

1
2
3 **1 A reference human genome dataset of the BGISEQ-500 sequencer**
4
5

6
7 2 Jie Huang^{1*}, Xinming Liang^{2*}, Yuankai Xuan^{3*}, Chunyu Geng², Yuxiang Li², Haorong
8
9 3 Lu², Shoufang Qu¹, Xianglin Mei³, Hongbo Chen¹, Ting Yu¹, Nan Sun¹, Junhua Rao²,
10
11 4 Jiahao Wang⁴, Wenwei Zhang², Ying Chen², Sha Liao², Hui Jiang², Xin Liu²,
12
13 5 Zhaopeng Yang^{1#}, Feng Mu^{2#} and Shangxian Gao^{1#}
14
15

16
17 6
18 7 *These authors contribute equally to the article

19 8 #Correspondence: Shangxian Gao, gaoshangxian@126.com, Feng Mu mufeng@genomics.cn
20 9 and Zhaopeng Yang yangzp@nifdc.org.cn

21 10 ¹ National Institutes for food and drug Control (NIFDC), Beijing 100050, P. R. China

22 11 ² BGI-Shenzhen, Shenzhen 518083, P. R. China

23 12 ³ State Food and Drug Administration Hubei Center for Medical Equipment Quality
24 13 Supervision and Testing, Wuhan 430000, P. R. China

25 14 ⁴ BGI-Qingdao, Qingdao 266555, China
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Abstract**

2 **Background:** BGISEQ-500 is a new desktop sequencer developed by BGI. Using
3 DNA nanoball (DNB) and combinational probe anchor synthesis (cPAS) developed
4 from Complete Genomics™ sequencing technologies, it generates short reads at a
5 large scale.

6 **Findings:** Here, we present the first human whole genome sequencing dataset of
7 BGISEQ-500. The dataset was generated by sequencing the widely used cell line,
8 HG001 (NA12878) in two sequencing runs of paired-end 50 bp (PE50) and two
9 sequencing runs of paired-end 100 bp (PE100). We also include examples of the raw
10 images from the sequencer for reference. Finally, we identified variations using this
11 dataset, estimated the accuracy of the variations and compared to that of the variations
12 identified from similar amounts of publicly available HiSeq2500 data.

13 **Conclusions:** We found similar SNP detection accuracy for the BGISEQ-500 PE100
14 data (false positive rate, FPR 0.00020%, and sensitivity 96.20%) comparing to the
15 PE150 HiSeq2500 data (FPR 0.00017% and sensitivity 96.60%), better than the PE50
16 data (FPR 0.0006% and sensitivity 94.15%). But for insertions and deletions (indels),
17 we found lower accuracy for BGISEQ-500 data (FPR 0.00069% and 0.00067% for
18 PE100 and PE50 respectively, sensitivity 88.52% and 70.93%) than the HiSeq2500
19 data (FPR 0.00032% and sensitivity 96.28%). Our dataset can serve as the reference
20 dataset providing basic information not just for future development, but also for all
21 research and applications based on the new sequencing platform.

22 **Keywords:** Genomics, sequencing, Next Generation Sequencing, BGISEQ-500

1 **Data Description**

2 Massively parallel sequencing technologies (also called as the second generation
3 sequencing) generate large amount of data with lower cost, shorter reads and higher
4 single base error rate comparing to Sanger sequencing technology[1]. With the large
5 amount of data and well-developed analysis tools, second generation sequencing data
6 can be used to effectively and accurately identify genomic variations[2]. Thus it has
7 been widely applied in both research and application[3]. Currently there are several
8 commercially available second generation sequencing platforms with differing
9 performance and data features[4, 5]. With more and more research areas and
10 applications to apply sequencing to, new sequencing platforms are being developed at
11 a rapid pace. The BGISEQ-500 sequencer was first announced by BGI in October,
12 2015. It was developed based on the Complete GenomicsTM sequencing technologies,
13 and applied DNA NanoBalls (DNBs) technology[6] for sequencing library
14 construction and combined primer anchor synthesis (cPAS) for sequencing. We
15 present here a dataset generated from the BGISEQ-500 sequencer, including examples
16 of the raw images and the final sequences. We also conducted variation calling using
17 this dataset and compared the variation calling result to that from other sequencers.
18 This dataset can be served as a useful reference for the community to develop
19 bioinformatics methods and sequencing based applications on this new sequencing
20 platform.

21 **DNA preparation**

22 NA12878 cell line (RRID:CVCL_7526) genomic DNA was ordered from the Coriell
23 Institute, and contained 50 µg per tube. The genomic DNA was quantified by Qubit
24 3.0 fluorometer (Life Technologies, Paisley, UK) and the integrity was qualified on
25 the 2% agarose gel to make sure the genomic DNA molecular was larger than 23kb
26 and not substantially degraded.

27 **Sequencing library preparation**

28 For the sequencing library construction, the NA12878 genomic DNA was fragmented

1 by ultrasound on Covaris E220 (Covaris, Brighton, UK) to DNA fragments between
2 50 bp~800bp according to the manufacturer's instructions. The fragmented DNA was
3 further selected to 100bp~300bp by AMPure XP beads (AGENCOURT). The selected
4 DNA fragments were then repaired to obtain a blunt end and modified at 3'end to get
5 a dATP as a sticky end. The dTTP tailed adapter sequence was ligated to both ends of
6 the DNA fragments. The ligation product was then amplified for 8cycles and
7 subjected to the following single strand circularization process. The PCR product was
8 heat-denatured together with a special molecule which was reverse complemented to
9 one special strand of the PCR product, and the single strand molecule was ligated
10 using DNA ligase. The remaining linear molecule was digested with the exonuclease,
11 finally obtaining a single strand circular DNA library (Figure 1a).

