

1 **Advanced Whole Genome Sequencing Using a Complete PCR-free Massively Parallel**

2 **Sequencing (MPS) Workflow**

3 Hanjie Shen<sup>1,2,3#</sup>, Pengjuan Liu<sup>1,2,3#</sup>, Zhanqing Li<sup>1,2,3#</sup>, Fang Chen<sup>1,2,3</sup>, Hui Jiang<sup>1,2,3</sup>,

4 Shiming Shi<sup>1</sup>, Yang Xi<sup>1,2,3</sup>, Qiaoling Li<sup>1,2,3</sup>, Xiaojue Wang<sup>2,3</sup>, Jing Zhao<sup>1,2,3</sup>, Xinming

5 Liang<sup>1</sup>, Yinlong Xie<sup>1</sup>, LinWang<sup>4</sup>, Wenlan Tian<sup>4</sup>, Tam Berntsen<sup>4</sup>, Yinling Luo<sup>1,2,3</sup>, Meihua

6 Gong<sup>1,2,3</sup>, Jiguang Li<sup>1,2,3</sup>, Chongjun Xu<sup>1,2,3,4</sup>, Sijie Dai<sup>1</sup>, Zilan Mi<sup>2,3</sup>, Han Ren<sup>2,3</sup>, Zhe Lin<sup>2</sup>,

7 Ao Chen<sup>2,3</sup>, Wenwei Zhang<sup>2,3</sup>, Feng Mu<sup>1</sup>, Xun Xu<sup>2,3</sup>, Xia Zhao<sup>1,2,3\*</sup>, Yuan Jiang<sup>2,3,4 \*</sup>

8 Radoje Drmanac<sup>1,2,3,4\*</sup>

9 1. MGI, BGI-Shenzhen, Shenzhen 518083, China

10 2. BGI-Shenzhen, Shenzhen 518083, China

11 3. China National Genebank, BGI-Shenzhen, Shenzhen 518120, China

12 4. Complete Genomics Inc., 2904 Orchard Pkwy, San Jose, CA 95134 USA

13 # These authors contributed equally to this work.

14 \*Correspondence:

15 Xia Zhao: [zhaoxia@genomics.cn](mailto:zhaoxia@genomics.cn), Yuan Jiang: [yjiang@completegenomics.com](mailto:yjiang@completegenomics.com), and

16 Radoje Drmanac: [rade@completegenomics.com](mailto:rade@completegenomics.com)

## 1 **Abstract**

2 Systematic errors could be introduced by amplification during MPS library preparation  
3 and cluster/array formation. Polymerase Chain Reaction (PCR)-free library preparation  
4 methods have previously demonstrated improved sequencing quality with PCR-amplified  
5 read-clusters, however we hypothesized that some some InDel errors are still introduced  
6 by the remaining PCR step. Here we sequenced PCR-free libraries on MGI's PCR-free  
7 DNBSEQ™ arrays to obtain for the first time a true PCR-free WGS (Whole Genome  
8 Sequencing). We used MGI's PCR-free WGS library preparation kits as recommended  
9 or with some modifications to make several NA12878 libraries. Reproducibly high quality  
10 libraries were obtained with low bias and less than 1% read duplication for both  
11 ultrasonic and enzymatic DNA fragmenting. In a triplicate analysis, over 99% SNPs and  
12 about 98% indels in each library were found in at least one of the other two libraries.  
13 Using machine learning (ML) methods (DeepVariant or DNAscope), variant calling  
14 performance (SNPs F-measure > 99.94%, InDels F-measure > 99.6%) exceeded the  
15 widely accepted standards. The F-measure of 15X PCR-free ML-WGS was comparable  
16 to or even better than 30X PCR WGS analyzed with GATK. Furthermore, PCR-free  
17 WGS libraries sequenced on PCR-free DNBSEQ™ platform have up to 55% less InDel  
18 errors compared to NovaSeq platform confirming that DNA clusters have PCR-  
19 generated errors. Enabled by the new PCR-free library kits, super high-throughput  
20 sequencer and ML-based variant calling, DNBSEQ™ true PCR-free WGS provides a  
21 powerful solution to improve accuracy while reducing cost and analysis time to facilitate  
22 future precision medicine, cohort studies and large population genome project.

1 **Keywords**

2 WGS, PCR-free, DNBSEQ™, Machine learning based variant calling

## 1 Introduction

2 MPS (also known as next-generation sequencing, NGS) technology has revolutionized  
3 basic biology and precision medicine during the past decade. Comparing to partially  
4 sequencing the genome via targeted panels such as whole-exome sequencing (WES),  
5 an entire genome sequencing, also known as whole genome sequencing (WGS),  
6 achieves better coverage uniformity and more reliable detection in single-nucleotide  
7 polymorphisms (SNPs), insertions and deletions (InDels), structural variants (SVs) and  
8 copy number variants (CNVs) at both coding and non-coding regions. There is an  
9 increasing clinical demand to apply WGS as one single test, especially in conditions  
10 when WES could potentially fail to detect all pathogenic variants in a large fraction of  
11 mendelian disorder cases [1,2,3] A variety of studies have already demonstrated the  
12 feasibility of WGS to investigate rare diseases such as inherited retinal disease [4],  
13 cancers such as liver cancer [5] and infectious diseases such as Mycobacterium  
14 tuberculosis (MTB) [6]. More importantly, according to the data from the NHGRI Genome  
15 Sequencing Program, sequencing cost has rapidly dropped to \$0.012 per Mb in July  
16 2017 [7]. WGS costs no more than thousand dollars nowadays, with faster turnaround

1 time and at greater depth, making it more economic to conduct large-scale projects such  
2 as the Genomics England 100,000 Genome Project.  
3  
4 Standard WGS workflow includes MPS library preparation, array/cluster generation, on  
5 chip sequencing, reads filtering, mapping, and variant calling (aka secondary analysis).  
6 Many efforts have been undertaken to further reduce the cost and turnaround time while  
7 improving the WGS data performance. For example, it is well known that PCR  
8 amplification in regular WGS library preparation causes uneven amplification and copy  
9 errors, resulting in coverage bias, GC bias and nucleotide misincorporation [8,9,10,11].  
10 To address this issue, an optimized WGS library protocol that eliminates the need for  
11 PCR called PCR-free WGS is developed. The PCR-free WGS improves read mapping  
12 quality and aids in *de novo* assembly of genomic regions containing extreme GC content  
13 (the percentage of G and C bases) [12]. Moreover, PCR-free protocol removes  
14 duplicated reads from library prep, which increases the read utility and variant calling  
15 confidence [13]. The last benefit of excluding PCR step in the WGS library prep is  
16 obviously shorter turnaround time and lower cost.

1

2 In addition to the new library construction chemistry, different innovative bioinformatics  
3 tools are applied to expedite data analysis without sacrificing accuracy. The steps of a  
4 standard analysis workflow are typically as following: 1) to trim and filter read data; 2) to  
5 align raw data to a reference genome, 3) to call germline or somatic variants, 4) to  
6 conduct tertiary analysis and generate report. The best Practices pipeline for germline  
7 short variant discovery developed at the Broad Institute with the Genome Analysis  
8 Toolkit (GATK) [14] is currently industry standard pipeline for variant calling on WGS.  
9 However, traditional GATK pipeline takes around two days for whole genome data  
10 processing on a standard 24 thread machine [15]. To achieve better detection accuracy  
11 and account for systematic errors in the WGS workflow, researchers have explored  
12 machine learning or deep learning based algorithms and developed several new  
13 analysis pipelines such as GATK CNNScoreVariants from the Broad Institute [16],  
14 Deepvariant from the Google Brain team [17], DNAscope from Sentieon [18] and  
15 Clairvoyante from R. Luo et al [19].

16

1 In this study, we presented a completely PCR-free MPS workflow by constructing PCR-  
2 free WGS libraries and sequencing on PCR-free on-chip arrays. Both PCR-free WGS  
3 sets (mechanic fragmentation with MGIEasy PCR-Free DNA Library Prep Set and  
4 enzymatic fragmentation with MGIEasy FS PCR-Free DNA Library Prep Set)  
5 demonstrated highly reproducible data quality even with 200ng or 50ng genomic DNA as  
6 input. More importantly, a significant improvement was achieved in InDels calling with  
7 GATK from an average F-score of 95.43% in three PCR libraries to that of 99.05% in  
8 three PCR-free WGS libraries. By incorporating machine learning based algorithms, the  
9 F-score of 15x PCR-free analyzed with DNAscope or DeepVariant can outperform that  
10 of 30x PCR analyzed with GATK in some scenarios. The complete PCR-free WGS  
11 significantly decreased FP and FN InDel reads, leading to a better InDel calling  
12 accuracy (99.3-99.6% F-score) than that of the Illumina's PCR-free libraries (98.0-99.3%)  
13 with lower genomic DNA as input. Additionally, an incredibly low duplication rate, less  
14 than 1%, was achieved in the DNBSEQ™ PCR-free WGS, resulting in better read utility  
15 and further reducing the sequencing cost for a standard 30x coverage. In summary, the  
16 advanced PCR-free WGS reported herein could lead to wider adoption of WGS in the  
17 genomic research and gradually in the clinical practice for rare disease urgent diagnosis.

1

## 2 **Results**

### 3 **A complete PCR-free MPS workflow**

4 PCR is frequently used to increase template amount during MPS library construction. It's  
5 also an essential step of 'bridge amplification' to generate identical copies on a flow cell  
6 surface [20]. Here, we describe an MPS workflow called DNBSEQ™ PCR-free which  
7 completely eliminates PCR amplification in both library and array preparations (Figure  
8 1a). DNBSEQ™ PCR-free starts with DNA fragmentation, followed by size selection  
9 using solid-phase reversible immobilization (SPRI) beads. A single-tube protocol is  
10 applied to conduct sequential multiple enzymatic reactions to attach a barcoded adapter  
11 to the DNA of interests. After removing excess adapters, single-stranded circle (ssCir) is  
12 formed and served as the template in rolling circle amplification (RCA) for DNB  
13 preparation. DNBs are then loaded onto the patterned flow slides and sequenced. In  
14 contrast to bridge amplification, RCA is a linear amplification from the original ssCir  
15 template and therefore errors would not accumulate[21,22]. Different from DNBSEQ™  
16 PCR with additional amplification after adapter ligation and bridge PCR or emulsion PCR



1 based dsDNA library sequencing, DNBSEQ™ PCR-free can completely avoid the  
2 enrichment of PCR errors in template amplification and library cloning process, and  
3 faithfully restores the original landscape of the genome.

4

5 Two sets from MGI (MGIEasy PCR-Free DNA Library Pre Set and MGIEasy FS PCR-  
6 Free DNA Library Pre Set) are used to prepare the DNBSEQ™ PCR-free libraries. The  
7 PCR-Free set is used with acoustic fragmented samples, while the FS PCR-Free one  
8 includes sequential reactions of enzymatic fragmentation, end repair/dA-tailing and  
9 adapter addition in a single tube. We compared the performance of both sets with nine  
10 libraries constructed from 1 µg NA12878 reference genomic DNA and sequenced with  
11 pair-end 150 bp read length. Summarized in figure 1b were the QC statistics including  
12 GC Content, Duplication Rate, Median Insert Size, and regions with > 10X coverage.  
13 Both sets showed highly reproducible performance with PCR-Free set having slightly  
14 smaller variation than FS PCR-Free set. We also tested different input amounts (1 µg,  
15 500 ng, 200 ng, and 50ng only for FS PCR-Free Set) with both sets and observed  
16 comparable performance (Figure 1S).

1

## 2 **Minimal GC bias for genomes with different GC contents**

3 The relationship between GC content and read coverage across a genome, namely GC  
4 bias, can be greatly affected by MPS library prep, cluster/array amplification, and  
5 sequencing. To evaluate the performance of the DNBSEQ™ PCR-free MPS workflow on  
6 GC bias, DNA samples from bacteria strains with GC contents of 38% and 62% were  
7 processed with both sets mentioned above. Libraries were made following by kit  
8 instruction and sequenced on MGISEQ-2000 with paired-end sequencing (2x150 bp).  
9 From Figure 2, the number of reads covering different regions was normalized by the  
10 mean read coverage and was then plotted against genomic GC contents of  
11 corresponding regions. Normalized coverage lower or higher than 1 indicates certain GC  
12 bias. The high GC and low GC genomes demonstrated fairly even coverage of reads  
13 across the genome with either acoustic shearing or enzymatic shearing. Overall, the  
14 DNBSEQ™ PCR-free MPS workflow demonstrated minimal GC bias in genomes with a  
15 variety of GC contents.

16

1 **Low duplicate rate and high variant calling F-score for GIAB with DNBSEQ™ PCR-free**

2 We compared the sequencing accuracy of MPS libraries with or without PCR using  
3 DNBSEQ™ sequencing technology [23]. After acoustic fragmentation, NA12878 DNA  
4 was processed following instructions of MGIEasy PCR-Free DNA Library Prep Set to  
5 construct three PCR-free WGS libraries. Meanwhile three PCR WGS libraries were  
6 prepared from the same DNA sample. Each library was sequenced individually on one  
7 lane of MGISEQ-2000 with paired-end sequencing (2x150 bp). (Methods and Materials)  
8 The total raw data of each lane was greater than 120G with GC content ranging from  
9 41.46% to 41.78% ( Table 1 ) .

10

11 The raw reads were down sampled from original full lane (approximately 46x) to create  
12 additional 30x and 15x depth dataset. After raw reads filtering, clean reads were aligned  
13 to the human reference genome with decoy sequence, hs37d5. The mapping quality  
14 matrix was summarized in Table 1. We compared all 3 depth datasets from the 6  
15 libraries (18 dataset in total). In 30x depth data, WGS PCR and PCR-free libraries had  
16 similarly high mapping rate at 99.8% and overall coverage higher than 99.1%. PCR-free

1 libraries showed slightly lower duplication rate and mismatch rate, around 1% and 0.7%  
2 respectively on average, while they were 1.5% and 0.9% in PCR libraries. Other depth  
3 dataset showed similar patterns. Theoretically, because of the “true PCR-free” workflow  
4 duplicate rate should be zero, but there are chances that the same DNB being read  
5 twice or multiple times caused by optical overflow ( aka “optical duplicates” ) or the  
6 same regions from different genome copy produced in DNA fragmentation step ( aka  
7 “natural duplicates” ) being sequenced and being marked as duplicate by QC tool  
8 incorrectly [24,25].

9

10 We next used three variant callers: GATK, DeepVariant, and DNAscope (Methods and  
11 Materials), to assess the accuracy of PCR or PCR-free methods (Table2, Table1S, and  
12 Figure 3). The GATK HaplotypeCaller has become the industry standard small variant  
13 caller due to its high accuracy and has achieved top performance in a variety of public  
14 and third-party benchmarks [15,26,27,28], while DeepVariant and DNAscope are two  
15 newly developed variant callers based on machine learning method [17,18,29,30]. It

1 should be noted that both machine learning variant callers used in this study were  
2 optimized for the DNBSEQ™ platform through the use of utilizing in-house training data.  
3  
4 Variant calling matrix is highly reproducible in all three replicates for both PCR and PCR-  
5 free libraries (Table1S). Table2 summarized the average number of three replicates for  
6 different depths. At 30x depth, GATK called and marked Passed Filter (named true  
7 positive, or TP for short) for an average of a total of 3,651,696 TP variants for PCR  
8 libraries, and 3,674,252 TP variants for PCR-free libraries, while DeepVariant and  
9 DNAscope showed higher sensitivity as they both detected additional TP variants for  
10 both PCR and PCR-free libraries. From all three callers, PCR-free libraries  
11 demonstrated slightly reduction in False Positive (FP) numbers of SNPs (from 2806 to  
12 2586 in GATK, from 2913 to 2111 in DeepVariant, from 1565 to 1254 in DNAscope) as  
13 well as reduction in False Negative (FN) numbers of SNPs (from 16671 to 10926 in  
14 GATK, from 4266 to 3109 in DeepVariant, from 4355 to 3202 in DNAscope), and  
15 dramatic reduction in FP InDels (from 20699 to 2766 in GATK, from 8008 to 2124 in  
16 DeepVariant, from 8632 to 1620 in DNAscope) as well as in FN InDels (from 23154 to

1 6345 in GATK, from 13082 to 3690 in DeepVariant, from 14301 to 3299 in DNAscope).

