci were found on all

175 the seven sequencing platforms, and this accounts for 87% to 91% of the discovered SNVs on

176 each platform (Table S6). We found 13,999 and 9,691 platform-specific SNVs on the MGI and

177 Illumina platforms, respectively. Interestingly, the number of singletons, variations found only

178 in one platform, was higher for the Illumina (~0.10 million SNVs on average) than MGI (~0.05

179 million SNVs on average; Table S7) sequencers. This means that the difference within the

180 Illumina platforms is greater than the difference between the MGI platforms. We also analyzed

181    the number of SNVs found in any six of the seven sequencing platforms, which we considered

182    false negatives. The HiSeq2000 had the largest number of false negatives (64,856 SNVs)

183    among the seven sequencing platforms. The two MGI platforms (MGISEQ-T7 and BGISEQ-

184    500) had 18,826 and 15,657 false negatives, respectively, and those of the NovaSeq6000

185    showed the smallest number of false negatives (6,999 SNVs). To investigate the relationship

186    between the sequencing platforms, an unrooted tree was constructed using a total of 1,036,417

187    loci where the genotypes of one or more platforms differ from the rest of the platforms (Fig. 3

188    and Table S8). We found that the two MGI platforms grouped together, and they are the closest

189    to the Illumina HiSeq2500 platform. The Illumina platforms were divided into two subgroups

190    in the tree: a long read length (151 bp) group, containing the HiSeq4000, HiSeqX10, and

191    NovaSeq6000 platforms and a short read length ($\leq$101 bp) group, containing the HiSeq2000

192    and HiSeq2500 platforms. Read length primarily affects the detection of variants through

193    alignment bias and alignment errors, which are higher for short reads because there is less

194    chance of a unique alignment to the reference sequence than with longer reads [23].

195    Since it was not possible to conduct standard benchmarking procedures and determine

196    error values for each platform in this study, we compared the variations called by the seven

197    whole-genome sequences with an SNP genotyping chip as an independent platform. Of the

198    total 950,585 comparable positions, more than 99.3% of the genotypes matched the WGS-

199    based genotypes from the seven sequencing platforms (Table S9). We found that 4,356 loci in

200    the SNP genotyping were inconsistent across all seven WGS-based genotyping results,

201    suggesting that these loci are probably errors in the SNP genotyping chip. With the exception

202    of HiSeq2000 and HiSeq4000, all the other platforms showed a similar concordance rate.

203

## Discussion

Our benchmark can provide a useful but rough estimation of the quality of short-read based whole-genome sequencers. We used the same individual's samples for all seven sequencing platforms but collected at different time points in the past nine years. Just one human sample cannot justify the variation that may occur among different individuals, extracted DNA molecules, and overall sequencing qualities. Furthermore, the sequencing quality may vary much depending on the version of the library preparation kit even, on the same platform [24]. These are clear limitations in our benchmarking, however, as our purpose was to compare two major platforms, namely Illumina and MGI, still, just one person sample can function as an intuitive index for researchers who consider purchasing large sequencers to generate a very large amount of data. Our method of statistical analysis does not allow us to conclude which of the seven sequencing instruments is the most accurate and precise as there is much variation in the sample preparation and sequencer specifications. Nevertheless, overall, the data generated by the Illumina and MGI sequencing platforms showed comparable levels of quality, sequencing uniformity, %GC coverage, and concordance rate with SNP genotyping, thus it can be broadly concluded that the MGI platforms can be used for a wide range of research tasks on a par with Illumina platforms at a lower cost [7].

## Materials and Methods

**Genomic DNA extraction and SNP genotyping**

Genomic DNA used for genotyping and sequencing were extracted from the peripheral blood of a Korean male sample donor (KOREF). The genomic DNA was extracted using the DNeasy Blood & Tissue kit (Qiagen, Valencia, CA) according to the manufacturer's recommendations. DNA quality was assessed by running 1 µl on the Bioanalyzer system (Agilent) to ensure size and analysis of DNA fragments. The concentration of DNA was assessed using the dsDNA BR assay on a Qubit fluorometer (Thermo Fisher). We conducted a genotyping experiment using the Illumina Infinium Omni1 quad chip according to the manufacturer's protocols. The Institutional Review Board (IRB) at Ulsan National Institute of Science and Technology approved the study (UNISTIRB-15-19-A).