12 **Sequencing**

13 We conducted sequencing according to the BGISEQ-500 protocol (Figure 1b). There
14 were three steps including making DNBs, loading DNBs and sequencing. For making
15 DNBs, a 6 ng single strand circular DNA library was first PCR amplified for 10
16 minutes in an 80 μ l reaction volume with pure water, buffer and DNB polymerase.
17 After the PCR reaction, 20 μ l DNBs stopping buffer was added to terminate the PCR
18 reaction. Finally, we used the Qubit® ssDNA Assay Kit to quantify the DNBs on a
19 Qubit® Fluorometer (concentration ≥ 10 ng/ μ L).

20 For loading DNBs, we first added 33 μ l DNBs loading buffer to DNBs product from
21 the last step, and the mixture was placed on the BGIDL-50 (the sample preparation
22 machine). Then we selected the DNBs loading process (Version: sample load 2.0) to
23 load DNB onto the sequencing chip, which included 96 minutes' loading time and 30
24 minutes' incubation at room temperature.

25 Finally, for sequencing, we followed to the BGISEQ-500 protocol. We selected
26 sequence control software Version 1.1.0.10003, sequence process Version 1.0.06 and
27 Zebrecall process Version 0.5.0.13875 (the base calling software, and a detailed
28 description can be found in the next section) for sequencing. Sequencing was initiated

1 after the sequencing reagents pre-loaded and sequencing chip installed, and this
2 process was finished in ~72 hours.

3 **Base calling and raw images**

4 During sequencing, four channels of 16-bit grey scale images were captured by high
5 resolution sCMOS with ~5.5 million pixels per image. About ~570K DNBs were
6 loaded onto the grid-patterned arrays of spots which were photolithographically
7 etched and surface modified on the sequencing chip. The spots were illuminated by
8 the lasers with different wavelengths. Intensity from neighboring channel would also
9 be observed due to crosstalk effect. The sequences of DNBs were base-called by the
10 software Zebra call (base calling software developed for BGISEQ-500). After
11 background subtraction and registration of images from 4 channels, intensities of
12 DNBs were extracted according to a template of grid-pattern. Correction within
13 channels and neighbor cycles was applied to increase the quality and stabilization.
14 The cross-talk between intensities from the four fluorophores is caused by the
15 imperfect wavelength filtering of optical filters isolating the bands of wavelength
16 from the four types of fluorophore molecules. A regression technique can identify
17 correlations in our intensity data and correct for them. For example, in order to correct
18 the crosstalk between 2 channels (C and G), the correction of the C background
19 intensity can be found by linear regression after eliminating DNBs that do have true
20 signal in the C channel. Such DNBs can be identified by searching for DNBs that
21 have the C intensity as the maximum of the four intensities and, to retain just DNBs
22 that are not too dim or noisy, we take only DNBs that have less than 80% of the C
23 intensity for the remaining three other intensities. This leaves us with reasonably well
24 performing DNBs that most likely do not contain C at the currently interrogated
25 position. Linear regression was then carried out for these background G intensities as
26 a function of the C intensities. All of the G intensities can in turn be corrected for this
27 cross-talk from the C channel by subtracting from them the expected background
28 intensity produced by a given C intensity. Such regression can be done for each

1 channel and simultaneously correcting for all of the correlations outlined above using
2 a multiple linear regression. After all correction steps, the base with highest
3 probability will be called according to the scale of intensities. When the whole
4 sequencing was finished, the binary file with bases and quality score were converted
5 into FASTQ format with Phred+33 quality score.

6 In order to map the base call quality to Phred+33 score, a prior probability model was
7 constructed by the scale of intensities from channels. Bases were separated to 10404
8 groups according to different parameters which may affect the confidence level of
9 base calling. The base-calling error probabilities(P) of each group was calculated by
10 the mismatch distribution from repeated sequencing of the standard reference genome.

11 Quality scores (Q) were calculated by the definition of Phred+33 quality scores:

$$Q = -10 \log_{10} P$$

12 A huge table was constructed and hard coded into the base call program to look up a
13 corresponding quality score by different parameters.

14 An example dataset of the images was included and the base calling process was
15 illustrated in Figure 2.

16 **Results**

17 **Sequencing data summary**

18 The sequencing data consists of four lanes, with two of PE100 and the other two of
19 PE50 (Table 1). First, we analyzed the sequencing quality by identifying the low
20 quality reads. Although previous studies revealed that raw data filtering would not
21 substantially affect variation calling result [7, 8], we found slight different
22 performances of variation calling using different raw data filtering criteria (Table S1).
23 Thus, we determined low quality reads as reads which had more than 10% bases with
24 sequencing quality lower than 10, and reads which had more than 1% Ns (ambiguous
25 bases). In this way, we identified 11.9% (9.2% low quality reads and 2.7% ambiguous
26 reads) low quality reads in PE100 data and 12.3% low quality raw reads in PE50 data

1 (5.4% low quality reads and 6.9% ambiguous reads). In order for comparison, we
2 selected similar amount of data (8 sequencing libraries and 16 lanes, PE150 reads,
3 ~98.5 G bp data) from a public Illumina HiSeq2500 dataset of this cell line generated
4 by GIAB (Genome in a Bottle)[9]. Using the same criteria for low quality
5 identification, we identified 7.95% low quality reads (7.7% low quality reads and 0.25%
6 ambiguous reads). Excluding these low quality reads, we then further analyzed the
7 reads quality by plotting the distributions of base quality scores and GC content
8 against those of the HiSeq2500 data (Figure 3). Thus we found higher proportion of
9 low quality reads, more stable base quality distribution along the reads (Figure 3 a-b)
10 and lower overall single base quality scores (Figure 3 c). And we observed some
11 secondary peak in the GC content distribution of BGISEQ-500 data, indicating higher
12 GC bias (Figure 3 d).