2 This leads to slight increase of SNP F-score (harmonic mean of recall and precision) and

3 significant increase of InDels F-score for PCR-free libraries, suggesting more precise

4 variant calling for all depths (Figure 3a). As the highest accuracy combination, PCR-free

5 data with DNAscope had the lowest FP SNPs, FP InDels, and FN InDels. As an

6 additional evaluation, in selected “difficult” genome regions such as repeat regions,

7 extreme GC regions, PCR-free generally showed better InDels F-score than PCR

8 libraries, that output more faithful genome sequences for applications (Figure 2S).

9 To increase the confidence in the performance of PCR-free WGS workflow across a

10 variety of sequencing depth, we generated additional 10x and 20x data point to conduct

11 a 10x-46x low depth test. As being demonstrated in Figure 3B, the reduction in coverage

12 (i.e. 10x, 15x and 20x) clearly affected the quality of variant calling from all methods.

13 Nevertheless, PCR-free method coupled with machine learning-based callers produced

14 more accurate calling than pipelines involving PCR library construction or GATK caller.

15 Importantly, the SNP F-scores of DNAscope and DeepVariant for 15x PCR-free libraries

16 were comparable to GATK for 30x PCR libraries (99.45%, 99.48% vs 99.70%), while

17 InDel F-scores of DNAscope and DeepVariant for 15x PCR-free libraries were

1 significantly higher than GATK for 30x PCR libraries (97.90%, 98.23% vs 95.43%),  
2 indicating the potential application to lower sequencing cost while maintain the variant  
3 detection accuracy level. Noticeably, here DeepVariant for 15X PCR-free libraries  
4 showed the highest accuracy among all callers.

5

#### 6 **Reproducibility of PCR-free libraries**

7 We also conducted consistency analysis, in order to determine the level of randomness  
8 introduced in library construction and sequencing step, and whether variant callers can  
9 help correct them. As a result, variant consistency of three replicates in PCR-free  
10 libraries was observed better than PCR libraries, especially for InDel consistency and  
11 the trend is similar across the three different variant callers (Figure 4). This is expected  
12 since PCR step inevitably introduces random errors during amplification. InDel  
13 consistency (represented by portion of variants shared by all replicates) of PCR-free  
14 libraries were found at GATK 84.2%, DeepVariant 86.5%, DNAscope 89.1% (three  
15 callers in average 86.6%), which were almost 20% higher than that of PCR libraries with  
16 GATK 63.2%, DeepVariant 66.9%, DNAscope 68.9% (three callers in average 66.3%).

1 We also observed more than 99% SNPs and approximately 98% InDels in each library  
2 overlapping with at least one of the other two libraries. Consistency of SNPs among all  
3 three replicates were similarly high from both PCR and PCR-free libraries (three callers  
4 in average 94.6% vs 95.2%).

5

6 On the other hand, we investigated the consensus calling on the same library among  
7 three callers. In average, the three pipelines for SNPs and InDels are 92.3% and 67.0%  
8 for PCR libraries as well as 92.6% and 82.7% for PCR-free libraries, indicating that  
9 PCR-free libraries generated more “clear” variants candidates that are less challenging  
10 for variant callers (Figure 3S).

#### 11 **Evaluation of PCR-free WGS performance on different sequencing platforms**

12 In addition to library construction kit itself, sequencing platform also brings in additional  
13 bias or systematic variation that causes different performance. Here we compared two  
14 PCR-free libraries made by both MGIEasy PCR-Free DNA Library Pre Set and MGIEasy  
15 FS PCR-Free DNA Library Pre Set and sequenced on MGISEQ-2000, with three  
16 datasets downloaded from Illumina Basespace website to represent performance of



1 TruSeq PCR-free libraries sequenced on HiSeq4000, HiSeqXTen, or Novaseq platforms.  
2 Included in the comparison were three additional datasets in order to provide further  
3 information: 1) library prepared with MGIEasy PCR-Free DNA Library Pre Set and  
4 sequenced on DNBSEQ-T7 [38]; 2) library prepared with MGIEasy PCR-Free DNA  
5 Library Pre Set with illumina's adapter and sequenced on Novaseq by a third party  
6 sequencing service provider; 3) library prepared with a research protocol that has yet  
7 been incorporated into the MGIEasy kit and sequenced on MGISEQ-2000. It should be  
8 noted that MGIEasy kit prepared libraries used 1 µg or 250ng DNA input, far less than  
9 the input amount of datasets downloaded from Illumina Basespace website. All FASTQ  
10 files were processed in a same pipeline for reads trimming/filtering and mapping, as well  
11 as variant calling using DNAscope pipeline (Methods and Materials), to minimize bias or  
12 variation introduced at secondary analysis stage.

13

14 From mapping QC matrix, all three datasets generated from MGISEQ-2000 showed  
15 significant lower duplicate rate at around or less than 1%, comparing to that from all  
16 Illumina's platforms with at least 10% duplicate rate, including the hybrid sample with

1 MGI library prep (Table 3). The “T7”dataset had 3.60% duplication rate with 250ng  
2 genomic DNA as input, still lower than ILMN's duplicate rate. This resulted in a more  
3 cost-effective read utility from DNBSEQ™. Meanwhile, with DNAscope's DNBSEQ™ and  
4 Illumina models applied separately, SNP calling accuracy (represented by F-score) of all  
5 samples reached a similarly high level, except for Hiseq4000 (Table 3). For InDels, the  
6 two pure MGI pipeline (library prep + sequencing) generated datasets and the “research  
7 library dataset” all outperformed the three ILMN datasets, indicating less errors (FP and  
8 FN) generated due to the completely elimination of PCR during both library construction  
9 and sequencing procedures. Comparing the top performer from both sides: the  
10 “research library dataset” and the “TruSeq and Novaseq dataset”, the “research library  
11 dataset” had similar SNP calling but significant lower false InDels calling - 1304 in FP,  
12 making a ~55% reduction from 2879 on “TruSeq and Novaseq dataset”, and 2201 FN as  
13 ~48% reduction as well. The inclusion of data from new sequencer DNBSEQ-T7  
14 provided us with an opportunity to place its accuracy performance among the sequencer  
15 matrix here. Although prioritizing scale and cost in designing this ultra-high throughput  
16 sequencer, its sequencing accuracy was not comprimized, still reached the level of  
17 “TruSeq and Novaseq dataset”.

1 These data collectively highlighted that the cost as well as data quality such as  
2 duplication rate InDel calling accuracy could be tremendously benefited from the  
3 complete DNBSEQ™ PCR-free WGS workflow.

4

## 5 **Discussion**

6 The widely adoption of WGS gives the credit to dramatic decreasing of the cost of  
7 sequencing, as well as the potential downside of WES, such as missing variants in the  
8 non-exome regions and incapable of detecting copy number variation. The development  
9 of machine-learning based secondary analysis tools also promote the usage of WGS, by  
10 shortening the analysis time from days to hours while reaching a higher variant calling  
11 accuracy.

12

### 13 **Advantages of DNBSEQ™ PCR-free**

14 In this study, we performed the PCR-free library prepare combined with DNBSEQ™  
15 excluding PCR from the entire sequencing process. This leads to a good performance of  
16 GC bias and coverage. We showed the PCR-free libraries had more accurate small

1 variant especially InDel calls at 30x and other sequencing depth dataset. This result  
2 confirmed a previous report that the PCR-free library had high quality of InDel calls [31].  
3 And this advantage expanded when the machine learning variant callers, such as  
4 DeepVariant and DNAscope, is applied to replace GATK. The fact that InDel F-score of  
5 PCR-free library called by both machine learning callers at 15x depth is significantly  
6 better than that of PCR library called by GATK at 30x implies a potential to use low  
7 depth sequencing data to conduct variant calling while not compromising any accuracy.  
8 Furthermore, true PCR-free WGS provided by good quality PCR-free libraries and  
9 DNBSEQ™ array is expected to improve detection of low frequency somatic mutations.

10

11 It is known that duplicate rate represents the proportion of the duplicate reads from all  
12 the sequenced data. In order to ensure the accuracy, the duplicated reads need to be  
13 removed for the following bioinformatics analysis. Therefore, for the same amount of raw  
14 data, lower duplicate rate means more usable clean data. The average duplicate rate of  
15 our PCR-free data on MGISEQ-2000 is about 1% , which is much lower than the PCR-  
16 free libraries (10.43%-13.62%) from the Illumina platforms. So for 30x WGS , 90-100G

1 raw base is normally enough on MGISEQ-2000 , however, about 20% more raw base  
2 ( 110-120G ) is requested on Illumina platforms. The superiority of MGISEQ-2000  
3 duplicate rate is due to the “true PCR-free” process, which strictly has no PCR step both  
4 in the library construction and sequencing workflow.

5

### 6 **Rapid WGS for infant genetic disorders diagnostic**

7 It was reported that genetic disorders are among the top cause of morbidity and mortality  
8 in infants, and the newly developed Rapid whole-genome sequencing (rWGS) shows the  
9 power to diagnose genetic disorders in a timely manner enabling healthcare providers to  
10 generate or change management plan accordingly and thus improves outcomes in  
11 acutely ill infants [32]. To facilitate this clinical utility, it is essential for the whole process  
12 from sample collection to diagnostic report review and signing can be complete within 2  
13 days. As a result, there is urgent need for pushing turn-around time of library  
14 construction, sequencing, secondary analysis down to 24 hours, to open up some space  
15 for candidate variant clinical annotation and board review.

1 Compared to traditional PCR involved library construction method, the PCR-free method  
2 skipped amplification and therefore significantly reduced the turn-around-time. Moreover,  
3 the enzymatic shearing method makes it more feasible for automation and further save  
4 the library construction hands-on time. The new sequencer DNBSEQ-T7 with super high-  
5 thoupout up to 6T per run can finish within 24h, which greatly shorten the sequencing time  
6 compared with MGISEQ-2000 and provides the possibility to realize rWGS. In the  
7 secondary analysis part, the most time-consuming steps are reads alignment and variant  
8 calling. While traditional BWA GATK pipeline takes more than desired computational  
9 time, some alternative pipelines have been proposed and developed to meet the  
10 accelerated speed requirement, including CPU based tools such as Strelka2 [33],  
11 DNaseq [34], and SpeedSeq [35], and FGPA&CPU based instrument such as DRAGEN  
12 system [36] and MegaBOLT system [37], as well as the GPU/TPU implemented  
13 Deepvariant [17]. All these advanced bioinformatics tools are compatible with data  
14 generated from MGI's library construction kit and MGI sequencing platforms.

15

16 **Structural Variation (SV) and Copy Number Variation (CNV)**

1 SV and CNV are two critical clinical utilities and for which researchers choose WGS  
2 instead of WES. Although not being designed as primary goal for this benchmark study,  
3 the data generated also allow us to investigate if PCR-free method improves SV and  
4 CNV detection accuracy. SV was called by DNAscope in default parameters and CNV  
5 was called by GATK 4.1.2 pipeline. From the result presented at Figure 4S it can be  
6 seen that PCR-free method likely helps improve both sensitivity and specificity of SV  
7 calling, comparing to PCR involved library construction method, as more SV events were  
8 called in PCR-free group while they reached higher consistency among three replicate  
9 samples. Key point in detecting SV events is correctly detecting breakpoints, which  
10 relies on sufficient coverage across targets and less error to generate false positive.  
11 Obviously, PCR-free libraries will benefit this process.

12 Germline CNV for all these 6 testing samples were called by GATK 4.1.2 and  
13 approximately 2500 CNV events (mainly deletions) were identified from each sample.  
14 When conducting 3-way comparison to analyze the repeatability among replicates, it can  
15 be seen PCR-free group showed slightly higher common called CNV events, but overall  
16 not much difference comparing to PCR group (Fig 5S).

1

## 2 **Clinical Utility brought by higher accuracy**

3 Clinical WGS starts to show its potential in rare disease diagnostic utility, as this  
4 application can quickly cover whole genome area, and finds clinical meaningful variants  
5 especially at UTR or promoter regions when panel or WES fail to cover. The increased  
6 variant detection accuracy from the utility of PCR-free library construction and machine  
7 learning based variant calling pipelines clearly increased WGS variant calling accuracy  
8 and therefore will definitely add value to the diagnostic application. On the other words,  
9 for the region PCR WGS fail to generate enough reads coverage or consistently  
10 generate wrong variant calling the PCR-free WGS will be able to provide with the correct  
11 SNP/InDel information. From the six DNBSEQ™ dataset and three ILMN dataset  
12 evaluated in this study, we found two example genes where all three PCR-free libraries  
13 showed accurate variant callings on the gene coding or UTR regions but the all three  
14 PCR libraries generated FN calling. These two genes were ATK1 and GNAS (Figure 6S,  
15 7S). Similarly, we also showed one example gene “MAF”(Figure 8S), where DNBSEQ™  
16 PCR-free method really stood out, as all 3 ILMN dataset and 2 MGI PCR dataset failed



1 to detect a C to CT insertion. These three genes code clinical meaningful proteins where  
2 a failed variant detection could lead to mis-diagnostic result. For example, gene ATK1  
3 (ATK serine / threonine kinase 1) is associated with multiple clinical phenotypes  
4 including Breast cancer (MIM #114480), Colorectal cancer (MIM #114500), Cowden  
5 syndrome 6 (MIM #615109), and Ovarian cancer (MIM #167000). The T to TC insertion  
6 at locus Chr14:105262025, and CG to GC SNP at Chr14:105262041 may cause  
7 malfunction and introduce the disease phenotype, and only PCR-free libraries are  
8 capable to catch the critical causing variant.

9

## 10 **Future improvement**

11 As a benchmark project, this study managed to show the current performance of tools  
12 and pipelines used in library prep, sequencing, and data analysis. As expected,  
13 emerging techs will continually push forward the upper limit of sequencing accuracy. For  
14 example, performance of new sequencer DNBSEQ-T7 [38] was briefly displayed in this  
15 study, which among other in developing sequencers, can bring down cost per genome or  
16 provide deeper coverage and further shorten the sequencing time when needed.

1 Weakness of the Structural Variation detection brought by the nature of short read  
2 sequencing can be greatly compensated by applying long fragment read (LFR)  
3 barcoding technologies [39,40,41]. For some clinical samples like cfDNA, it's impractical  
4 to obtain 200ng for library preparation. With developed method based on the MGI PCR-  
5 Free sets with pooling sequencing strategy, we successfully generated good data from  
6 10ng cfDNA in some unpublished studies (Data not shown).