**Illumina paired-end library construction and sequencing**

High-molecular weight genomic DNA was sheared using a Covaris S2 ultra sonicator system, in order to get appropriate sizes. Libraries with short inserts of 500 bp for HiSeq2000, 400 bp for HiSeq2500 and HiSeq4000, and 450 bp for HiSeqX10 and NovaSeq6000 for paired-end reads were prepared using TruSeq DNA sample prep kit following the manufacturer's protocol. Products were quantified using the Bioanalyzer (Agilent, Santa Clara, CA, USA) and the raw data were generated by each Illumina platform. Further image analysis and base calling were conducted with the Illumina pipeline using default settings.

**MGI paired-end library construction and sequencing**

244  The KOREF genomic DNA was fragmented by Frag enzyme (MGI) to DNA fragments

245  between 100 bp and ~1,000 bp suitable for PE100 sequencing according to the manufacturer's

246  instructions (MGI FS DNA library prep set, cat no; 1000005256). The fragmented DNA was

247  further selected to be between 300 bp and ~500 bp by DNA clean beads (MGI). The selected

248  DNA fragments were then repaired to obtain a blunt end and modified at the 3'end to get a

249  dATP as a sticky end. The dTTP tailed adapter sequence was ligated to both ends of the DNA

250  fragments. The ligation product was then amplified for seven cycles and subjected to the

251  following single-strand circularization process. The PCR product was heat-denatured together

252  with a special molecule that was reverse-complemented to one special strand of the PCR

253  product, and the single-strand molecule was ligated using DNA ligase. The remaining linear

254  molecule was digested with the exonuclease, finally obtaining a single-strand circular DNA

255  library. We sequenced the DNA library using BGISEQ-500 and MGISEQ-T7 with a pair-end

256  read length of 100bp.

257

258  **Raw data preprocessing**

259  We used the FastQC v0.11.8 [12] to assess overall sequencing quality for MGI and Illumina

260  sequencing platforms. PCR duplications (reads were considered duplicates when forward read

261  and reverse read of the two paired-end reads were identical) were detected by the PRINSEQ

262  v0.20.4 (PRINSEQ, RRID:SCR_005454) [25]. The random sequencing error rate was

263  calculated by measuring the occurrence of 'N' bases at each read position in raw reads. Reads

264  with sequencing adapter contamination were examined according to the manufacturer's adapter

265  sequences          (Illumina          sequencing          adapter          left          =

266  "*GATCGGAAGAGCACACGTCTGAACTCCAGTCAC*", Illumina sequencing adapter right =

267 "*GATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT*", MGI sequencing adapter left =

268 "*AAGTCGGAGGCCAAGCGGTCTTAGGAAGACAA*", and MGI sequencing adapter right =

269 "*AAGTCGGATCGTAGCCATGTCGTTCTGTGAGCCAAGGAGTTG*"). We conducted base

270 quality filtration of raw reads using the NGS QC Toolkit v2.3.3 (cutoff read length for high

271 quality 70; cutoff quality score, 20) (NGS QC Toolkit, RRID:SCR_005461) [26]. We used

272 clean reads after removing low-quality reads and adapter containing reads for the mapping step.


273


274 **Mapping, variant calling, and coverage calculation**

275 After the filtering step, clean reads were aligned to the human reference genome (GRCh38)

276 using BWA-MEM v0.7.12, and duplicate reads were removed using Picard v2.6.0 (Picard,

277 RRID:SCR_006525) [27]. After removing duplicate reads, we down-sampled the deduplicated

278 clean reads of all the sequencing platforms to 24× coverage according to the amount of the

279 deduplicated clean reads of HiSeq2000 for a fair comparison. Realignment and base score

280 recalibration of the bam file was processed by GATK v3.3. Single nucleotide variants, short

281 insertions, and deletions were called with the GATK (Unifiedgenotyper, options --

282 output_mode EMIT_ALL_SITES --genotype_likelihoods_model BOTH). The resulting

283 variants were annotated with the dbSNP (v153) database [28]. Coverage was calculated for

284 each nucleotide using SAMtools v1.9 (SAMTOOLS, RRID:SCR_002105) [29]. We defined a

285 specific covered region based on the 100 bp non-overlapping windows by calculating the

286 average depth of the windows. We used more than five times the difference with an average

287 depth in each window between MGI and Illumina platforms. GC coverage for raw reads and

288 the genome was calculated by the average %GC of the 100bp non-overlapping windows.