13 **Variation calling and false positive/negative ratios estimation**

14 In order to further depict the data quality and test applications of the new sequencing
15 platform, we carried out variation calling using this dataset. We adapted the widely
16 used pipeline (BWA[10] and GATK[11-13], an illustration of the pipeline and key
17 parameters can be found in Figure 4a) for variation calling. We observed higher
18 mapping rate, similar sequencing coverage and similar sequencing uniformity of the
19 two BGISEQ-500 datasets compared to the HiSeq2500 dataset (Table 2). The lower
20 unique mapping rate probably reflected the shorter read length of the dataset (2×50 bp
21 and 2×100 bp comparing to 2×150 bp). We also observed slightly higher duplication
22 rate and comparable mismatch rate in the BGISEQ-500 PE100 dataset comparing to
23 the HiSeq2500 data (Table 2).

24 In total, we identified ~3.4 million SNPs using the BGISEQ-500 datasets (3.45
25 million for PE50 data and 3.48 million for PE100 data), more than 3.6 million SNPs
26 identified using HiSeq2500 data (Table 3). While for indels (insertion and deletions),
27 we identified 842,058 from BGISEQ-500 PE100 data, comparing to 553,842
28 identified from BGISEQ-500 PE50 data. Using the HiSeq2500 data, we identified

1 733,797 indels. The SNPs identified using BGISEQ-500 datasets were similar to
2 those identified from HiSeq2500 data in different features including dbSNP rate,
3 proportion of SNPs in different regions related to genes and Ti/Tv
4 (transition/transversion) ratio, which indirectly reflected the SNP accuracy. We also
5 observed similar situation for indels.

6 Further to assess the accuracy of the variations, we used the high confident variations
7 previously identified in NA12878 provided by GIAB (Genome in A Bottle)[14].
8 Using the methods provided by GIAB, we estimated the false positive rates and
9 sensitivity for BGISEQ-500 PE50 and PE100 data compared to those of HiSeq2500
10 data (Table4). The SNP sensitivity was lower for the BGISEQ-500 datasets (96.20%
11 for PE100 and 94.15% for PE50) than HiSeq2500 data (96.60%). And the SNP false
12 positive rate (FPR) was similar for the BGISEQ-500 PE100 data (0.00020%)
13 compared to HiSeq2500 data (0.00017%), and lower than the BGISEQ-500 PE50 data
14 (0.0006%). For indels, BGISEQ-500 PE100 data resulted in worse performance with
15 lower sensitivity (88.52%) than the HiSeq2500 PE150 data with sensitivity of 96.28%.
16 In contrast, HiSeq2500 PE150 data shows lower FPR (0.00032%) than BGISEQ-500
17 PE100 data (0.00069%). The BGISEQ-500 PE50 data resulted in sensitivity of 70.93%
18 and FPR of 0.00067%. The difference performances of indel calling might also be
19 caused by read length difference (50 or 100 bp comparing to 150 bp), in addition to
20 sequencing quality, mapping accuracy, etc.

21 Furthermore, to depict variation calling accuracy in different genomic regions, we
22 compared the false negative rate (FNR), FPR and sensitivity in different genome
23 context given by GIAB (Figure S1). For the coding sequences, data from the two
24 platforms have similar FNR, FPR and sensitivity (3.85% vs. 2.52%, 0.00012% vs.
25 0.00015% and 96.15% vs. 97.48% accordingly). For the regions which are difficult to
26 sequence, including some of the promoters [15], substantially high GC content (>55%)
27 regions, substantially low GC content (<30%) regions, regions with multiple
28 variations (more than 1 variations within 50 bp), regions with compound variations,

1 repeats and segmental duplications, BGISEQ-500 data has a higher FNR, lower
2 sensitivity and lower FPR (Figure S1).

3 **Discussion**

4 Using the new sequencer, BGISEQ-500, we obtained one run of PE50 data and the
5 other run of PE100 data. The raw data were ~135.5 Gbp and ~153.6 Gbp respectively,
6 and were generated from two chips (~72 hours). Thus the sequencing throughput and
7 turnaround time were comparable to HiSeq2500 sequencer Rapid mode v1 (~80 Gbp
8 per single flow cell and ~40 hours). Both the single base quality and read quality
9 (reflected by duplication rate, mapping rate and unique mapping rate) were basically
10 comparable to those of the HiSeq2500 data. Furthermore, the variation calling result
11 was similar to that identified using similar amounts of HiSeq2500 data, further
12 reflecting that the sequencer can be used in different research and applications. With
13 Future improvements over data quality, sequencing length, different and optimized
14 insert sizes of the paired reads, as well as specially modified or designed
15 software/bioinformatics tools, the performance can be further improved. In the
16 meantime, quality of the whole genome sequencing data also reflected feasibility of
17 applying this sequencing platform for other sequencing purposes including
18 transcriptome, epigenome, metagenome etc. From this first reference dataset of
19 sequencing data from BGISEQ-500 sequencer, we provided an overview and some
20 basic information for the new sequencing platform. This dataset can serve as reference
21 for all the research using the BGISEQ-500 sequencing platform. And we anticipate it
22 to help stimulating the further technical improvement and development of novel tools
23 for accurately analyzing this data.