7 On the analysis side, current GATK best practice pipeline was developed and tuned  
8 based on Illumina data and did not officially support DNBSEQ™ generated data at least  
9 by the time of December 2018 according to GATK team's response in forum [42]. It is  
10 known that certain error correction steps in the pipeline such as BQSR and VQST was  
11 not developed for DNBSEQ™ data and thus may generate less than optimized result  
12 comparing to applying on Illumina sequencer generated data. Our suspect also comes  
13 from two recent published benchmark studies [43,44], where DNBSEQ™ data analyzed  
14 by Strelka2 and DeepVariant showed comparable accuracy with Illumina data, but  
15 DNBSEQ™ data analyzed by GATK returned a worse accuracy.

16

1 Both DeepVariant and DNAscope rely on the proper model training, which requires  
2 sufficient sample numbers from the library construction method and sequencing platform  
3 same as the testing samples would be generated through. In this study we do not  
4 believe the requirement was fully met, especially for DNAscope. For example, only 30x  
5 and higher depth dataset was included into training set but testing on 15x depth data  
6 was conducted. The negative effect of DNAscope lacking proper 15x depth training  
7 dataset can be obvious when comparing its 10x and 15x accuracy result with  
8 DeepVariant, whose training dataset included 15x depth data. Another point worth  
9 noticing is that using a single model for all library kit/assay and sequencers could  
10 sacrifice accuracy for individual cases. It is desirable for users to train individual model  
11 for each individual case (i.e., combination of sample prep kit, library construction kit, and  
12 sequencing platform), to achieve the optimal variant calling accuracy. With all the above-  
13 mentioned improvements for future WGS cohort studies, the unprecedented data  
14 generation speed and quality would help to answer difficult genetic questions and move  
15 the genomics field into a new era of broad clinical use.

## 1 **Methods/Experimental**

### 2 **DNA preparation**

3 Genomic DNA of NA12878 cell line (RRID:CVCL\_7526) was purchased from the Coriell  
4 Institute. Genomic DNA was quantified by using Qubit™ 3 Fluorometer (Life  
5 Technologies, Paisley, UK). Size and quality of genomic DNA was confirmed by running  
6 0.8% of agarose gel.

7

### 8 **PCR-free library preparation**

9 Acoustic fragmentation PCR-free libraries and enzymatic fragmentation PCR-free  
10 libraries were respectively constructed using MGIEasy PCR-Free DNA Library Prep Set  
11 (MGI, cat. no. 1000013453) and MGIEasy FS PCR-Free DNA Library Prep Set (MGI, cat.  
12 No. 1000013455) respectively.

13 For acoustic PCR-free libraries, genomic DNA was fragmented to 100-600 bp with peak  
14 size at 350-400 bp using LE220 Ultrasonicator (Covaris, Woburn, MA, USA). For FS  
15 PCR-free libraries, genomic DNA was fragmented to 100-1000 bp with peak size at 350-

1 475 bp using enzymatic shearing method. Subsequently, fragmented DNA with a size  
2 range of 200-450bp was selected using MGIEasy DNA Clean Beads (MGI, cat. no.  
3 1000005279) and attached with DNBSEQ™ adapters following the set instruction. We  
4 also followed the protocols from the kit to make single-stranded DNA (ssDNA) circles  
5 and quantified them on Qubit™ 3 Fluorometer (Life Technologies, Carlsbad, CA, USA).

6 The library preparation procedure for research libraries was similar to that in MGIEasy  
7 PCR-Free DNA Library Prep Set, except for size selection and single strand  
8 degeneration method.

9

#### 10 **PCR library preparation**

11 PCR libraries were prepared with the same procedure of DNA fragmentation, end repair  
12 and adapter ligation using MGIEasy PCR-Free DNA Library Prep Set (Cat. no.  
13 1000013453), which was as described above. After adapter ligation, the reaction  
14 product was purified with MGIEasy DNA Clean Beads (MGI, cat. no. 1000005279) and  
15 the ligation products were subject for PCR amplification following instruction from KAPA  
16 HiFi HotStart ReadyMix ( KAPA BIOSYSTEMS , KK2602 ) . A total of 6 cycles (95°C 3  
17 min; 6 cycles, 98°C 20 s, 60°C 15 s, 72°C 60 s; 72°C 10 min; 4°C Forever) in a volume

1 of 100 ul was carried out. After beads purification using MGIEasy DNA Clean Beads and  
2 quantification using Qubit™ 3 Fluorometer, 1 pmol PCR product was taken for the single  
3 strand molecule circularization following the ssCir formation protocol from the MGIEasy  
4 kit.

5

## 6 **Sequencing**

7 Whole genome sequencing was performed on a DNBSEQ™ platforms MGISEQ-2000  
8 for PE150, and DNBSEQ-T7 for PE150. Before sequencing, 75 fmol ssDNA of PCR-  
9 free library or 40fmol single strand circle DNA of PCR library was made into DNB (DNA  
10 Nanoball) following kit instruction from MGISEQ-2000RS High-throughput Sequencing  
11 Set (FCL PE150) (MGI, cat. no. 1000012555 ) . 75 fmol ssDNA of PCR-free library was  
12 made into DNB following kit instruction from DNBSEQ-T7RS High-throughput  
13 Sequencing Set (PE150) (MGI, cat. no. 1000016106). Subsequently, loading DNBs  
14 onto the sequencing chip was performed, and PE150 sequencing was conducted on  
15 DNBSEQ™ platforms using MGISEQ-2000RS or DNBSEQ-T7RS High-throughput  
16 Sequencing Set.

1

## 2 **GC bias analysis**

3 To explore GC bias performance, we sequenced two bacteria samples (Table 2S) on  
4 MGISEQ-2000 for PE150. Filtered reads were aligned to reference genomes by  
5 Burrows-Wheeler aligner (BWA) [45]. To investigate the relationship between GC bias  
6 and read coverage, we scanned genome with a sliding window of the default size (100  
7 bases). GC content and average read coverage were calculated within each window.  
8 Read coverage was normalized to the mean value so that the results would not scale  
9 with the total amount of data. In addition, we eliminated the data points whose coverage  
10 was more than twice the mean coverage because they likely represented repeats.  
11 Finally, we fit the remaining data points by a straight line and defined the slope as the  
12 degree of GC bias in the real data.

13

## 14 **Reads filtering and mapping**

15 As the first step, raw reads sequenced from PCR or PCR-free libraries were debarcoded  
16 by Seqtk [46] with default parameters. Then split reads were pre-processed by

1 SOAPnuke to generate filtered reads [47]. During this filtering process reads containing  
2 more than 10% of 'N' or 50% of the base quality score lower than 12 were removed.  
3 Meanwhile, adapters were trimmed off, which was followed by alignment of all reads  
4 against human reference genome with decoy sequencing hs37d5 (or reference genome  
5 sequences of two bacterias for GC bias analysis) using Burrows-Wheeler aligner (BWA)  
6 with default parameters. The output SAM file was converted to BAM file and sorted by  
7 Samtools [48]. Lastly, duplicates were marked by Picard [49] to make both BAM files  
8 variant calling ready as input for GATK, DeepVariant and DNAscope variant calling  
9 pipeline.

10

## 11 **Running GATK**

12 SNP and InDel calling were performed according to GATK (version 3.3) Best Practice  
13 [50]. Reads around InDel were re-aligned and base quality scores were recalibrated.  
14 HaplotypeCaller was used to call variants in gVCF mode on each chromosome.  
15 Genotyping on the gVCF files was performed by using GenotypeGVCFs with parameters  
16 as following: `-stand_call_conf 30` and `-stand_emit_conf 10`. SNPs and InDels were



1 separated using SelectVariants tool. Variants quality score recalibrates (VQSR) was  
2 carried out to filter low quality variants. SNP annotation --ts\_filter\_level was used for  
3 calculation and filtered at a 99.9% level, while for InDel, --ts\_filter\_level was used for  
4 calculation and filtered at a 99.9% level of the true sensitivity.

5

## 6 **Running DeepVariant**

7 Taking advantages of the state-of-the-art deep-learning technique for image  
8 classification, DeepVariant (V0.7.0 in this study) can achieve a higher accuracy for  
9 bioinformatics analysis. The GIAB (Genome In A Bottle) truth set and corresponding  
10 fastq reads were utilized as training dataset to train a Convolutional Neural Network  
11 (CNN) model. As an alternative for GATK HaplotypeCaller, DeepVariant accepts aligned  
12 reads (e.g. BAM file) as input. In DeepVariant, candidate variants are carefully filtered  
13 along the genome and classified into three genotype states, homozygous reference  
14 (hom-ref), heterozygous (het) or homozygous alternate (hom-alt), with the previously  
15 trained CNN model.

1 To achieve best calling performance, we fine-tuned the CNN model in DeepVariant  
2 using a set of PCR-free data, including 30x and 15x DNBSEQ™ PCR-free sequencing  
3 data of HG001 and HG005 samples. The fine-tuned model is accessible at [51].

4

#### 5 **Running DNAscope**

6 Sentieon DNAscope (version 201808.01 and 201808.05 were used in this study)  
7 uniquely combines the well-established methods from haplotype-based variant callers  
8 with machine learning to achieve improved accuracy. As a successor to GATK  
9 HaplotypeCaller, DNAscope uses an identical logical architecture of active region  
10 detection, local haplotype assembly, and read-likelihood calculation (Pair-HMM) to  
11 produce variant candidates. DNAscope outputs additional informative variant  
12 annotations used by the machine learning model. Annotated variant candidates are then  
13 passed to a machine learning model for variant genotyping resulting in improvements in  
14 both variant calling and genotyping accuracy.

15

1 For this study, DNBSEQ™ model for DNAscope was constructed using publicly available  
2 data from the HG001 and HG005 Genome in a Bottle samples downloaded from the  
3 NIST GiaB FTP site along with proprietary 30x HG001 samples. Illumina model for  
4 DNAscope was trained using a subset of the GiAB HG001 and HG005 data as well.  
5 None of the tested samples were used during model training. Training was performed  
6 across all chromosomes with the exception of chromosome 20.

7

#### 8 **Variant accuracy evaluation**

9 All VCF files generated from this benchmark study were collected for accuracy  
10 evaluation. Firstly, they were separated into SNP and InDel subgroups and each  
11 subgroup was then compared against NIST truth set using Vcfeval from RTGtools [52]  
12 to calculate a F-score as representation of accuracy.

13

#### 14 **Acknowledgements**

1 We would like to acknowledge the ongoing contributions and support of all MGI-  
2 Shenzhen and BGI research employees, and also acknowledge Frank Hu from Sentieon  
3 for commenting on the manuscript.

4

## 5 **Funding**

6 This work was supported by Shenzhen Peacock Plan.No.KQTD20150330171505310,  
7 the National Key R&D Program of China (2018YFC0910201), and the Key R&D  
8 Program of Guangdong Province (2019B020226001).

9

## 10 **Availability of data and materials**

11 The data reported in this study are available in the CNGB Nucleotide Sequence  
12 Archive.(<https://db.cngb.org/cnsa>; accession number CNP0000602,CNP0000466. ) All  
13 the other data used here are included within the published article and its Additional files.

14

## 1   **References**

- 2   1.   Meienberg J, Zerjavic K , Keller I, et al. New insights into the performance of  
3       human whole-exome capture platforms. *Nucleic Acids Res*, 2015, 43(11): e76.
- 4   2.   Belkadi A, Bolze A, Itan Y, et al. Whole-genome sequencing is more powerful than  
5       whole-exome sequencing for detecting exome variants. *Proc Natl Acad Sci U S A*,  
6       2015, 112(17): 5473-5478.
- 7   3.   Meienberg J, Bruggman R, Oexle K, et al. Clinical sequencing: is WGS the better  
8       WES? *Hum Genet*, 2016, 135: 359–362.
- 9   4.   Carss KJ, Arno G, Erwood M, Stephens J, Sanchis-Juan A, Hull S, Megy K, Grozeva  
10       D, Dewhurst E, Malka S, et al. Comprehensive Rare Variant Analysis via Whole-  
11       Genome Sequencing to Determine the Molecular Pathology of Inherited Retinal  
12       DiseaseAn integrated map of structural variation. *Am J Hum Genet*, 2017, 100, 75-  
13       90.
- 14   5.   Fujimoto A, Furuta M, Totoki Y, et al. Whole-genome mutational landscape and  
15       characterization of noncoding and structural mutations in liver cancer. *Nat Genet*,  
16       2016, 48(5): 500-9.

- 1 6. Satta G, Lipman M, Smith GP, et al. Mycobacterium tuberculosis and whole-  
2 genome sequencing: how close are we to unleashing its full potential? *Clin Microbiol*  
3 *Infect*,2018, 24(6): 604-609.
- 4 7. <https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data>
- 5 8. Kobschull JM, Zador AM. Sources of PCR-induced distortions in high-  
6 throughput sequencing data sets. *Nucleic Acids Res*, 2015, 43(21): e143.
- 7 9. Aird D, Ross MG, Chen WS, et al. Analyzing and minimizing PCR amplification bias  
8 in Illumina sequencing libraries. *Genome Biol*, 2011, 12(2): R18.
- 9 10. Chen YC, Liu T, Yu CH, et al. Effects of GC Bias in Next-Generation-Sequencing  
10 Data on De Novo Genome Assembly. *PLoS One*,2013, 8(4): e62856.
- 11 11. Ross MG, Russ C, Costello M, et al. Characterizing and measuring bias in  
12 sequence data. *Genome Biology*, 2013, 14(5): R51.
- 13 12. Kozarewa I, Ning Z, Quail MA, et al. Amplification-free Illumina sequencing  
14 preparation facilitates improved mapping and assembly of (G+C)-biased genomes.  
15 *Nat Methods*, 2009, 6(4):291-5.

- 1 13. Jones MB, Highlander SK, Anderson EL, et al. Library preparation methodology can  
2 influence genomic and functional predictions in human microbiome research. *Proc*  
3 *Natl Acad Sci U S A*, 2015, 112 (45) 14024-14029.
- 4 14. McKenna A, Hanna M, Banks E, et al. The Genome Analysis Toolkit A MapReduce  
5 framework for analyzing next-generation DNA sequencing data. *Genome Res*, 2010,  
6 20(9): 1297-303.
- 7 15. Heldenbrand JR, Baheti S, Bockol MA, et al. Performance benchmarking of  
8 GATK3.8 and GATK4. *bioRxiv*, 348565v1; doi: <https://doi.org/10.1101/348565>.
- 9 16. Friedman S. Deep learning in GATK4 [Internet]. GATK | Blog. 2017. Available from:  
10 <https://software.broadinstitute.org/gatk/blog?id=10996>.
- 11 17. Poplin R, Chang PC, Alexander D, et al. A universal SNP and small-InDel variant  
12 caller using deep neural networks. *Nat Biotechnol*, 2018, 36(10):983-987.
- 13 18. <https://github.com/Sentieon/sentieon-dnascopy-ml>
- 14 19. Luo R, Sedlazeck FJ, Lam TW, et al. A multi-task convolutional deep neural network  
15 for variant calling in single molecule sequencing. *Nat Commun*, 2019, 10(1):998.