289

**Variant comparison and concordance rate with SNP genotyping**

The chromosome position and genotype of each variant called from each sequencing platform was used to identify the relationship between seven sequencing platforms. We compared 1,036,417 loci found on one or more platforms for locations where genotypes were determined on all the seven platforms. An unrooted tree was generated using FastTree v2.1.10 (FastTree, RRID:SCR_015501) [30] with the generalized time-reversible (GTR) model. For calculating the concordance rate between SNP genotyping and WGS-based genotype, the coordinates of SNP genotyping data were converted to GRCh38 assembly using the UCSC LiftOver tool [31]. We removed unmapped positions and indel markers and used only markers that were present on the autosomal chromosomes.

# Availability of Supporting Data and Materials

All sequences generated in this study, including the HiSeq2000, HiSeq2500, HiSeq4000, HiSeqX10, NovaSeq6000, BGISEQ-500, and MGISEQ-T7 sequencing reads, were deposited in the NCBI Sequence Read Archive database under BioProject PRJNA600063. All data will be hosted and distributed from http://biosequencer.org.

# Additional Files

Additional file 1: **Figure S1**. Distribution of nucleotide quality across seven sequencing platforms. **Figure S2**. Base quality filtration statistics for seven sequencing platforms. **Figure S3**. Random error ratio for seven sequencing platforms. **Figure S4**. Insert-size distributions for

311　seven sequencing platforms. **Figure S5**. The coverage distribution of two MGI and five

312　Illumina platforms. **Figure S6**. Depth distribution of chromosome 8. **Figure S7**. Depth

313　distribution of chromosome 12. **Figure S8**. Depth distribution of chromosome 18. **Figure S9**.

314　Depth distribution of chromosome 20. **Figure S10**. GC distribution of platform-specific

315　covered regions. **Figure S11**. The GC composition distribution of the human genome and

316　sequencing reads. **Table S1**. Base quality summary. **Table S2**. Duplicate reads, random error

317　base, and adapter read rate. **Table S3**. The putatively erroneous $K$-mers ($\leq 3$ $K$-mer depth) for

318　seven sequencing platforms. **Table S4**. Statistics of clean reads for seven sequencing platforms.

319　**Table S5**. Statistics for platform-specific covered regions. **Table S6**. The number of shared

320　SNVs for seven sequencing platforms. **Table S7**. The number of SNVs that were singleton or

321　not found in a specific platform. **Table S8**. Genotype concordance rate among seven

322　sequencing platforms. **Table S9**. Genotype comparison between SNP genotyping and WGS.

323

# 324　**List of abbreviations**

325　PE: paired-end;

326　WGS: whole-genome sequencing;

327　BWA: burrows-wheeler aligner;

328　SNVs: single nucleotide variants;

329　indels: insertions and deletions;

330　Ti/Tv: transition/transversion;

331　GATK: Genome Analysis ToolKit;

## Competing Interests

## Funding

## Authors' contributions

J.B. supervised and coordinated the project. J.B. and Y.S.C. conceived and designed the experiments. H.M.K., S.J., O.C., J.H.J., H.Y.L., and Y.Y. conducted the bioinformatics data processing and analyses. H.M.K., S.J., D.M.B., and J.B. wrote and revised the manuscript. A.B. and H.S.K. reviewed and edited the manuscript. All authors have read and approved the final manuscript.

## Acknowledgments

# References

354    1.    Wetterstrand K. DNA sequencing costs: data: data from the NHGRI Genome
355          Sequencing    Program    (GSP).    2018,    Available    online    at:
356          https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data
357          (12 March 2020, date last accessed)

358    2.    Huddleston J, Chaisson MJP, Steinberg KM, et al. Discovery and genotyping of
359          structural variation from long-read haploid genome sequence data. Genome Res
360          2017;**27**(5):677-85.

361    3.    Goodwin S, McPherson JD and McCombie WR. Coming of age: ten years of next-
362          generation sequencing technologies. Nat Rev Genet 2016;**17**(6):333-51.

363    4.    Drmanac R, Sparks AB, Callow MJ, et al. Human genome sequencing using unchained
364          base reads on self-assembling DNA nanoarrays. Science 2010;**327**(5961):78-81.

365    5.    Huang J, Liang X, Xuan Y, et al. A reference human genome dataset of the BGISEQ-
366          500 sequencer. Gigascience 2017;6(5):1-9.

367    6.    Chen J, Li X, Zhong H, et al. Systematic comparison of germline variant calling
368          pipelines cross multiple next-generation sequencers. Sci Rep. 2019;9 1:9345.
369          doi:10.1038/s41598-019-45835-3.