24

1 **Availability of supporting data**

2 The BGISEQ-500 sequences described in this article are available in the GigaDB
3 repository (PE 50[16] and PE 100 [17]), and the European Nucleotide Archive under
4 accession number ERP017158. This GigaDB entry also contains examples of the raw
5 image data including images of all the sequencing cycles in a small region and images
6 of the first and last 10 cycles of the whole flowcell [16]. Future data will also be
7 updated via the GigaDB repository with versions indicated.

8
9 **Abbreviations**

10 bp: base-pair; DNBs: DNA nanoballs; FNR: false negative rate; FPR: false positive
11 rate; GIAB: Genome in A Bottle; PE50: pair-end 50 bp; PE100: pair-end 100 bp;
12 SNPs: Single Nucleotide Polymorphisms; indels: insertions and deletions

13
14 **Competing interests**

15 JH, YX, SQ, XM, HC, TY, NS, ZY and SG are involved in the beta test of the
16 BGISEQ-500 sequencer. X Liang, CG, YL, HL, HJ, X Liu and FM are involved in the
17 BGISEQ-500 sequencer development, library construction technology optimization,
18 base calling software development, or alpha and beta tests.

19
20 **Authors' contributions**

21 JH, ZY, FM and SG designed the project. YX, SQ and CG conducted sample
22 preparation and sequencing library construction. HL, XM, HC, TY and NS conducted
23 sequencing. X Liang, JR, JW, YL, X Liu, HJ, JR, JW, WZ, YC and SL conducted data
24 analysis. X Liu, X Liang, YL, CG, HL, JH and HJ wrote the manuscript.

25
26 **Acknowledgements**

27 This work was supported by following funding: the National High Technology
28 Research and Development Program ("863" Program) of China (Project No.

1 2011AA02A115), the Technology Innovation and Developing Plan of Shenzhen
2 (Project No. CXZZ20140904154910774).

3
4
5 3

6
7 **4 Author details**

8
9 ¹ National Institutes for food and drug Control (NIFDC), Beijing 100050, P. R. China

10
11 ² BGI-Shenzhen, Shenzhen 518083, P. R. China

12
13 ³ State Food and Drug Administration Hubei Center for Medical Equipment Quality
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

8 Supervision and Testing, Wuhan 430000, P. R. China

9 ⁴BGI-Qingdao, Qingdao 266555, China

References

1. Metzker ML: **Sequencing technologies - the next generation.** *Nat Rev Genet* 2010, **11**(1):31-46.
2. Wang J, Wang W, Li R, Li Y, Tian G, Goodman L, Fan W, Zhang J, Li J, Guo Y *et al*: **The diploid genome sequence of an Asian individual.** *Nature* 2008, **456**(7218):60-65.
3. Goodwin S, McPherson JD, McCombie WR: **Coming of age: ten years of next-generation sequencing technologies.** *Nat Rev Genet* 2016, **17**(6):333-351.
4. Mardis ER: **Next-generation sequencing platforms.** *Annu Rev Anal Chem (Palo Alto Calif)* 2013, **6**:287-303.
5. Quail MA, Smith M, Coupland P, Otto TD, Harris SR, Connor TR, Bertoni A, Swerdlow HP, Gu Y: **A tale of three next generation sequencing platforms: comparison of Ion Torrent, Pacific Biosciences and Illumina MiSeq sequencers.** *BMC Genomics* 2012, **13**:341.
6. Drmanac R, Sparks AB, Callow MJ, Halpern AL, Burns NL, Kermani BG, Carnevali P, Nazarenko I, Nilsen GB, Yeung G *et al*: **Human genome sequencing using unchained base reads on self-assembling DNA nanoarrays.** *Science* 2010, **327**(5961):78-81.
7. Hwang S, Kim E, Lee I, Marcotte EM: **Systematic comparison of variant calling pipelines using gold standard personal exome variants.** *Scientific reports* 2015, **5**:17875.
8. Liu Q, Guo Y, Li J, Long J, Zhang B, Shyr Y: **Steps to ensure accuracy in genotype and SNP calling from Illumina sequencing data.** *BMC genomics* 2012, **13 Suppl 8**:S8.
9. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, Weng Z, Liu Y, Mason CE, Alexander N *et al*: **Extensive sequencing of seven human genomes to characterize benchmark reference materials.** *Sci Data* 2016, **3**:160025.
10. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler transform.** *Bioinformatics* 2009, **25**(14):1754-1760.
11. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M *et al*: **The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.** *Genome Res* 2010, **20**(9):1297-1303.
12. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del Angel G, Rivas MA, Hanna M *et al*: **A framework for variation discovery and genotyping using next-generation DNA sequencing data.** *Nat Genet* 2011, **43**(5):491-498.
13. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, Jordan T, Shakir K, Roazen D, Thibault J *et al*: **From FastQ data to high confidence variant calls: the Genome Analysis Toolkit**

1 **best practices pipeline.** *Curr Protoc Bioinformatics* 2013, **43**:11 10 11-33.

2 14. Zook JM, Chapman B, Wang J, Mittelman D, Hofmann O, Hide W, Salit M:
3 **Integrating human sequence data sets provides a resource of benchmark**
4 **SNP and indel genotype calls.** *Nat Biotechnol* 2014, **32**(3):246-251.

5 15. Ross MG, Russ C, Costello M, Hollinger A, Lennon NJ, Hegarty R, Nusbaum
6 C, Jaffe DB. **Characterizing and measuring bias in sequence data.** *Genome*
7 *Biol.* 2013 May 29;14(5):R51.