- 1 20. Adessi C1, Matton G, Ayala G, et al. Solid phase DNA amplification characterisation  
2 of primer attachment and amplification mechanisms. *Nucleic Acids Res*, 2000,  
3 28(20): E87.
- 4 21. Blanco L, Bernad A, Lázaro JM, et al. Highly efficient DNA synthesis by the phage phi  
5 29 DNA polymerase. Symmetrical mode of DNA replication. *J Biol Chem*, 1989,  
6 264(15): 8935-40.
- 7 22. Ali MM, Li F, Zhang Z, Zhang K, et al. Rolling circle amplification: a versatile tool for  
8 chemical biology, materials science and medicine. *Chem Soc Rev*, 2014, 43(10):  
9 3324-41.
- 10 23. <https://en.mgitech.cn/article/detail/285.html>
- 11 24. Zhou X, Rokas A. Prevention, diagnosis and treatment of high-throughput  
12 sequencing data pathologies. *Mol Ecol*, 2014, 23(7): 1679-700.
- 13 25. Bansal V. A computational method for estimating the PCR duplication rate in DNA  
14 and RNA-seq experiments. *BMC Bioinformatics*, 2017, 18(Suppl 3): 43.
- 15 26. Cornish A, Guda C. A Comparison of Variant Calling Pipelines Using Genome in a  
16 Bottle as a Reference. *Biomed Res Int*, 2015, 2015: 456479.

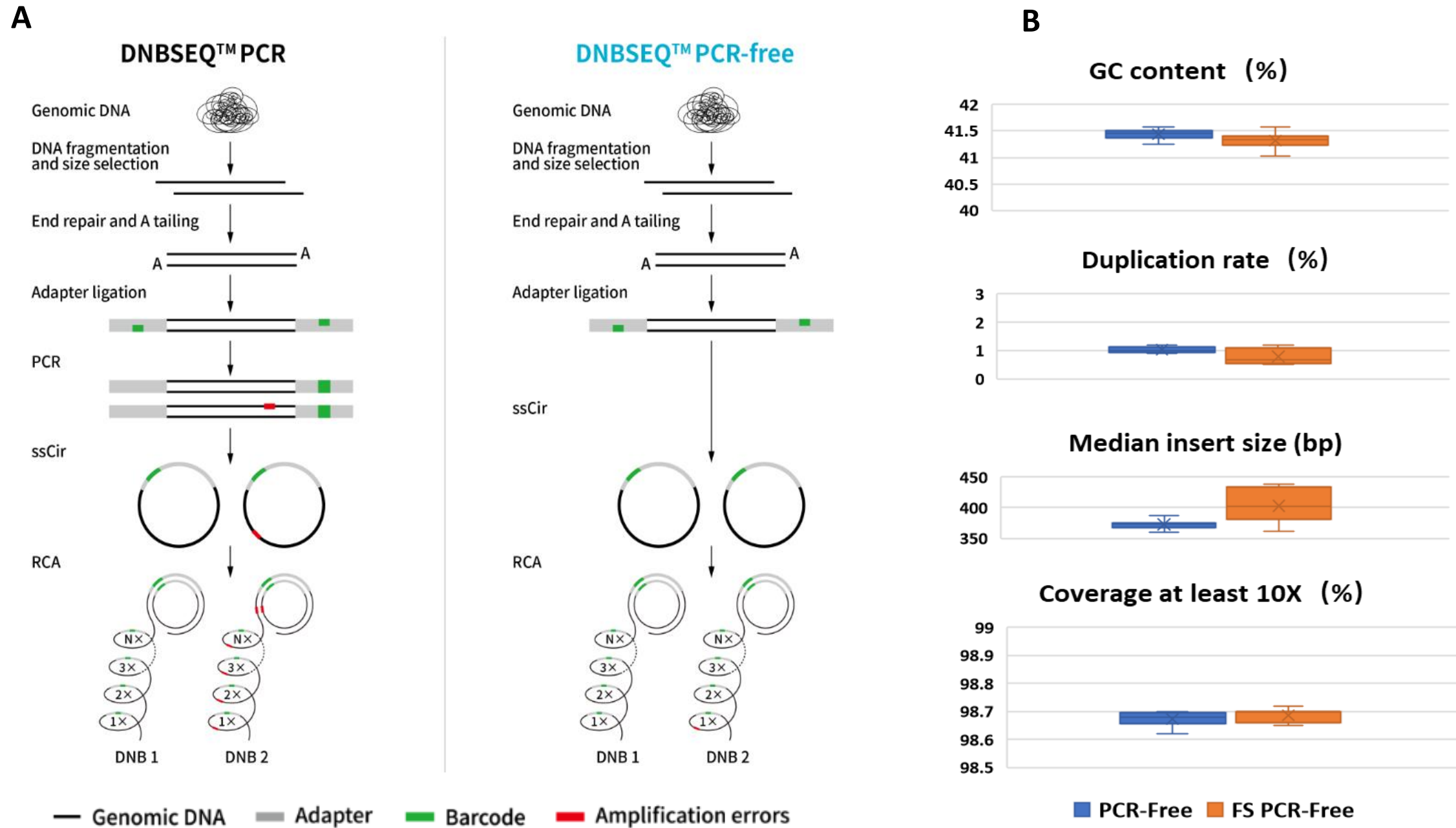


- 1 27. Hwang S, Kim E, Lee I, et al. Systematic comparison of variant calling pipelines  
2 using gold standard personal exome variants. *Sci Rep*, 2015, 5: 17875.
- 3 28. Zook JM, Catoe D, McDaniel J, et al. Extensive sequencing of seven human  
4 genomes to characterize benchmark reference materials. *Sci. Data*, 2016, 3:160025.
- 5 29. Supernat A, Vidarsson OV, Steen VM, et al. Comparison of three variant callers for  
6 human whole genome sequencing. *Sci Rep*, 2018, 8(1): 17851.
- 7 30. <https://precision.fda.gov/challenges/truth/results>
- 8 31. Fang H, Wu Y, Narzisi G, et al. Reducing INDEL calling errors in whole genome and  
9 exome sequencing data. *Genome Med*, 2014, 6(10): 89.
- 10 32. Willig LK, Petrikin JE, Smith LD, et al. Whole-genome sequencing for identification  
11 of Mendelian disorders in critically ill infants: a retrospective analysis of diagnostic  
12 and clinical findings. *Lancet Respir Med* , 2015, 5: 377–87.
- 13 33. Kim S, Scheffler K, Halpern AL, et al. Strelka2: fast and accurate calling of germline  
14 and somatic variants. *Nat Methods*, 2018, 15(8): 591-594.
- 15 34. Kendig KI, Baheti S, Bockol MA, et al. Sentieon DNaseq Variant Calling Workflow  
16 Demonstrates Strong Computational Performance and Accuracy. *Front Genet*, 2019,  
17 10: 736.

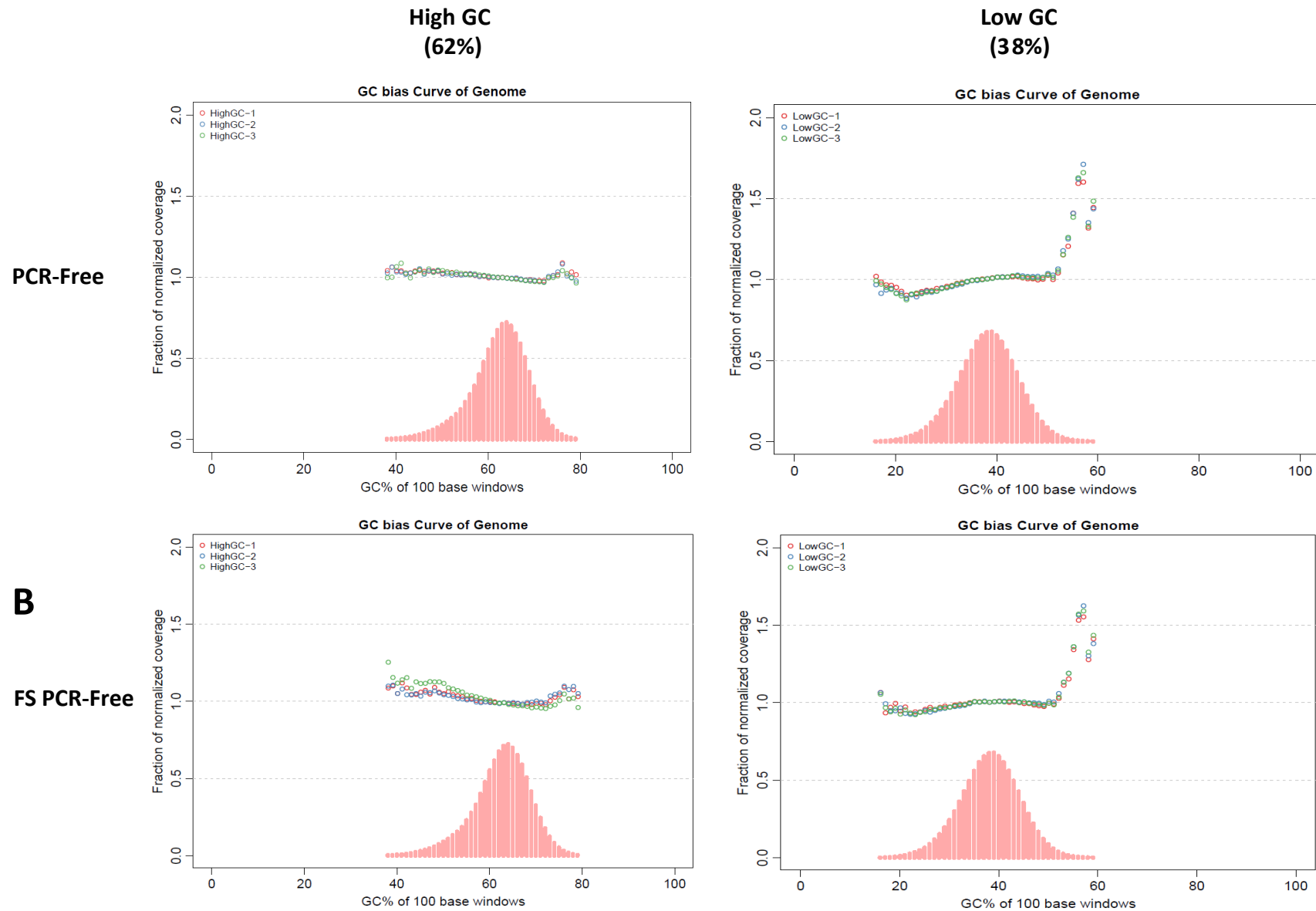
- 1 35. Chiang C, Layer RM, Faust GG, et al. SpeedSeq: ultra-fast personal genome  
2 analysis and interpretation. *Nat Methods*, 2015,12(10): 966-8.
- 3 36. Goyal A, Kwon HJ, Lee K, et al. Ultra-Fast Next Generation Human Genome  
4 Sequencing Data Processing Using DRAGENTM Bio-IT Processor for Precision  
5 Medicine. *Open Journal of Genetics*, 2017, 7, 9-19.
- 6 37. <https://en.mgitech.cn/article/detail/germline.html>
- 7 38. <https://en.mgitech.cn/product/detail/DNBSEQ-T7.html>
- 8 39. Peters BA, Kermani BG, Sparks AB, et al. Accurate whole-genome sequencing and  
9 haplotyping from 10 to 20 human cells. *Nature*, 2012, 487(7406): 190-5.
- 10 40. Zheng GXY, Lau BT, Schnall-Levin M, et al. Haplotyping germline and cancer  
11 genomes with high-throughput linked-read sequencing. *Nature Biotech*, 2016, 34:  
12 303–311.
- 13 41. Wang O, Chin R, Cheng XF, et al. Efficient and unique cobarcoding of second-  
14 generation sequencing reads from long DNA molecules enabling cost-effective and  
15 accurate sequencing, haplotyping, and de novo assembly. *Genome Res*, 2019,  
16 29: 798-808.

- 1 42. <https://gatkforums.broadinstitute.org/gatk/discussion/23202/does-sequencing->
- 2 platform-have-effect-on-the-variant-detection
- 3 43. <https://blog.dnanexus.com/2018-05-31-training-and-applying-genomic-deep->
- 4 learning-models/
- 5 44. Chen J, Li X, Zhong H, et al. Systematic comparison of germline variant
- 6 calling pipelines cross multiple next-generation sequencers. *Sci Rep*, 2019,
- 7 9(1):9345.
- 8 45. Li H. Aligning sequence reads, clone sequences and assembly contigs with bwa-
- 9 mem. *arXiv*. 1303. 3997.
- 10 46. <https://github.com/lh3/seqtk>
- 11 47. Chen Y, Chen Y, Shi C, et al. SOAPnuke: A MapReduce acceleration-supported
- 12 software for integrated quality control and preprocessing of high-throughput
- 13 sequencing data. *GigaScience*, 2018, 7: 1–6.
- 14 48. Li H, Handsaker B, Wysoker A, et al. The sequence alignment/map format and
- 15 samtools. *Bioinformatics*, 2009, 25(16): 2078-9.
- 16 49. <https://github.com/broadinstitute/picard>

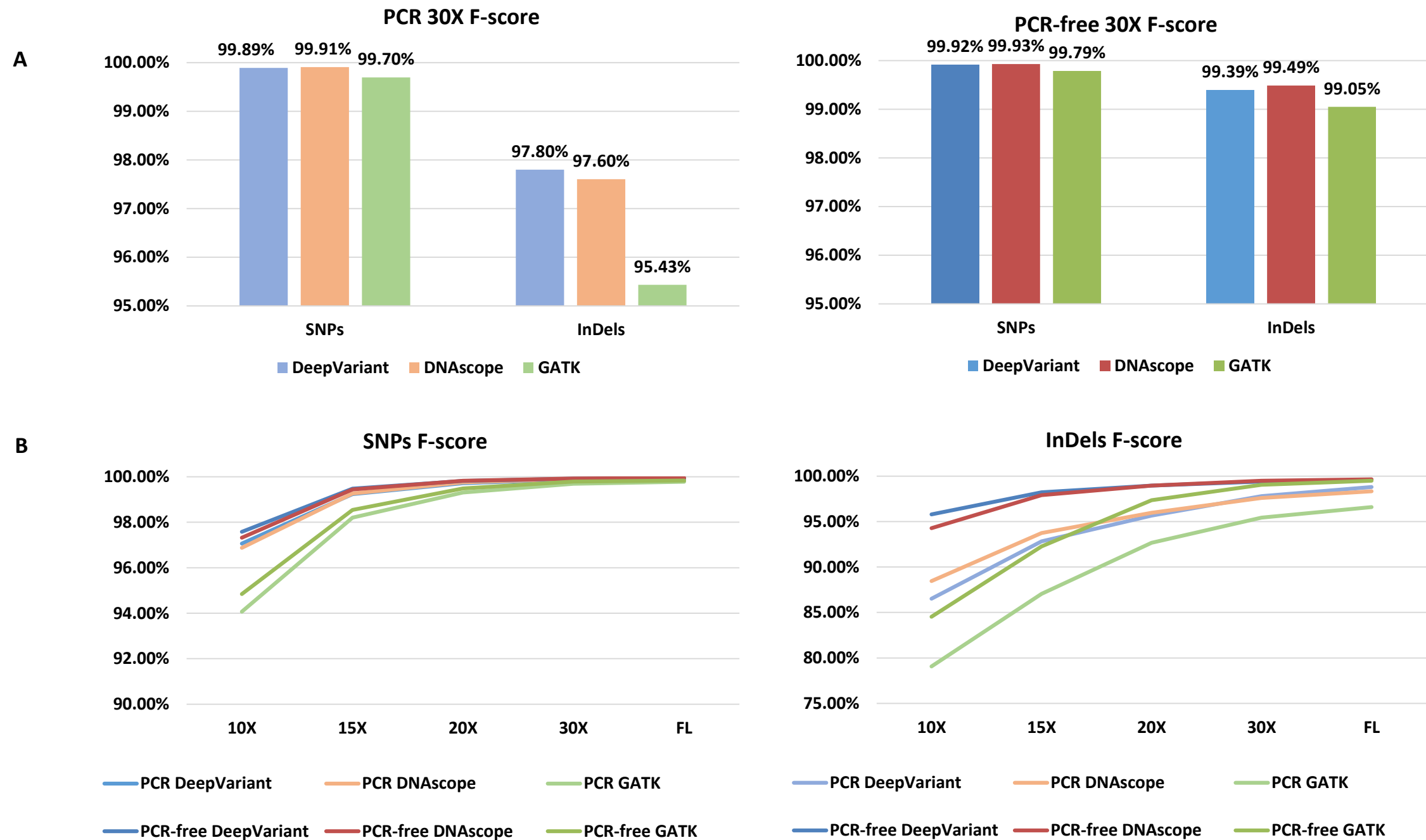
- 1 50. McKenna A, Hanna M, Banks E, et al. The genome analysis toolkit: A MapReduce
- 2 framework for analyzing next-generation DNA sequencing data. *Genome Res*, 2010,
- 3 20: 1297–303.
  