370    7.    Jeon SA, Park JL, Kim JH, et al. Comparison of the MGISEQ-2000 and Illumina HiSeq
371          4000 sequencing platforms for RNA sequencing. Genomics Inform. 2019;17 3:e32.
372          doi:10.5808/GI.2019.17.3.e32.

373    8.    Fang C, Zhong H, Lin Y, et al. Assessment of the cPAS-based BGISEQ-500 platform
374          for    metagenomic    sequencing.    Gigascience.    2018;7    3:1-8.
375          doi:10.1093/gigascience/gix133.

376    9.    Lu R, Zhao X, Li J, et al. Genomic characterisation and epidemiology of 2019 novel

377  coronavirus: implications for virus origins and receptor binding. Lancet. 2020;395

378  10224:565-74. doi:10.1016/S0140-6736(20)30251-8.

379  10.  Kim D, Lee JY, Yang JS, et al. The Architecture of SARS-CoV-2 Transcriptome. Cell.

380  2020;181 4:914-21 e10. doi:10.1016/j.cell.2020.04.011.

381  11.  Cho YS, Kim H, Kim HM, et al. An ethnically relevant consensus Korean reference

382  genome is a step towards personal reference genomes. Nat Commun. 2016;7:13637.

383  doi:10.1038/ncomms13637.

384  12.  Andrews S. FastQC: a quality control tool for high throughput sequence data.

385  Babraham Bioinformatics, Babraham Institute, Cambridge, United Kingdom, 2010.

386  13.  Zhao L, Xie J, Bai L, et al. Mining statistically-solid k-mers for accurate NGS error

387  correction. BMC Genomics. 2018;19 Suppl 10:912. doi:10.1186/s12864-018-5272-y.

388  14.  Liu B, Shi Y, Yuan J, et al. Estimation of genomic characteristics by analyzing k-mer

389  frequency in de novo genome projects. arXiv preprint arXiv:13082012. 2013.

390  15.  Dohm JC, Lottaz C, Borodina T, et al. Substantial biases in ultra-short read data sets

391  from high-throughput DNA sequencing. Nucleic Acids Res. 2008;36 16:e105.

392  doi:10.1093/nar/gkn425.

393  16.  Kozarewa I, Ning Z, Quail MA, et al. Amplification-free Illumina sequencing-library

394  preparation facilitates improved mapping and assembly of (G+C)-biased genomes. Nat

395  Methods. 2009;6 4:291-5. doi:10.1038/nmeth.1311.

396  17.  Aird D, Ross MG, Chen WS, et al. Analyzing and minimizing PCR amplification bias

397  in Illumina sequencing libraries. Genome Biol. 2011;12 2:R18. doi:10.1186/gb-2011-

398  12-2-r18.

399  18.  Oyola SO, Otto TD, Gu Y, et al. Optimizing Illumina next-generation sequencing

400  library preparation for extremely AT-biased genomes. BMC Genomics. 2012;13:1.

401  doi:10.1186/1471-2164-13-1.

402    19.    Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-
403           MEM. arXiv preprint arXiv:13033997. 2013.

404    20.    DePristo MA, Banks E, Poplin R, et al. A framework for variation discovery and
405           genotyping using next-generation DNA sequencing data. Nat Genet. 2011;43 5:491-8.
406           doi:10.1038/ng.806.

407    21.    Van der Auwera GA, Carneiro MO, Hartl C, et al. From FastQ data to high confidence
408           variant calls: the Genome Analysis Toolkit best practices pipeline. Curr Protoc
409           Bioinformatics. 2013;43:11 0 1- 0 33. doi:10.1002/0471250953.bi1110s43.

410    22.    McKenna A, Hanna M, Banks E, et al. The Genome Analysis Toolkit: a MapReduce
411           framework for analyzing next-generation DNA sequencing data. Genome Res. 2010;20
412           9:1297-303. doi:10.1101/gr.107524.110.

413    23.    Patch AM, Nones K, Kazakoff SH, et al. Germline and somatic variant identification
414           using BGISEQ-500 and HiSeq X Ten whole genome sequencing. PLoS One. 2018;13
415           1:e0190264. doi:10.1371/journal.pone.0190264.

416    24.    Rhodes J, Beale MA and Fisher MC. Illuminating choices for library prep: a
417           comparison of library preparation methods for whole genome sequencing of
418           Cryptococcus neoformans using Illumina HiSeq. PLoS One. 2014;9 11:e113501.
419           doi:10.1371/journal.pone.0113501.