8 16. Huang, J; Liang, X; Xuan, Y; Geng, C; Li, Y; Lu, H; Qu, S; Mei, X; Chen, H;
9 Yu, T; Sun, N; Rao, J; Wang, J; Zhang, W; Chen, Y; Liao, S; Jiang, H; Liu, X;
10 Yang, Z; Mu, F; Gao, S (2016): **BGISEQ-500 sequencer first reference**
11 **dataset** GigaScience Database. <http://dx.doi.org/10.5524/100252>

12 17. Huang, J; Liang, X; Xuan, Y; Geng, C; Li, Y; Lu, H; Qu, S; Mei, X; Chen, H;
13 Yu, T; Sun, N; Jiang, H; Liu, X; Yang, Z; Mu, F; Gao, S (2017): **An updated**
14 **reference human genome dataset of the BGISEQ-500 sequencer**
15 GigaScience Database. <http://dx.doi.org/10.5524/100274>

16 18. Andrews S: **FastQC, a quality control tool for high throughput sequence**
17 **data.** 2016.

18 19. **Genome in a Bottle Consortium high-confidence VCF and BED datasets**
19 ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/NISTv3.3/NA1287
20 8_GIAB_highconf_CG-III-FB-III-GATKHC-Ion-Solid-10X_CHROM1-X_v3.3_highco
21 <nf.bed>

1 **Figure Legends**

2 **Figure 1. Flowchart of library construction and sequencing.** The library
3 construction includes fragmentation, size selection, end-repair and A-tailing, adaptor
4 ligation and PCR amplification and splint circularization (a). The sequencing includes
5 making DNBs, loading DNBs and sequencing (b).

6
7 **Figure 2. Raw image data processing on the BGISEQ-500 platform. a.**
8 **Registration of images from different channels.** Relative coordinates will be
9 calculated according to the pattern layout of DNBs. **b. Intensity correction between**
10 **channels and cycles.** Correction of the optical and chemical interferences on different
11 channels and the neighbor cycles was applied. **c. Connecting called bases to FASTQ.**
12 Bases from all cycles will be collected and converted to FASTQ format. Phred score
13 calculation and statistics will be applied during the conversion.

14
15 **Figure 3. Quality control of the dataset after data filtering.** Base-wise quality
16 score distributions of the first read (a, from left to right, BGISEQ-500 PE50,
17 BGISEQ-500 PE100 and HiSeq2500 PE150) and the second read (b, from left to right,
18 BGISEQ-500 PE50, BGISEQ-500 PE100 and HiSeq2500 PE150). For each position
19 along the reads, the quality scores of all reads were used to calculate the mean,
20 median and quantile values thus the box plot can be shown. The overall quality
21 score distribution of BGISEQ-500 and HiSeq2500 data (c). GC content distribution of
22 the BGISEQ-500 and HiSeq2500 data (d). FastQC [18] was used for the calculation.

23
24 **Figure 4. Variation calling based on the dataset.** The major steps included data
25 filtering, alignment and variation calling, and the major parameters are also indicated.

1 **Tables**

2 **Table 1. Summary of the dataset*.**

Sequencing Type	Read ($\times 10^6$)	Bases (Gbp)	GC content	>Q20	>Q30
PE50	2,379	118.94	41.62%	96.00%	87.02%
PE100	1,159	115.88	41.28%	96.39%	87.13%

3 *This dataset was from two runs of the BGISEQ-500 sequencer (PE50 and PE100).
 4 ‘>Q20/Q30 percentage’ indicates the percent of bases with quality score ($-10 \times \lg(\text{error rate})$)
 5 higher than 20 and 30 (indicating error rates of 1% and 1‰ respectively).

7 **Table 2. Mapping statistics of the dataset*.**

Metrics	BGISEQ-500	BGISEQ-500	HiSeq2500
	PE50	PE100	PE150
Clean reads	2,378,725,921	1,136,008,901	708,941,148
Clean bases (bp)	118,936,296,050	113,600,890,100	104,923,289,904
Mapping rate	97.87%	99.22%	99.05%
Unique rate	93.17%	96.47%	97.06%
Duplicate rate	6.26%	2.47%	1.52%
Mismatch rate	0.34%	0.58%	0.56%
Average sequencing depth	37.57	37.44	34.52
Coverage	99.28%	99.12%	99.06%
Coverage at least 4×	98.90%	98.69%	98.60%
Coverage at least 10×	97.97%	97.81%	97.83%
Coverage at least 20×	95.78%	96.06%	94.81%

8 *The statistics shown here are calculated based on the clean reads (raw reads after filtering,
 9 the two platforms’ data went through the same filtering process). Unique mapping rate
 10 indicates proportion of reads with unique alignment in the genome.

1

2 **Table 3. Variation statistics of the dataset*.**

	BGISEQ-500	BGISEQ-500	HiSeq2500
	PE50	PE100	PE150
SNPs	3,451,124	3,477,642	3,609,606
1000genome and dbSNP	3,242,083	3,288,653	3,347,441
1000genome specific	1,260	420	693
dbSNP specific	180,935	179,967	243,256
dbSNP rate	99.19%	99.74%	99.48%
Novel	26,846	8,602	18,216
Homozygous	1,426,328	1,433,490	1,472,063
Heterzygous	2,024,796	2,044,152	2,137,543
Synonymous	19,880	20,012	20,860
Ti/Tv	2.0462	2.065	2.0427
dbSNP Ti/Tv	2.0608	2.0693	2.0503
Novel Ti/Tv	0.8948	0.9775	1.0544
Indels	553,842	842,058	733,797
1000genome and dbSNP	260,157	320,741	314,161
1000genome specific	7,007	22,919	20,049
dbSNP specific	211,846	326,984	285,834
dbSNP rate	85.22%	76.92%	81.77%
Novel	74,832	171,414	113,753
Homozygous	206,163	295,492	300,013
Heterzygous	347,679	546,566	433,784

3 *1000genome and dbSNP equals the number of SNPs that are found in both 1000 genome
4 and dbSNP databases (version 147 was used), 1000genome specific equals the number of
5 SNPs that are only found in 1000 genomes database. dbSNP rate equals the number of SNPs
6 found in dbSNP database/total detected SNPs. Novel SNP equals the number of SNPs that are

1 not found in SNP database. Ti/Tv equals the ratio of SNP type are transition/SNP type are
 2 transversion.