- 4 51. [https://github.com/MGI-tech-bioinformatics/MGI\\_DeepVariant\\_model](https://github.com/MGI-tech-bioinformatics/MGI_DeepVariant_model)
  
- 5 52. <https://github.com/RealTimeGenomics/rtg-tools>



**Fig1. DNBSEQ™ WGS PCR vs PCR-free workflows and general performance of PCR-free libraries.** (a) MPS library construction workflows of WGS PCR and PCR-free libraries. RCA (rolling circle amplification) is used to increase signal intensity during array formation, which is followed by sequencing DNBs with DNBSEQ™ technology. Individual copies from the same DNB are replicated independently using the same ssCir template. Therefore, amplification errors cannot accumulate. Black, genomic DNA; Gray rectangle, adapter; green, barcode; red, amplification errors. (b) Two sets of 9 replicates from 1  $\mu$ g NA12878 DNA were processed with MGIEasy PCR-Free DNA Library Prep Set (blue) or MGIEasy FS PCR-Free DNA Library Prep Set (orange). The GC content, Duplication rate, Median Insert size, and regions with >10x Coverage were calculated and plotted. The error bars represent the standard deviations across the replicates.

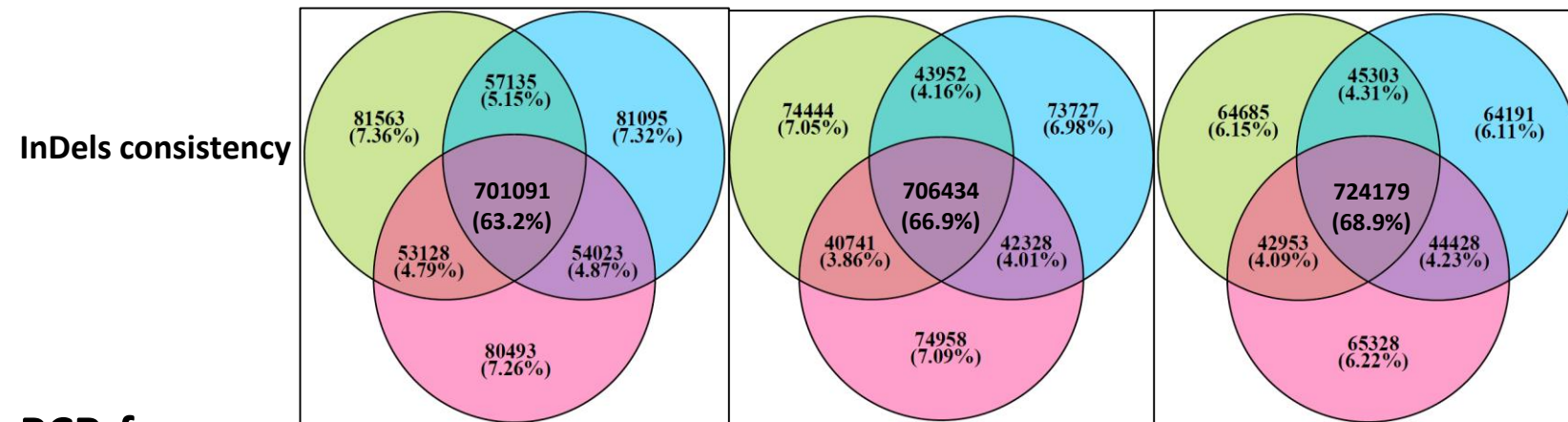
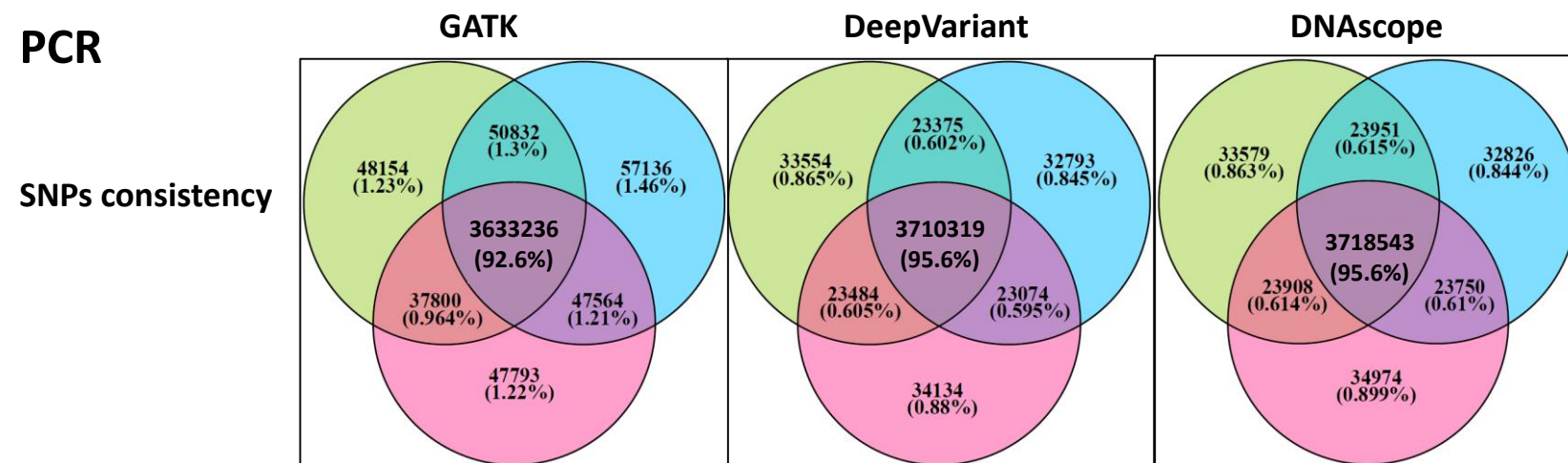
**A**

**Figure 2. Coverage of the microbial genomes, *Olsenella profusa* (62% GC) (left) and *Bacillus megaterium* (38% GC) (right) with (A) MGIEasy PCR-Free DNA Library Pre Set and (B) MGIEasy FS PCR-Free DNA Library Pre Set. Read coverage across the range of the GC content, calculated in 100 bp windows (pink bars) and normalized coverage (colored dots). 3 replicates (red, blue and green dots) were included in the normalized coverage VS GC content analysis.**

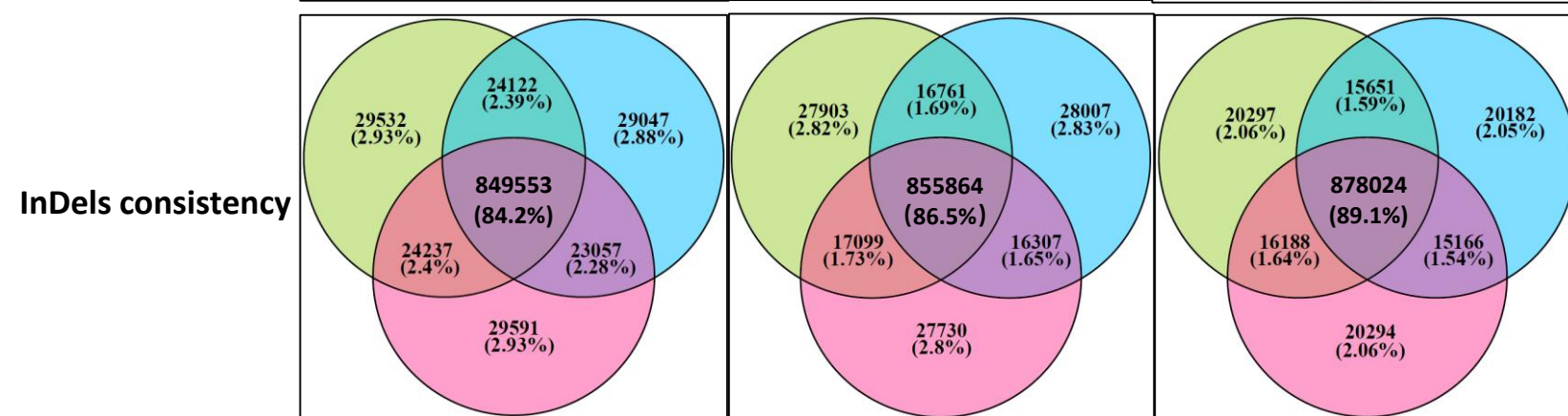
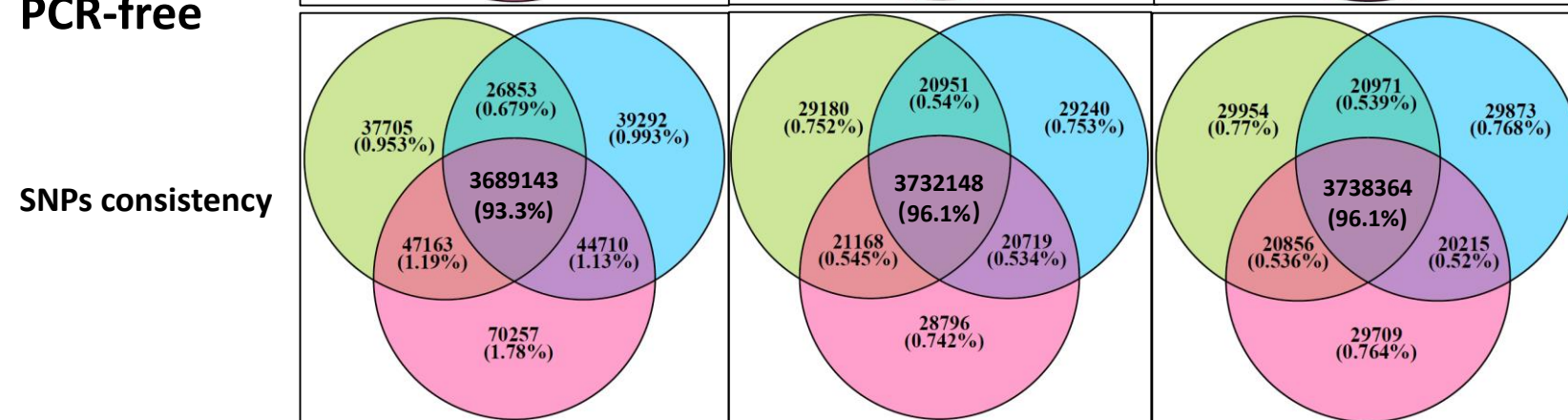


**Figure 3. variant calling performance in PCR vs PCR-free libraries** (A) F-Scores of 30x sequencing PCR and PCR-free libraries were compared for the accuracy performance of different variant callers. (B) Down-sampled dataset representing different sequencing depth were processed to exam the tolerance to shallow data from different variant callers. PCR-free libraries (Dark color) and machine learning variant callers (blue and red) showed good accuracy at >15X, and even better than PCR+GATK combination at 30X depth.

**PCR**



**PCR-free**



**Figure4. Variant consistency of 3 replicates of PCR and PCR-free libraries** Consistency analysis were conducted on the 3 libraries generated by the same library kit and variant calling pipelines. Venn-Diagram were generated to show the common shared variants and the unique variants.



Method	PCR-1			PCR-2			PCR-3			PCR-free-1			PCR-free-2			PCR-free-3		
	15x	30x	FL	15x	30x	FL	15x	30x	FL	15x	30x	FL	15x	30x	FL	15x	30x	FL
Clean reads (x10 <sup>6</sup> )	296.64	633.73	940.37	296.80	633.76	954.92	296.45	632.78	949.34	295.93	631.36	926.38	296.03	632.98	944.92	296.13	631.49	964.32
Clean bases (Gb)	44.49	95.05	141.05	44.52	95.06	143.23	44.46	94.91	142.40	44.39	94.70	138.95	44.40	94.94	141.73	44.42	94.72	144.64
Insert size peak (bp)	383	383	383	374	374	374	375	375	375	374	374	375	375	375	375	375	375	375
GC content (%)	41.64	41.64	41.64	41.78	41.78	41.78	41.49	41.49	41.49	41.53	41.53	41.53	41.48	41.48	41.48	41.46	41.47	41.46
Mapping rate (%)	99.91	99.82	99.82	99.91	99.84	99.84	99.89	99.82	99.82	99.89	99.83	99.83	99.72	99.65	99.65	99.93	99.86	99.86
Unique rate (%)	99.42	95.41	94.78	99.48	95.53	94.92	99.35	95.23	94.51	99.69	95.76	95.38	99.71	95.82	95.44	99.77	95.95	95.61
Duplicate rate (%)	0.58	1.56	2.24	0.52	1.44	2.10	0.65	1.72	2.49	0.31	1.07	1.50	0.29	1.03	1.45	0.23	0.91	1.31
Mismatch rate (%)	0.89	0.87	0.86	0.91	0.89	0.88	0.95	0.93	0.92	0.74	0.74	0.73	0.74	0.74	0.73	0.71	0.70	0.70
Average seq depth (X)	15.37	30.83	45.46	15.39	30.89	46.27	15.34	30.71	45.74	15.37	30.81	45.04	15.35	30.86	45.91	15.4	30.91	47.04
Coverage (%)	99.07	99.13	99.16	99.07	99.13	99.16	99.08	99.14	99.17	99.1	99.16	99.18	99.1	99.16	99.18	99.1	99.16	99.18
Coverage at least 4X (%)	98.65	98.96	99.03	98.64	98.97	99.03	98.65	98.97	99.04	98.76	99.01	99.06	98.75	99.01	99.06	98.74	99.00	99.06
Coverage at least 10X (%)	89.78	98.58	98.83	89.61	98.58	98.84	89.55	98.58	98.84	90.5	98.70	98.89	90.44	98.69	98.89	90.49	98.69	98.90