420    25.    Schmieder R and Edwards R. Quality control and preprocessing of metagenomic
421           datasets. Bioinformatics. 2011;27 6:863-4. doi:10.1093/bioinformatics/btr026.

422    26.    Patel RK and Jain M. NGS QC Toolkit: a toolkit for quality control of next generation
423           sequencing data. PLoS One. 2012;7 2:e30619. doi:10.1371/journal.pone.0030619.

424    27.    Institute B: Picard: A set of command line tools (in Java) for manipulating high-
425           throughput sequencing (HTS) data and formats such as SAM/BAM/CRAM and VCF.
426           http://broadinstitute.github.io/picard/ (2018).

427  28.  Sherry ST, Ward MH, Kholodov M, et al. dbSNP: the NCBI database of genetic

428      variation. Nucleic Acids Res. 2001;29 1:308-11. doi:10.1093/nar/29.1.308.

429  29.  Li H, Handsaker B, Wysoker A, et al. The Sequence Alignment/Map format and

430      SAMtools. Bioinformatics. 2009;25 16:2078-9. doi:10.1093/bioinformatics/btp352.

431  30.  Price MN, Dehal PS and Arkin AP. FastTree 2--approximately maximum-likelihood

432      trees    for    large    alignments.    PLoS    One.    2010;5    3:e9490.

433      doi:10.1371/journal.pone.0009490.

434  31.  Kuhn RM, Haussler D and Kent WJ. The UCSC genome browser and associated tools.

435      Brief Bioinform. 2013;14 2:144-61. doi:10.1093/bib/bbs038.

436

## Figures

**Figure 1. Distribution of *K*-mer frequency for 21-mers using raw reads from seven sequencing platforms.** The x-axis represents *K*-mer depth, and the y-axis represents the proportion of *K*-mer, as calculated by the frequency at that depth divided by the total frequency at all depths.

**Figure 2. GC-bias plots for seven sequencing platforms.** Unbiased coverage is represented by a horizontal dashed line at relative coverage = 1. A relative coverage below 1 indicates lower than expected coverage and above 1 indicates higher than expected coverage.

**Figure 3. An unrooted tree for seven sequencing platforms showing the similarity of the variant calling.** Numbers of nodes denote bootstrap values based on 1,000 replicates.

## Tables

### Table 1. Raw read statistics for seven sequencing platforms

| | Illumina platforms | | | | | MGI platforms | |
|---|---|---|---|---|---|---|---|
| Metrics | HiSeq2000 | HiSeq2500 | HiSeq4000 | HiSeqX10 | NovaSeq6000 | BGISEQ-500 | MGISEQ-T7 |
| Production date | 2012 | 2015.03 | 2015.10 | 2015.12 | 2019.04 | 2017.04 | 2019.09 |
| Quality range | Illumina 1.5+ | Illumina 1.8+ | Illumina 1.8+ | Illumina 1.8+ | Illumina 1.8+ | Illumina 1.8+ | Illumina 1.8+ |
| # of Total read | 1,044M | 1,500M | 629M | 833M | 833M | 1,171M | 1,035M |
| Read length (bp) | 90 PE | 101 PE | 151 PE | 151 PE | 151 PE | 100 PE | 100 PE |
| Total bases | 94 Gb | 151.5 Gb | 95 Gb | 125.8 Gb | 125.8 Gb | 117.1 Gb | 103.4 Gb |
| Sequencing depth (×, based on 3 Gb) | 31.31 | 50.52 | 31.65 | 41.94 | 41.94 | 39.04 | 34.49 |