3 **Table 4. Performances of variation calling of dataset*.**

Variant type	Metrics	BGISEQ-500 PE50	BGISEQ-500 PE100	HiSeq2500 PE150
SNPs	True Positive	3,006,132	3,071,579	3,084,449
	False Positive	15,203	6,907	4,318
	False Negative	186,825	121,379	108,508
	Precision	99.50%	99.78%	99.86%
	Sensitivity	94.15%	96.20%	96.60%
	FPR	0.00060%	0.00020%	0.00017%
	FNR	5.85%	3.80%	3.40%
indels	True Positive	261,867	326,810	355,728
	False Positive	16,931	22,246	7,981
	False Negative	107,311	42,391	13,751
	Precision	93.93%	93.63%	97.81%
	Sensitivity	70.93%	88.52%	96.28%
	FPR	0.00067%	0.00069%	0.00032%
	FNR	29.7%	11.48%	3.72%

4 *Above first four metrics are calculated by using rtg-tools software. True Positive (TP) is the
 5 number of SNPs that are found in high-confidence reference dataset, False Positive (FP) is the
 6 number of SNPs that are not found in reference dataset, False Negative (FN) is the number of
 7 SNPs that are found in high-confidence reference dataset but are not found in reference
 8 dataset. Precision is $TP/(TP+FP)*100$. Sensitivity is $TP/(TP+FN)*100$. FPR is $FP/(all$
 9 $high-confident\ region\ length-TP-FN)*100$, where all high-confident region length equals
 10 252,9164,928bp that comes from GIAB released high confidence variants datasets [19]. FNR
 11 is $FN/(FN+TP)*100$.