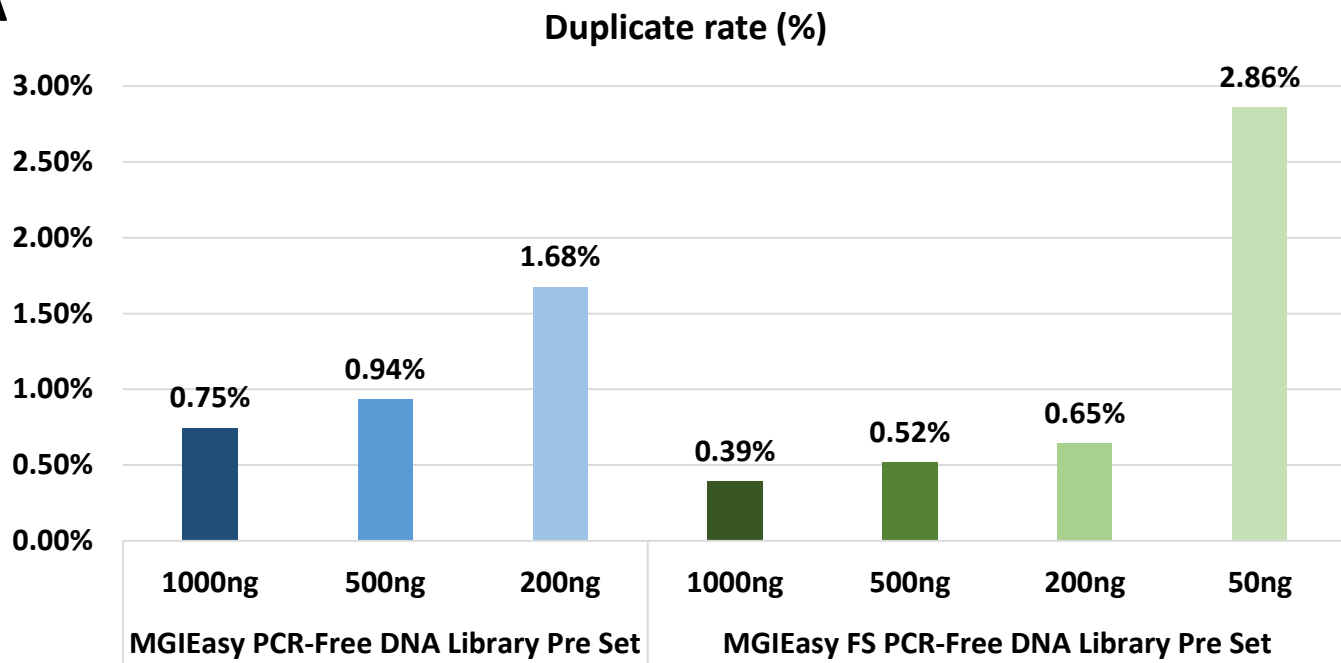
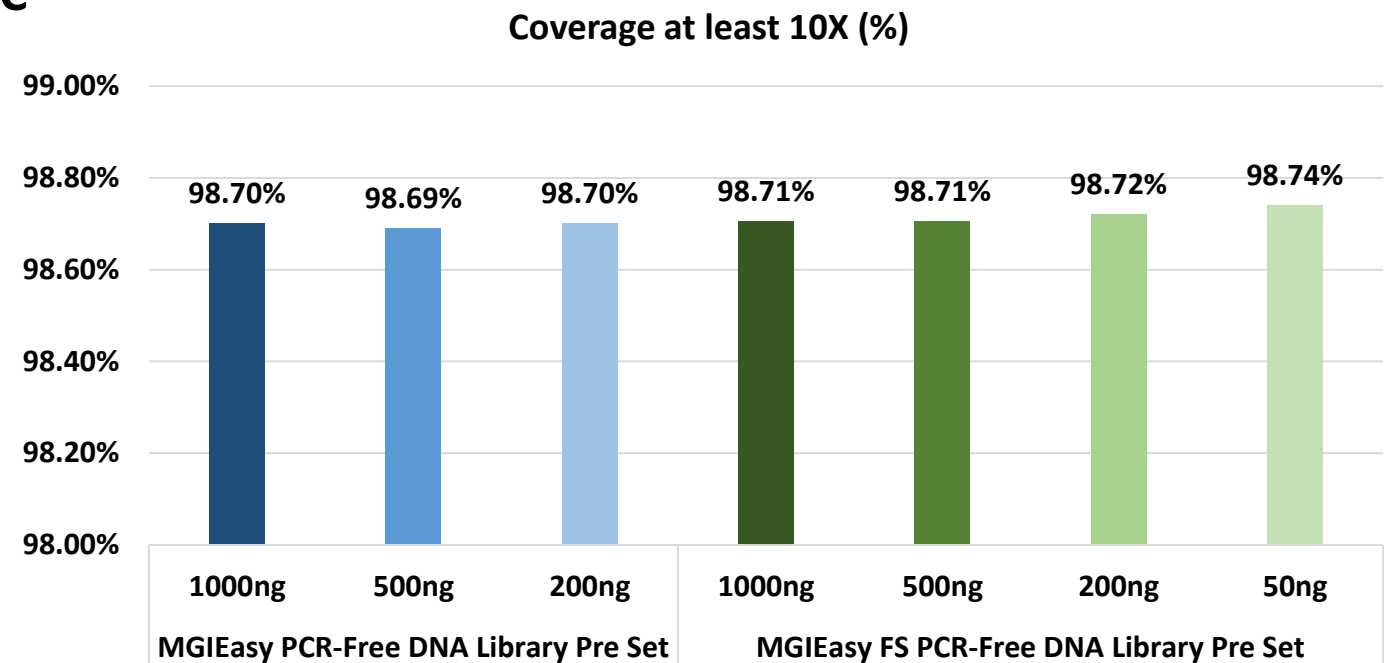
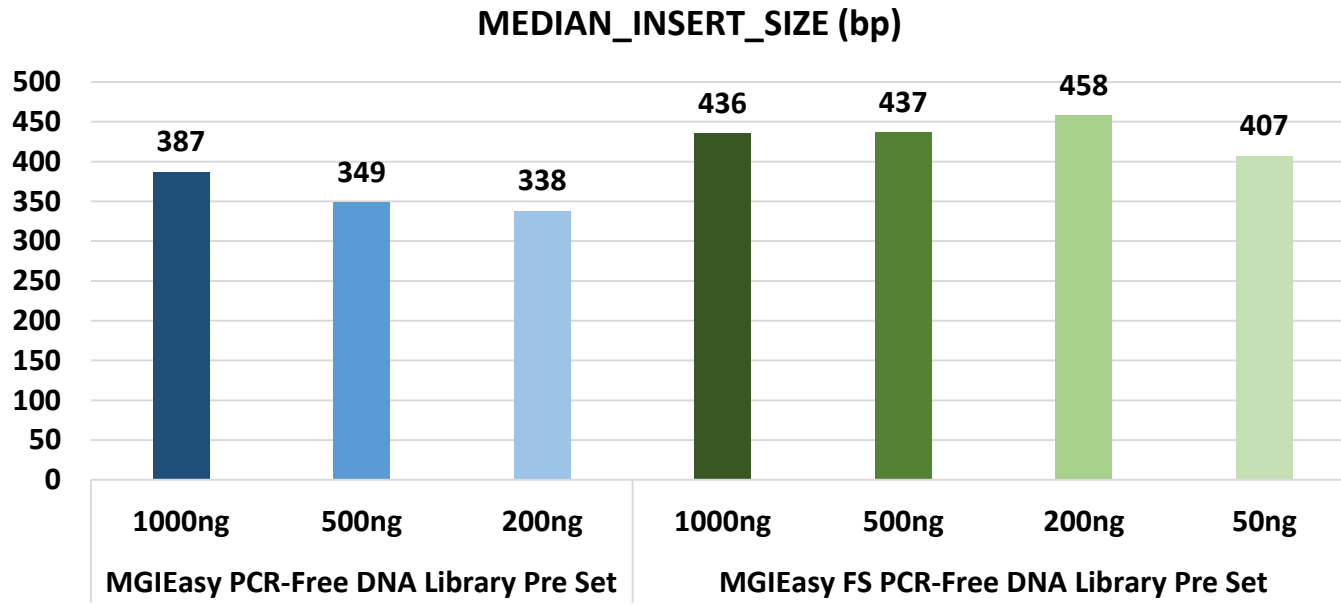
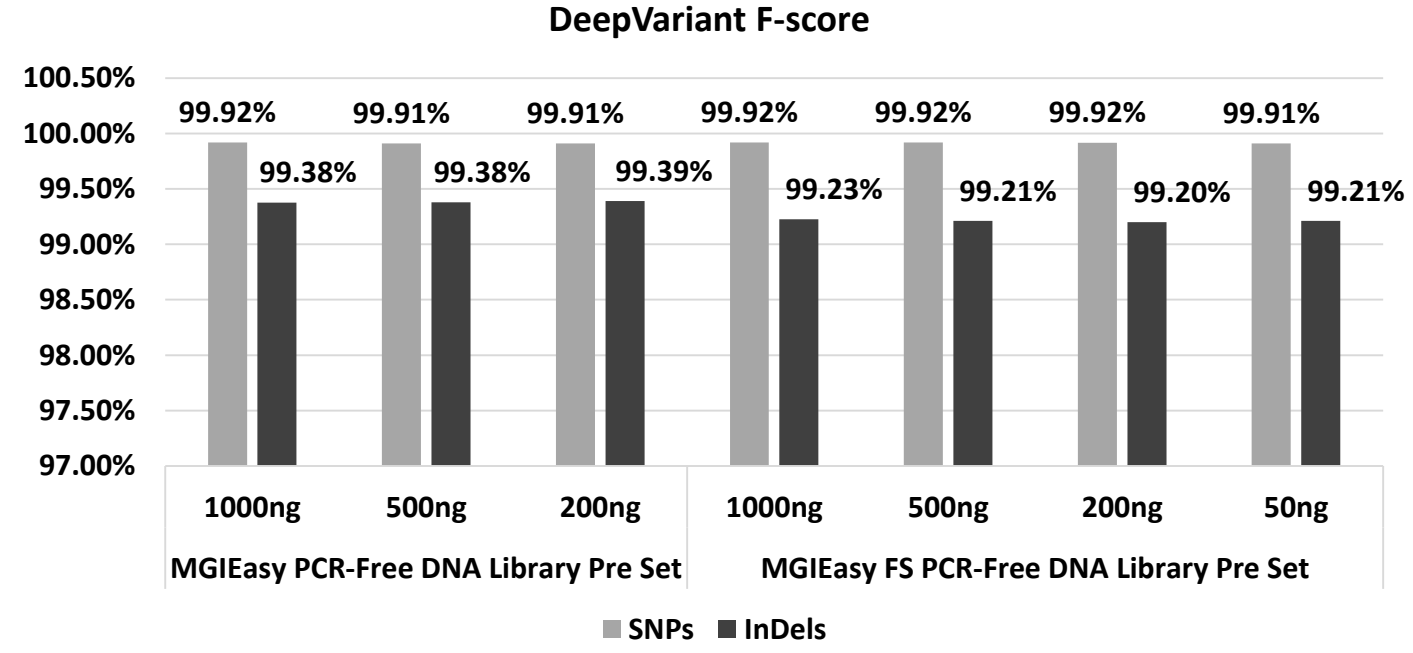
**Table1. Mapping performance of three replicates for PCR and PCR-free libraries** QC matrix from 3 PCR and 3 PCR-free libraries was collected and showed. Each sample was down sampled to 15x and 30x as well as FL (full lane). “FL” represents Full Lane sequencing data amount approximately equals 46x coverage.

Pipeline	Variant type	Method Depth	PCR			PCR-free		
			15x	30x	FL	15x	30x	FL
GATK	SNPs	True positive	3106478	3193586	3198933	3126614	3199331	3202461
		False positive	9615	2806	2479	8544	2586	2347
		False negative	103778	16671	11324	83643	10926	7796
		Precision	99.69%	99.91%	99.92%	99.73%	99.92%	99.93%
		Sensitivity	96.77%	99.48%	99.65%	97.39%	99.66%	99.76%
		F-score	98.21%	99.70%	99.79%	98.55%	99.79%	99.84%
	InDels	True positive	397779	458110	466343	429105	474921	477918
		False positive	34787	20699	17891	19730	2766	1686
		False negative	83485	23154	14922	52160	6345	3349
		Precision	91.96%	95.67%	96.30%	95.60%	99.42%	99.65%
		Sensitivity	82.65%	95.19%	96.90%	89.16%	98.68%	99.31%
		F-score	87.05%	95.43%	96.60%	92.27%	99.05%	99.48%
DeepVariant	SNPs	True positive	3174776	3205991	3207021	3184757	3207148	3207769
		False positive	12952	2913	1871	7837	2111	1691
		False negative	35481	4266	3236	25499	3109	2488
		Precision	99.59%	99.91%	99.94%	99.75%	99.93%	99.95%
		Sensitivity	98.90%	99.87%	99.90%	99.21%	99.90%	99.92%
		F-score	99.24%	99.89%	99.92%	99.48%	99.92%	99.93%
	InDels	True positive	436804	468237	474129	470682	477616	478556
		False positive	22906	8008	4389	6310	2124	1615
		False negative	44546	13082	7185	10695	3690	2745
		Precision	95.02%	98.32%	99.08%	98.67%	99.56%	99.66%
		Sensitivity	90.74%	97.28%	98.51%	97.78%	99.23%	99.43%
		F-score	92.83%	97.80%	98.80%	98.23%	99.39%	99.55%
DNAscope	SNPs	True positive	3171728	3205902	3207318	3180864	3207055	3207911
		False positive	8351	1565	1021	6034	1254	901
		False negative	38529	4355	2939	29393	3202	2346
		Precision	99.74%	99.95%	99.97%	99.81%	99.96%	99.97%
		Sensitivity	98.80%	99.86%	99.91%	99.08%	99.90%	99.93%
		F-score	99.27%	99.91%	99.94%	99.45%	99.93%	99.95%
	InDels	True positive	443656	466964	471037	466852	477969	479103
		False positive	21653	8632	5773	5617	1620	1202
		False negative	37609	14301	10229	14413	3299	2164
		Precision	95.35%	98.19%	98.79%	98.81%	99.66%	99.75%
		Sensitivity	92.18%	97.03%	97.88%	97.00%	99.31%	99.55%
		F-score	93.74%	97.60%	98.33%	97.90%	99.49%	99.65%

**Table2. Average variant calling performance of three replicates for PCR and PCR-free libraries using three variant callers** Variant called from each library and variant caller was evaluated by Vcfeval tool in RTGtools against NIST truth set at high confident regions. Then average values from the same library construction method were generated and showed in table. "FL" represents Full Lane sequencing data amount approximately equals 46x coverage.