### Table 2. Mapping and coverage statistics

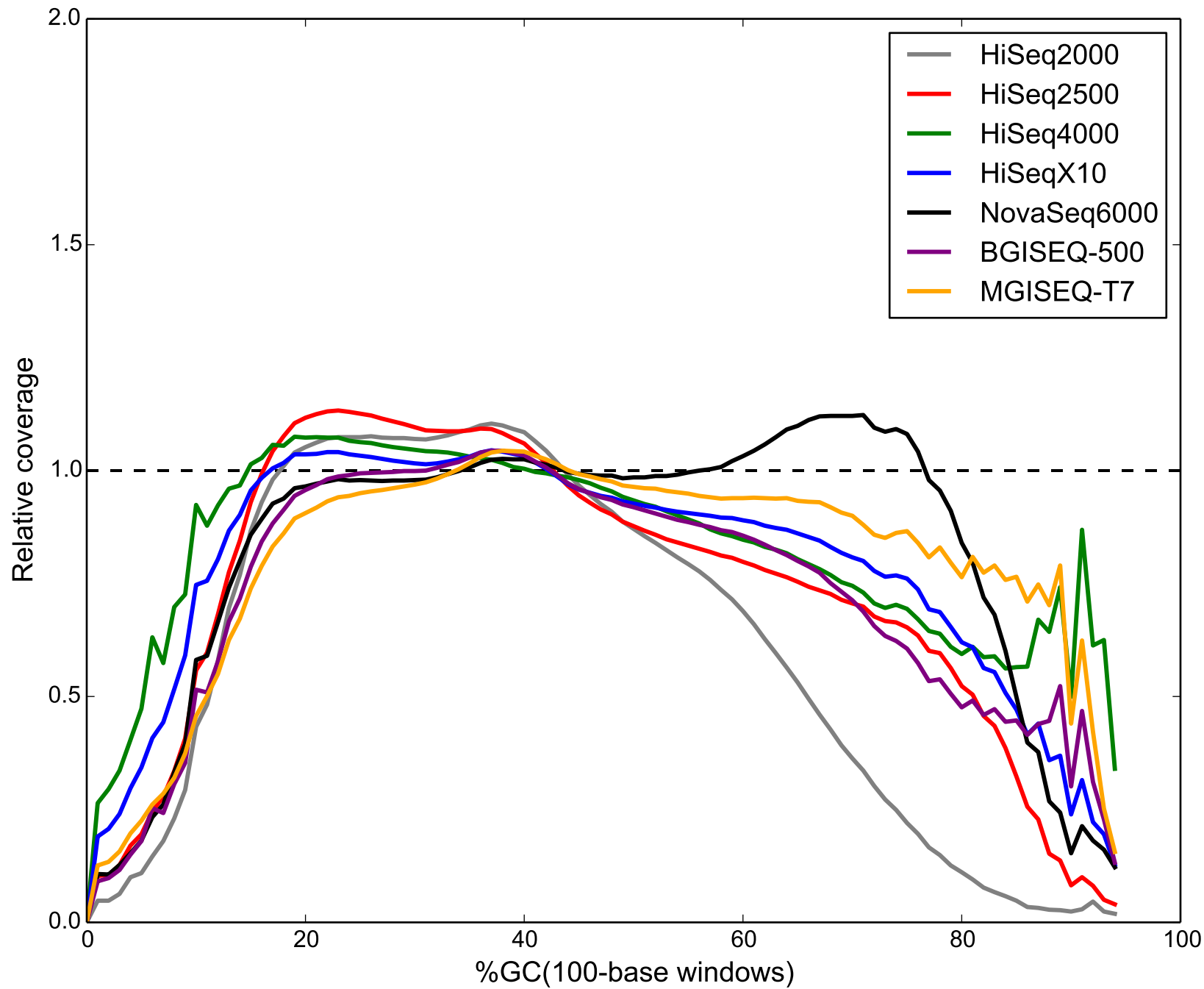| Metrics | HiSeq2000 | HiSeq2500 | HiSeq4000 | HiSeqX10 | NovaSeq6000 | BGISEQ-500 | MGISEQ-T7 |
|---|---|---|---|---|---|---|---|
| # of clean reads | 935,951,974 | 1,050,028,628 | 512,891,970 | 705,987,420 | 706,000,000 | 1,060,837,856 | 991,021,996 |
| Read length | 90 | 101 | 151 | 151 | 151 | 100 | 100 |
| Clean bases (Gb) | 84.23 | 106.05 | 77.45 | 106.60 | 106.6 | 106.08 | 99.1 |
| Clean read depth (based on 3 Gb, ×) | 28.08 | 35.35 | 25.82 | 35.53 | 35.54 | 35.36 | 33.03 |
| Mapping rate | 99.986% | 99.999% | 99.990% | 99.999% | 99.9996% | 99.983% | 99.999% |
| Properly mapped rate* | 96.67% | 98.30% | 97.24% | 96.91% | 97.15% | 97.44% | 98.17% |
| Duplicate rate | 15.35% | 3.01% | 3.19% | 5.08% | 3.39% | 2.56% | 8.77% |
| Duplicate clean read depth (×) | 23.90 | 34.29 | 24.99 | 33.73 | 34.33 | 34.46 | 30.14 |
| Down-sampled depth (×) | 23.90 | 23.90 | 23.90 | 23.90 | 23.90 | 23.90 | 23.90 |
| Coverage | 99.68% | 99.82% | 99.71% | 99.81% | 99.76% | 99.83% | 99.83% |
| Coverage at least 5× | 98.62% | 99.30% | 98.37% | 99.30% | 99.19% | 99.34% | 99.24% |
| Coverage at least 10× | 94.63% | 96.65% | 93.98% | 97.05% | 96.89% | 97.05% | 96.61% |
| Coverage at least 15× | 85.10% | 88.54% | 85.08% | 90.23% | 90.36% | 90.11% | 89.36% |