Figure 1
a

1 ug Genomic DNA
≥ 23kb



1. Fragmentation



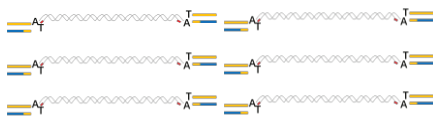
2. Size selection



3. End repair and A-tailing



4. Adaptor ligation and PCR



5. Splint circularization



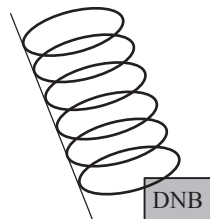
Click here to download Figure 1.pdf 

b

Single strand
circle DNA



1. Making DNBs



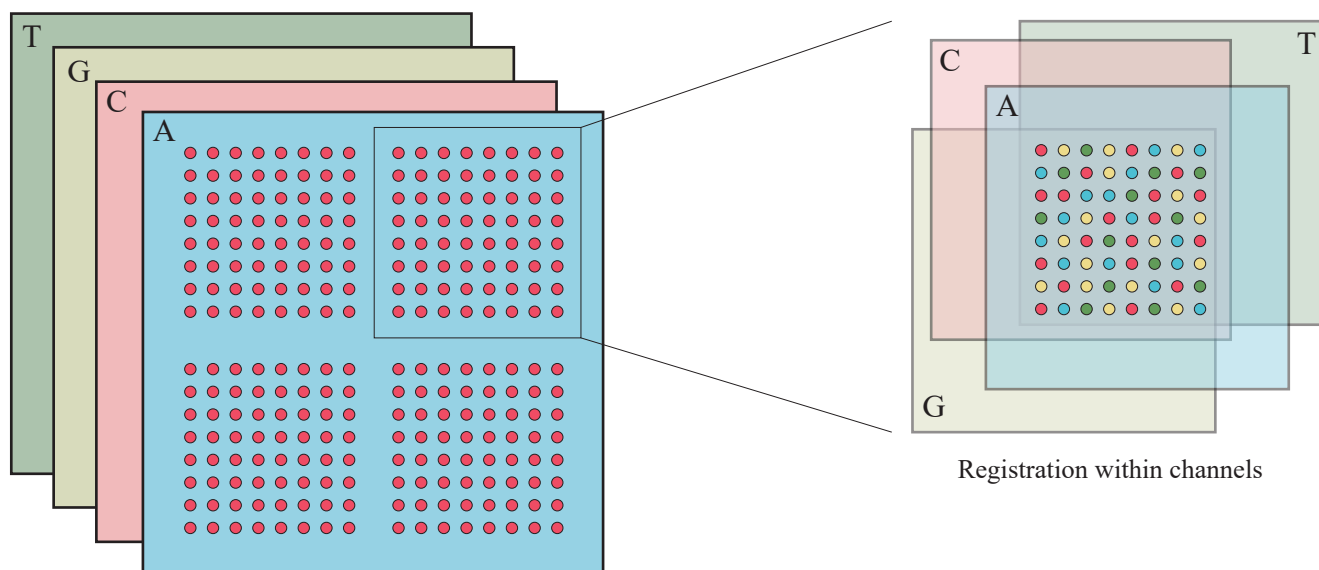
2. Loading DNBs



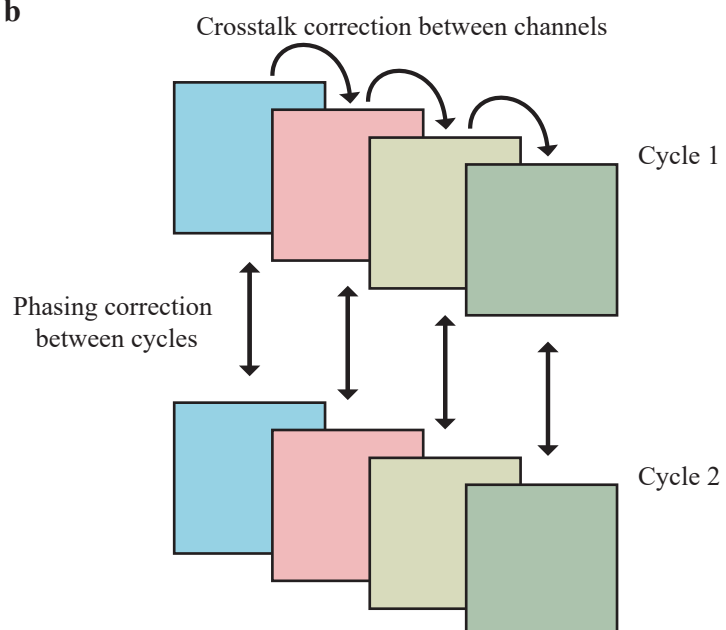
3. Sequencing



a



b



c

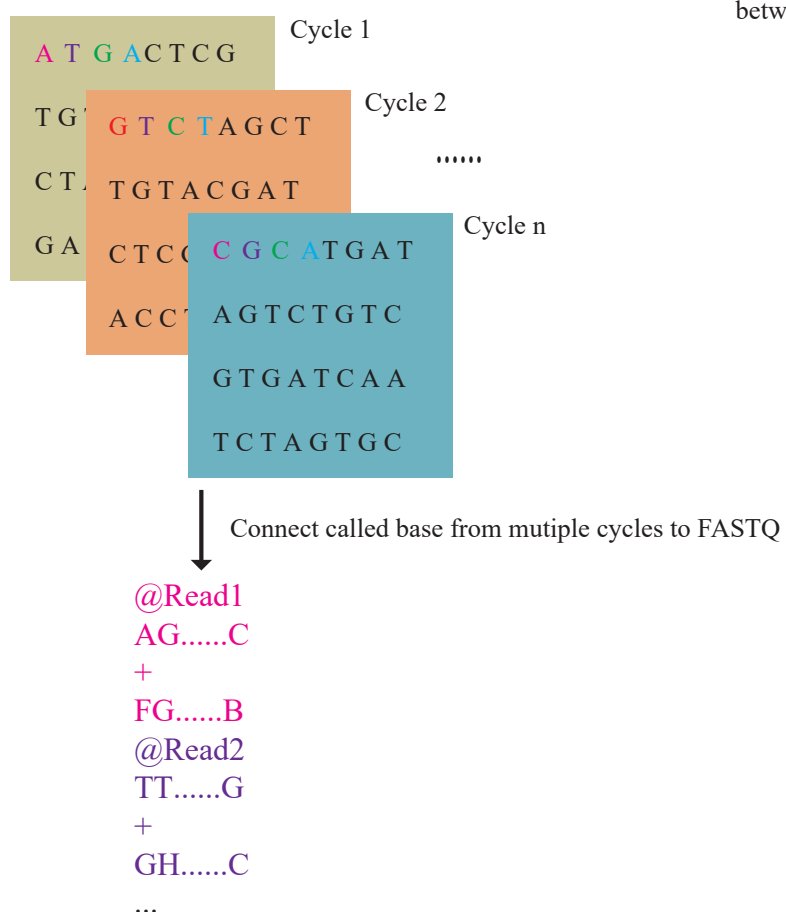
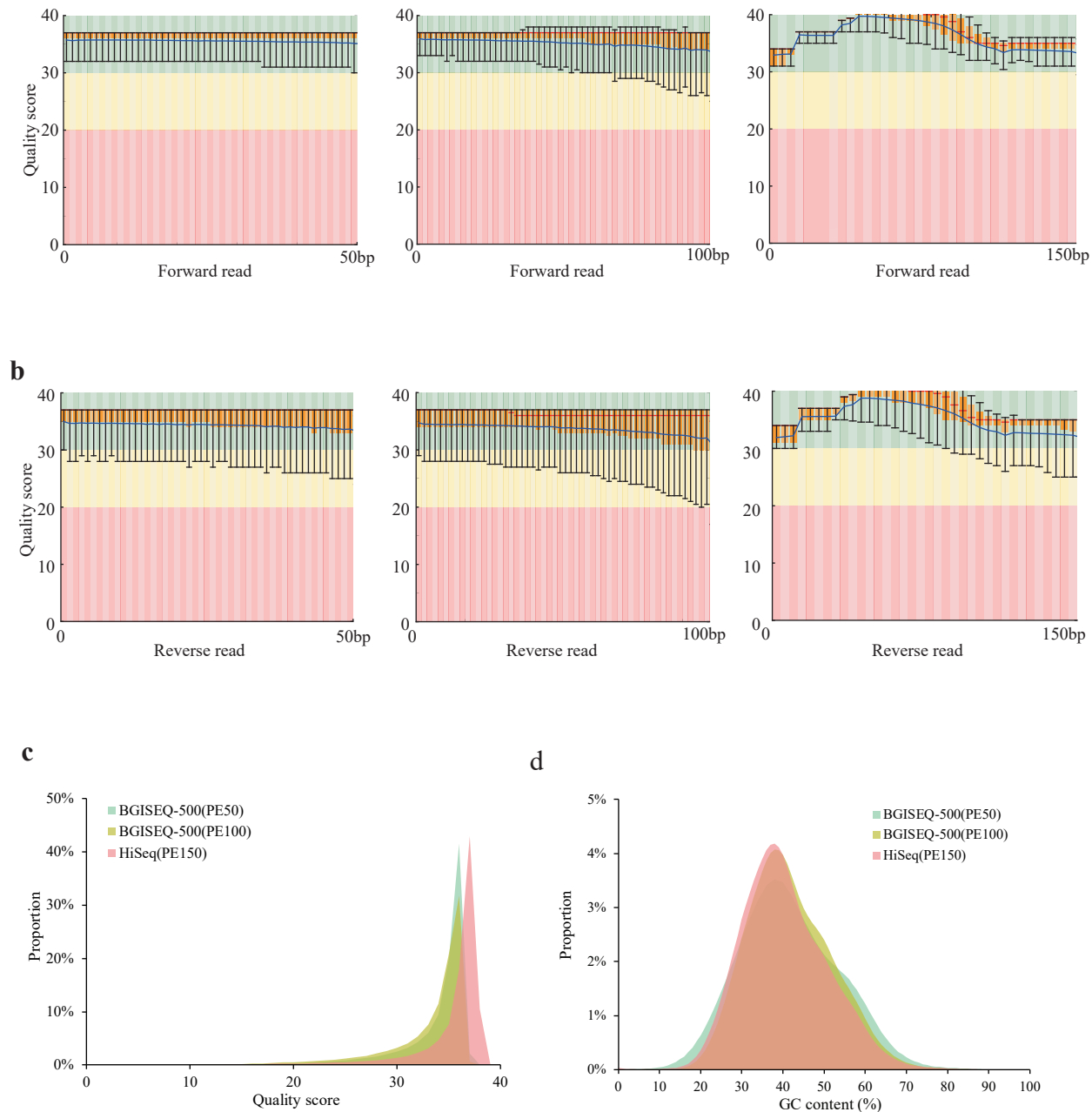
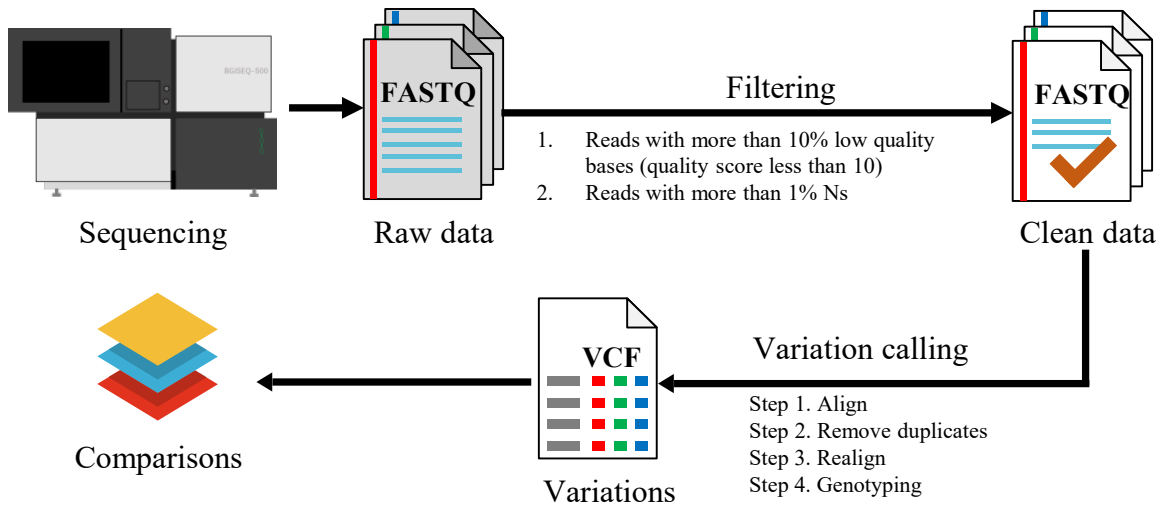
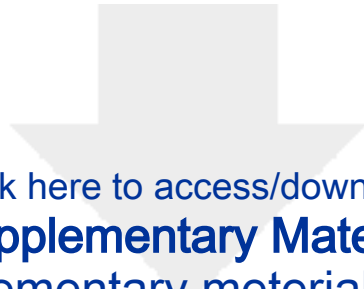


Figure 3

[Click here to download Figure 3.pdf](#)





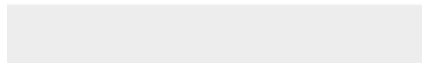


Click here to access/download
Supplementary Material
Supplementary materials.docx





Click here to access/download
Supplementary Material
RebutallLetter.docx



Dear Editor,

Thanks for potentially accepting our manuscript entitled *A reference human genome dataset of the BGISEQ-500 sequencer* for possible publication as a datanote in *GigaScience*.

We have revised our manuscript again according to the suggestions of the reviewer #2.

During the revision, we revised the mistakes in grammar and sentences as pointed out by the reviewer, and we also added the variation calling performance assessment of different data filtering criteria to the supplementary material according to the reviewer's advice. Thus we have addressed all the questions raised by Reviewer #2 as you can find out in the rebuttal letter with point-to-point response letter.

The revised manuscript has ~2,400 words in the main text (including the abstract), with four figures and four tables, and one supplementary file (with detailed parameters for the analysis, the accuracy assessment in different genome content as well as performances of two platforms under different filtering threshold). Please let us know if other information needed, and we are looking forward to your response.

Xin Liu,

BGI Research