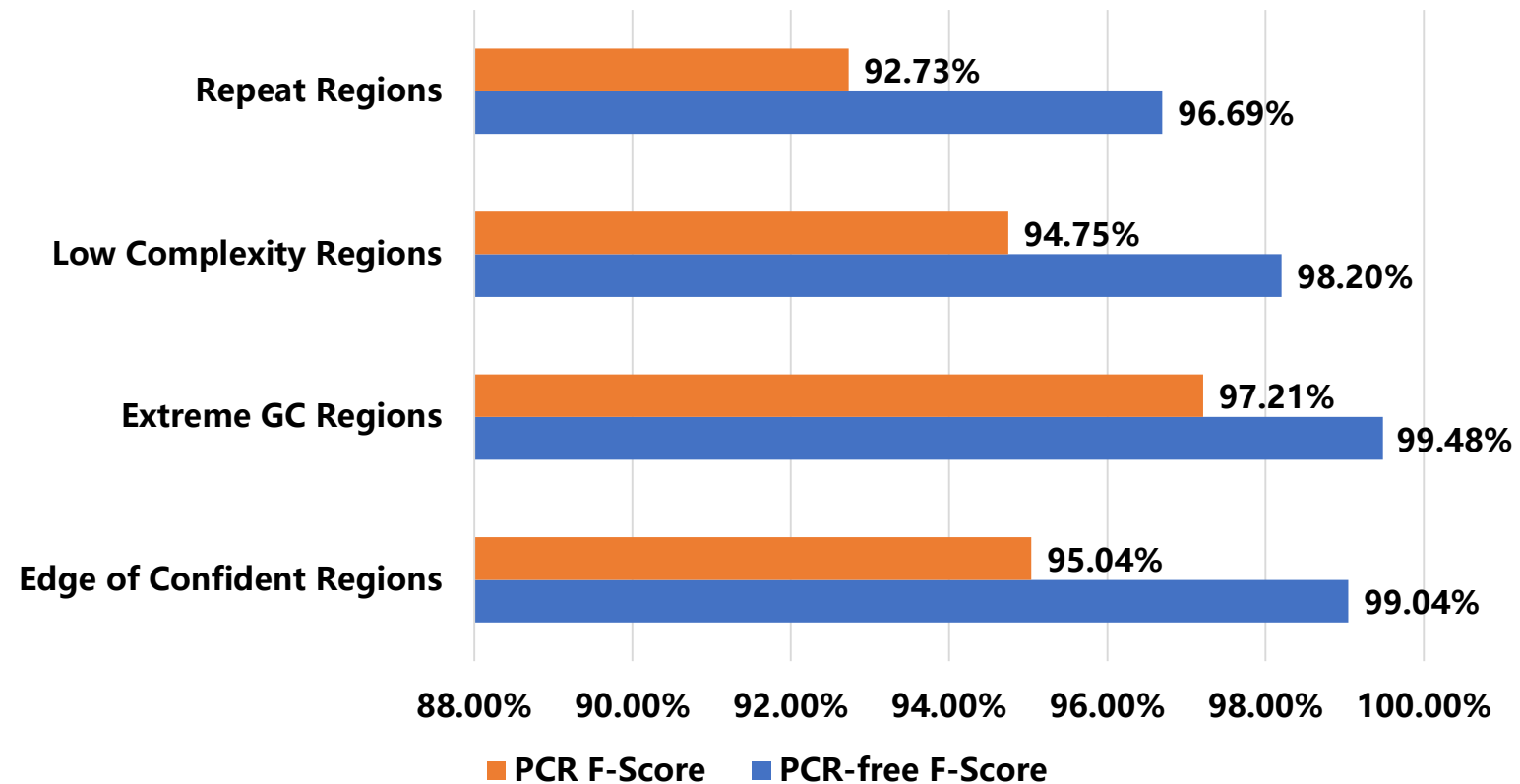
Library Kit	MGIEasy PCR-Free DNA Library Prep Set	MGIEasy FS PCR-Free DNA Library Prep Set	MGIEasy PCR-Free DNA Library Prep Set	MGIEasy PCR-Free DNA Library Prep Set	PCR-Free research lib	TruSeq DNA PCR-Free Library Prep Kits			
DNA input	250ng	1µg	1µg	1µg	1µg	2µg	2µg	2µg	
Sequence platform	DNBSEQ-T7	MGISEQ-2000	MGISEQ-2000	Novaseq	MGISEQ-2000	Hiseq4000	xTen	Novaseq	
<b>Average seq depth (X)</b>	30.46	30.68	30.91	30.75	30.6	30.02	31.5	30.43	
<b>Mapping rate (%)</b>	99.99	99.89	99.86	99.82	99.83	99.03	99.85	99.94	
<b>Duplicate rate (%)</b>	3.6	0.6	0.91	18.19	0.76	13.62	12.31	10.43	
<b>Coverage at least 10X (%)</b>	98.64	98.73	98.69	98.83	98.82	98.8	98.91	98.87	
<b>SNPs</b>	<b>True positive</b>	3206204	3207301	3207056	3207540	3207839	3206782	3207834	3208037
	<b>False positive</b>	1545	1209	1255	1331	1048	2704	1280	1262
	<b>False negative</b>	4053	2956	3201	2717	2418	3475	2423	2220
	<b>Precision</b>	99.95%	99.96%	99.96%	99.96%	99.97%	99.92%	99.96%	99.96%
	<b>Sensitivity</b>	99.87%	99.91%	99.90%	99.92%	99.92%	99.89%	99.92%	99.93%
	<b>F-score</b>	99.91%	99.94%	99.93%	99.94%	99.95%	99.90%	99.94%	99.95%
<b>Indels</b>	<b>True positive</b>	476673	478047	478004	477283	479066	469353	477420	477017
	<b>False positive</b>	2266	1882	1595	2479	1304	7642	2598	2879
	<b>False negative</b>	4592	3218	3265	3984	2201	11913	3846	4250
	<b>Precision</b>	99.53%	99.61%	99.67%	99.48%	99.73%	98.40%	99.46%	99.40%
	<b>Sensitivity</b>	99.05%	99.33%	99.32%	99.17%	99.54%	97.52%	99.20%	99.12%
	<b>F-score</b>	99.29%	99.47%	99.49%	99.33%	99.64%	97.96%	99.33%	99.26%

**Table3. Mapping and Variant calling performance of PCR-free libraries sequenced on different sequencing platforms** DNBSEQ-T7, MGISEQ-2000 data was generated following kit instructions. Illumina Hiseq4000, HiSeq xTen, Novaseq data was downloaded from Illumina Basespace website. To compare the sequence platform only, one PCR-free library was generated by MGIEasy FS PCR-Free DNA Library Prep Set with illumina's adapter and sequenced on Novaseq by a third party sequencing service provider. To further improve the library preparation method, one additional PCR-free library was prepared in house using a research protocol for the library construction and sequenced on MGISEQ-2000. Variants from each library was called using DNAscope pipeline and accuracy was evaluated by vcfeval tool in RTGtools against NIST truth set at high confident regions.

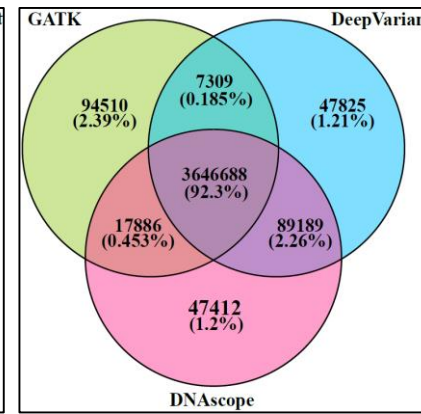
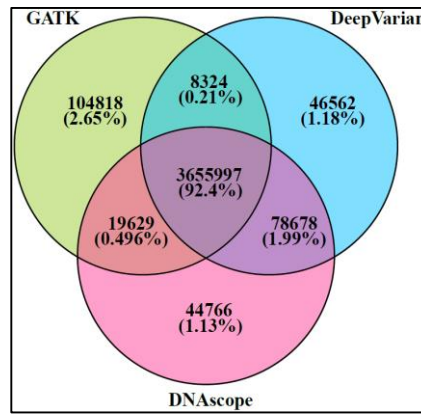
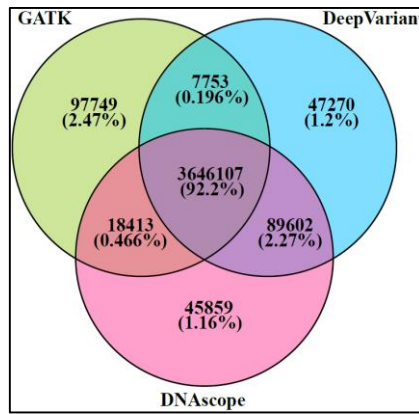
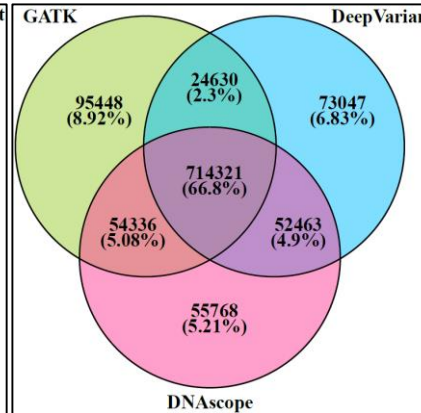
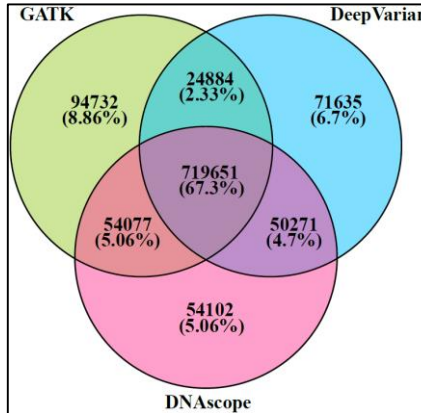
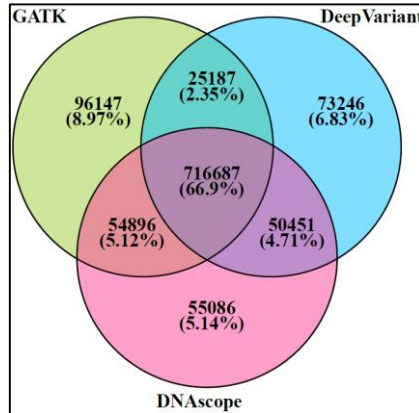
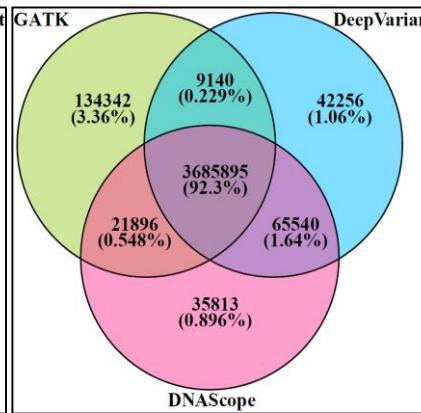
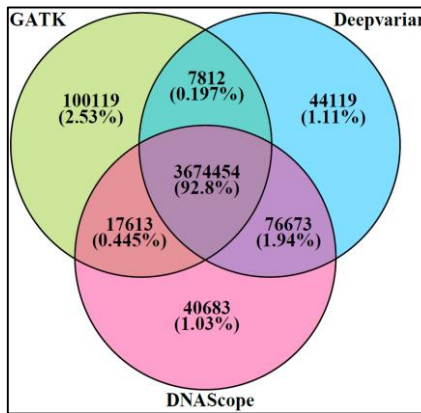
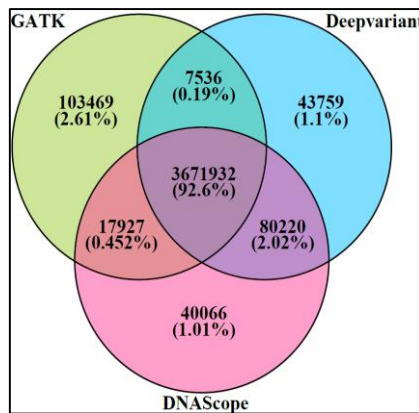
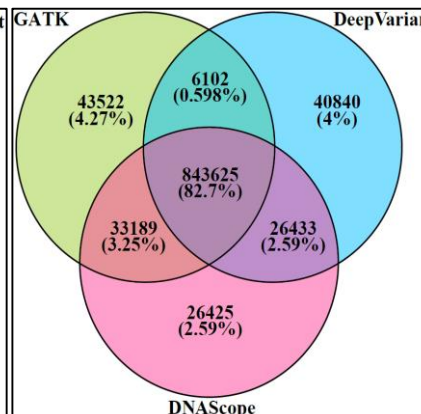
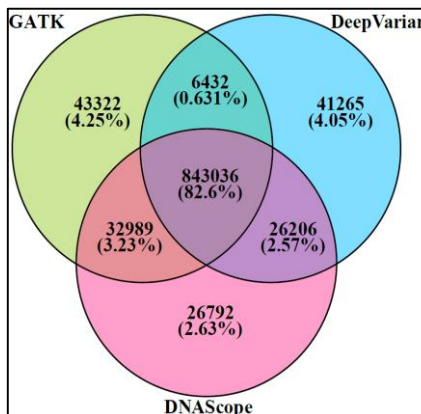
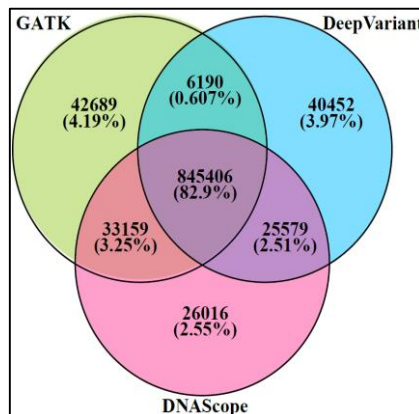
**A****C****B****D**

**Fig1S.** The data performance with different library inputs using either ultrasonic shearing (MGIEasy PCR-Free DNA Library Prep Set) or enzymatic shearing method (MGIEasy FS PCR-Free DNA Library Prep Set) Each condition was repeated three times, and the average value was calculated for display.

### InDels F1-score In Selected Genome Regions



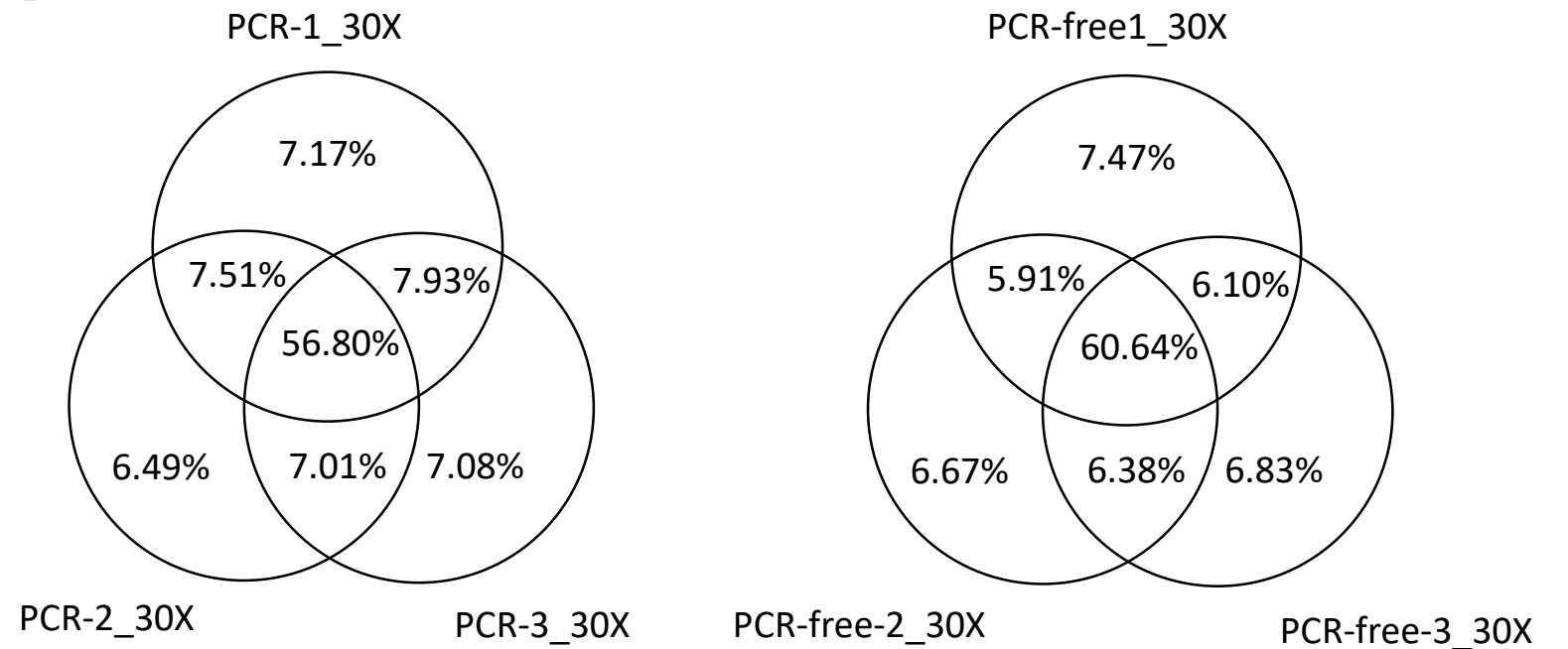
**Figure 2S. InDels F-score in Selected Genome Regions** In comparison between one PCR sample and PCR-free sample, a number of genome subsets regions were selected and F1-scores generated from Sentieon DNAscope were compared within each region. InDels showed more significant differences between two samples, and especially in 4 categories above. These categories are selected from GA4GH stratification regions with category name changed for easier understand. Repeat Regions matches category "lowcmp\_AllRepeats\_51to200bp\_gt95identity\_merged" ; Low Complexity Regions matches category "lowcmp\_Human\_Full\_Genome\_TRDB\_hg19\_150331\_all\_merged" ; Extreme GC Regions matches category "gclt30orgt55" ; Edge of Confident Regions matches category "TS\_boundary" .