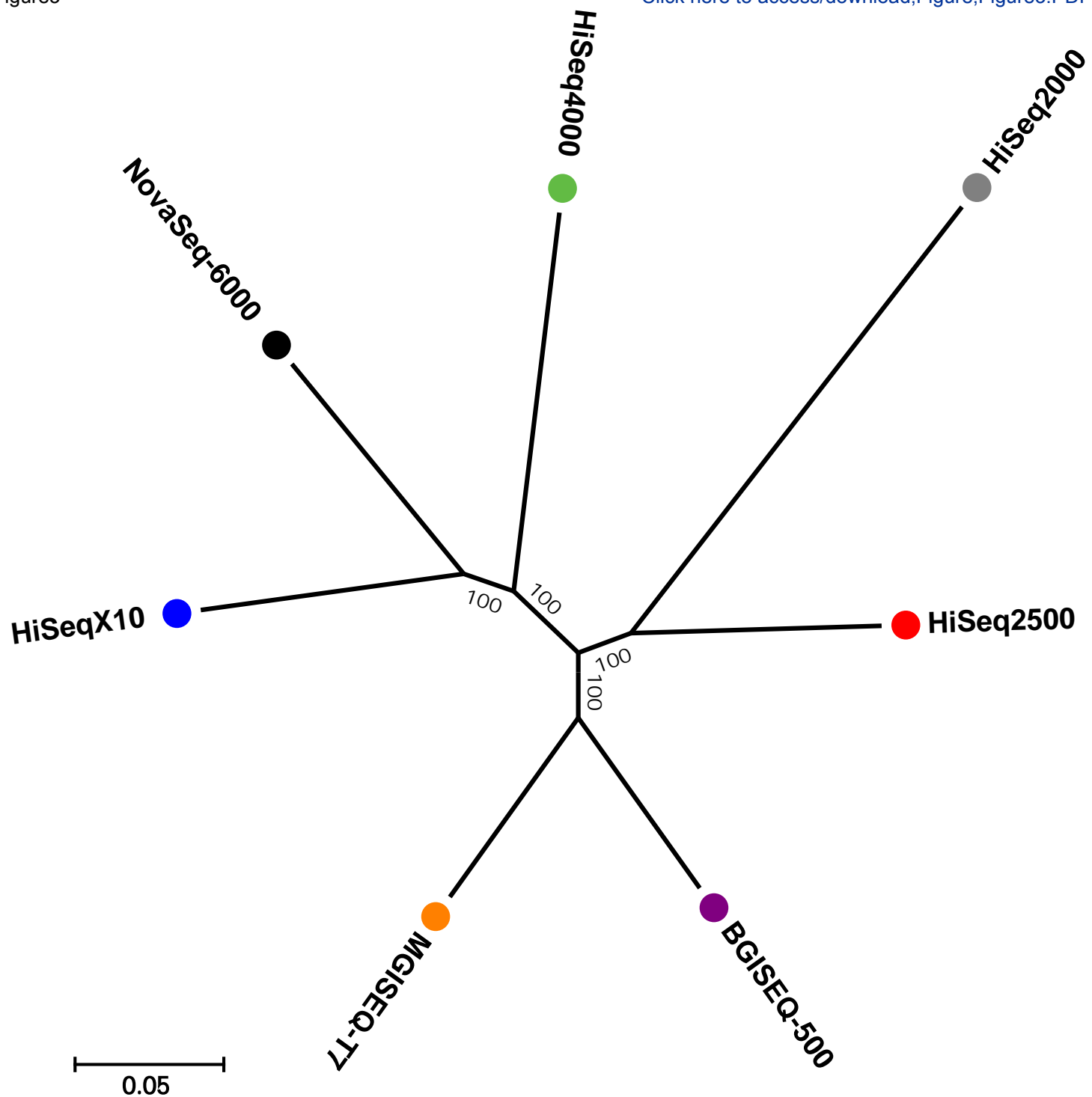* Both of the read mates are in the correct orientation.

459 **Table 3. Variant statistics of Illumina and MGI sequencing platforms.**

| | Metrics | HiSeq2000 | HiSeq2500 | HiSeq4000 | HiSeqX10 | NovaSeq6000 | BGISEQ-500 | MGISEQ-T7 |
|---|---|---|---|---|---|---|---|---|
| Reference homozygous | | 2,839,358,003 | 2,855,619,759 | 2,855,062,233 | 2,864,272,103 | 2,861,198,782 | 2,851,898,568 | 2,853,066,635 |
| # of no call positions | | 80,241,142 | 63,980,549 | 64,532,078 | 55,244,498 | 58,311,103 | 67,747,107 | 66,584,361 |
| No call rate | | 2.74% | 2.19% | 2.21% | 1.89% | 1.99% | 2.32% | 2.28% |
| SNVs | Total SNVs | 4,133,925 | 4,132,468 | 4,138,296 | 4,216,589 | 4,223,612 | 4,088,645 | 4,082,103 |
| | Total SNVs in dbSNP | 4,094,212 | 4,114,993 | 4,112,253 | 4,198,005 | 4,184,100 | 4,070,101 | 4,064,986 |
| | dbSNP rate | 99.04% | 99.58% | 99.37% | 99.56% | 99.06% | 99.55% | 99.58% |
| | Singleton | 159,429 | 78,109 | 98,574 | 100,158 | 104,052 | 52,127 | 51,978 |
| | Singleton in dbSNP | 126,762 | 68,673 | 78,361 | 89,094 | 73,177 | 41,092 | 41,743 |
| | dbSNP rate for Singleton | 79.51% | 87.92% | 79.49% | 88.95% | 70.33% | 78.83% | 80.31% |
| | Homozygous | 1,703,616 | 1,690,878 | 1,704,813 | 1,708,639 | 1,714,752 | 1,688,328 | 1,689,834 |
| | Heterozygous | 2,430,309 | 2,441,590 | 2,433,483 | 2,507,950 | 2,508,860 | 2,400,317 | 2,392,269 |
| | Het/Hom ratio | 1.43 | 1.44 | 1.43 | 1.47 | 1.46 | 1.42 | 1.42 |
| | Ti/Tv ratio | 1.91 | 1.92 | 1.9 | 1.88 | 1.85 | 1.92 | 1.92 |
| Indels | Total Indels | 526,504 | 546,918 | 491,899 | 689,357 | 708,062 | 703,873 | 631,163 |
| | Total Indels in dbSNP | 524,738 | 544,866 | 489,777 | 686,916 | 705,553 | 701,802 | 629,314 |
| | dbSNP rate | 99.66% | 99.62% | 99.57% | 99.65% | 99.65% | 99.71% | 99.71% |
| | Singleton | 7,864 | 7,444 | 8,094 | 17,036 | 23,596 | 41,384 | 12,092 |
| | Singleton in dbSNP | 7,612 | 7,259 | 7,915 | 16,784 | 23,303 | 41,183 | 11,964 |
| | dbSNP rate for Singleton | 96.80% | 97.51% | 97.79% | 98.52% | 98.76% | 99.51% | 98.94% |

460

Figure1

Figure2

Figure3

Click here to access/download
**Supplementary Material**
Sequencing_Platform_Comparison_Supplementary_revision_final.docx

Respond to reviewers

Click here to access/download
**Supplementary Material**
Sequencing_Platform_Comparison_RevisionNote_final.docx

**GIGA-D-20-00072**

Dear *GigaScience* editors,

Thank you for considering our manuscript for publication in *GigaScience*.

The reviewer #1' criticisms on "the sequencing data of seven sequencing platforms need to have the same genome coverage" was useful to improve the quality of our analyses and manuscripts.

To accommodate the reviewer #1' criticisms, we matched the seven sequencing platforms to the same genome coverage and re-analyzed the downstream analyses, including variant comparison, platform-specific covered regions, and concordance rate with SNP genotyping, to remove the bias due to the different genome coverage. As a result, we could compare the seven sequencing platforms more objectively.

We have also edited the manuscript to accommodate the reviewers' concerns on clarity and better presentation of the results. Please see the detailed point-by-point revision notes that are submitted on-line.

We hope that our revision will be suitable for your journal's standards.

Sincerely yours,

Dan M. Bolser & Jong Bhak
dan@geromics.co.uk; jongbhak@genomics.org