**PCR-1****PCR-2****PCR-3****SNPs  
Consistency****InDels  
Consistency****PCR-free-1****PCR-free-2****PCR-free-3****SNPs  
Consistency****InDels  
Consistency**

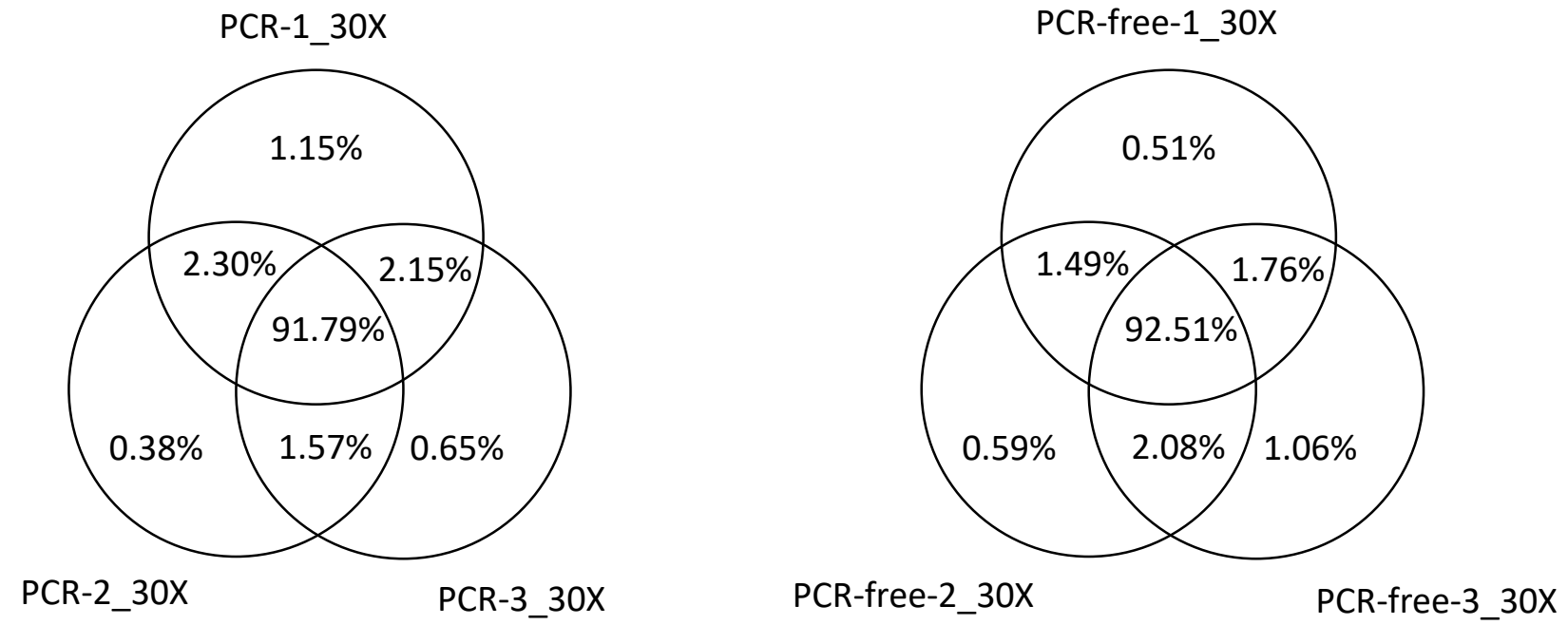
**Fig 3S. Consistency of three pipelines for SNPs and InDels.** Consistency analysis were conducted on the 3 variant calling pipelines on the same library. Venn-Diagram were generated to show the common shared variants and the unique variants.

**A**

Samples	TRA	DEL	DUP	INV	INS
PCR-1_30X	203	3528	175	227	99
PCR-2_30X	205	3488	138	224	99
PCR-3_30X	204	3477	172	225	95
PCR-free-1_30X	293	3684	189	231	104
PCR-free-2_30X	288	3625	182	254	106
PCR-free-3_30X	251	3689	168	256	123

**B**

**Figure 4S.** (A) SV events detected by DNAscope in 6 testing samples. SV events were reported into one of the five categories: Translocation, Deletion, Duplication, Inversion, and Insertion. PCR-free libraries showed higher number of reported SV events, and (B) higher consistency among all 3 samples as well.

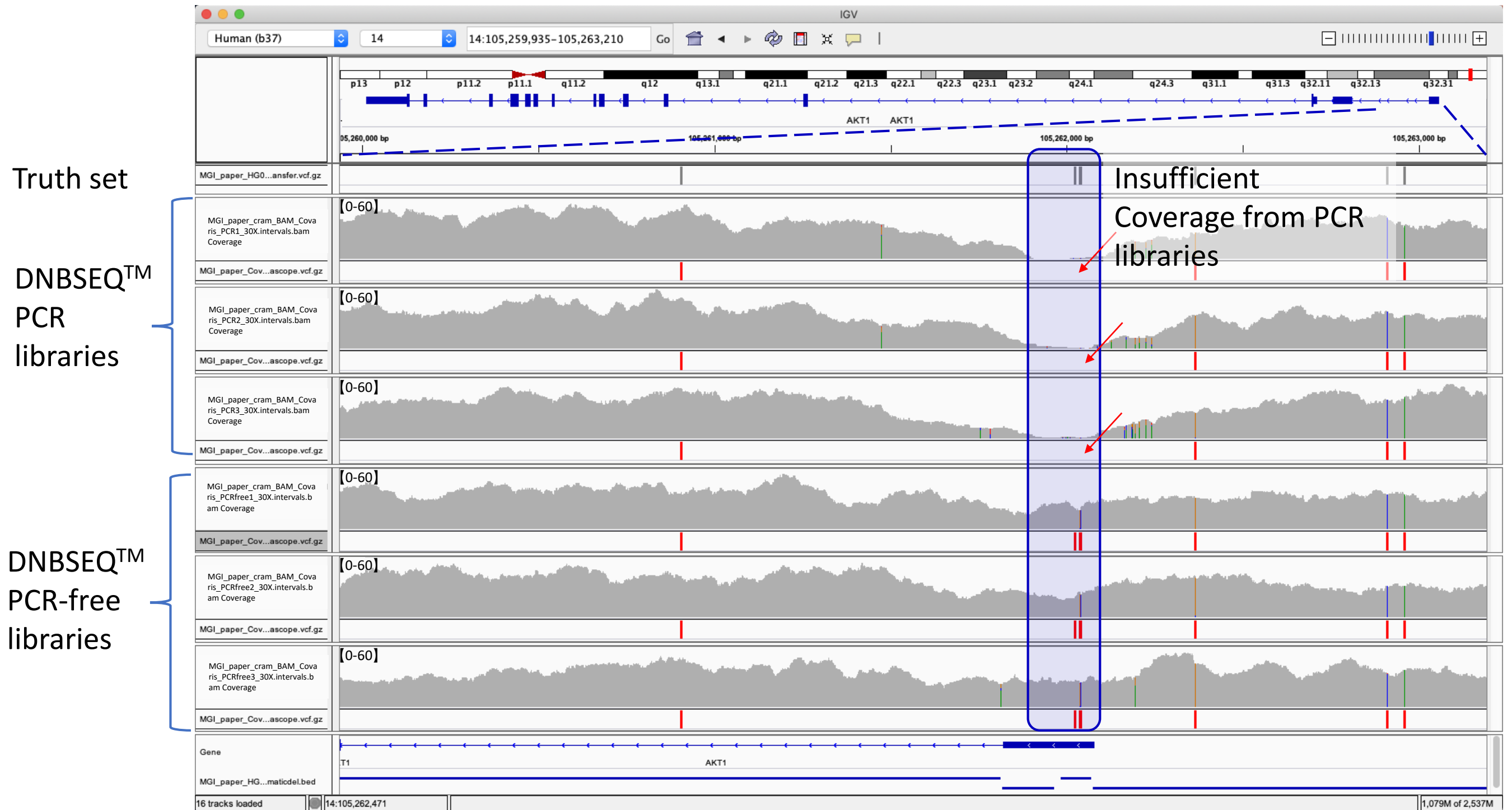


**Figure 5S.** CNV events detected by DNAscope in 6 testing samples. 3-way comparison was conducted to analysis consistency among replicates.

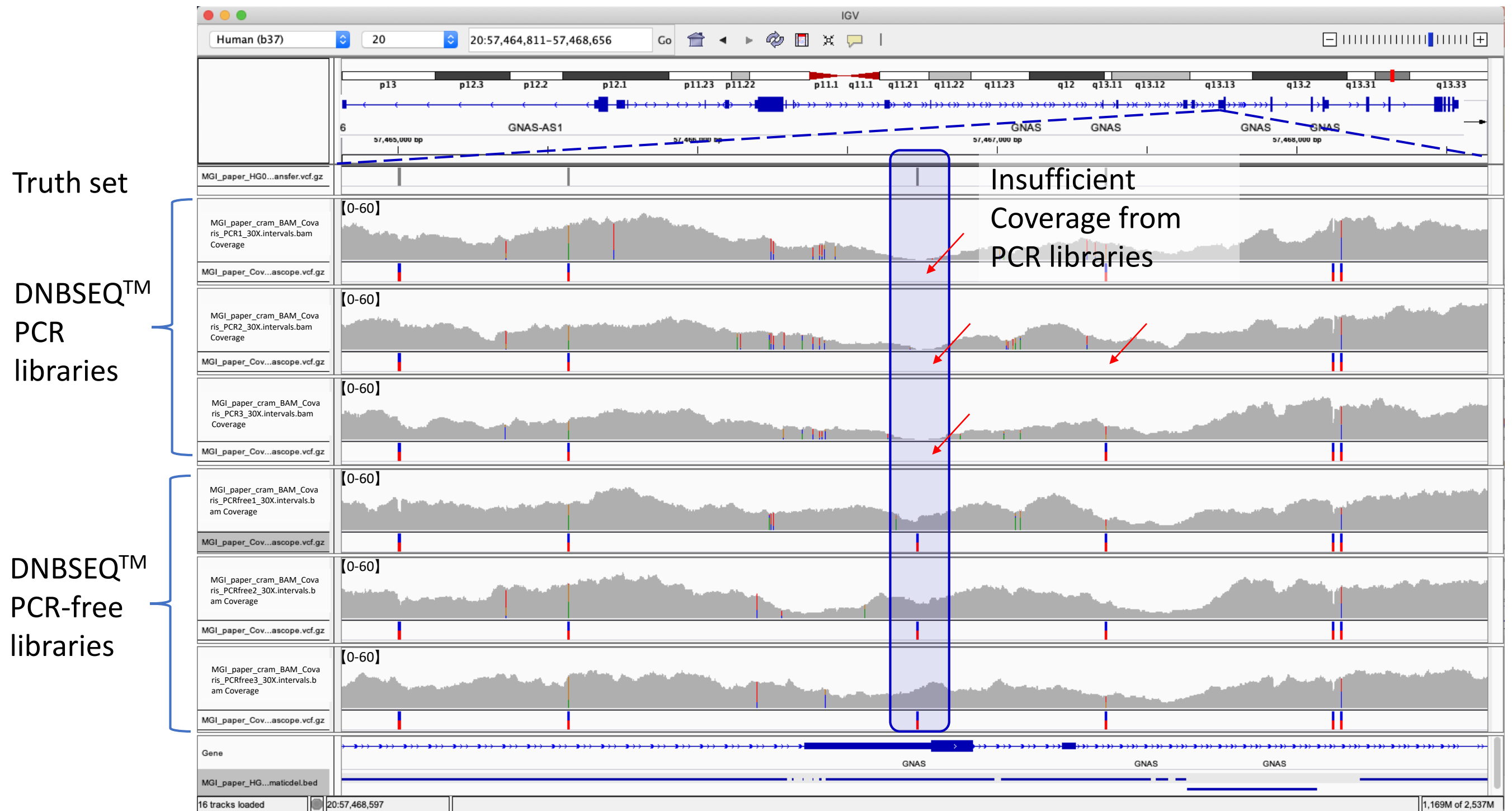


Pipline	Variant type	Method Depth	PCR-1			PCR-2			PCR-3			PCR-free-1			PCR-free-2			PCR-free-3		
			15x	30x	FL	15x	30x	FL	15x	30x	FL	15x	30x	FL	15x	30x	FL	15x	30x	FL
GATK	SNPs	True positive	3110618	3193126	3200655	3108081	3195584	3197789	3100734	3192049	3198355	3128229	3198582	3201565	3123297	3198921	3202891	3128317	3200489	3202928
		False positive	9634	2770	2750	9352	2939	2205	9858	2708	2483	8708	2403	2570	8559	2268	2385	8364	3086	2087
		False negative	99638	17131	9602	102175	14673	12468	109522	18208	11902	82028	11675	8692	86960	11336	7366	81940	9768	7329
		Precision	99.69%	99.91%	99.91%	99.70%	99.91%	99.93%	99.68%	99.92%	99.92%	99.72%	99.92%	99.92%	99.73%	99.93%	99.93%	99.73%	99.90%	99.93%
		Sensitivity	96.90%	99.47%	99.70%	96.82%	99.54%	99.61%	96.59%	99.43%	99.63%	97.44%	99.64%	99.73%	97.29%	99.65%	99.77%	97.45%	99.70%	99.77%
		F-score	98.27%	99.69%	99.81%	98.24%	99.73%	99.77%	98.11%	99.67%	99.78%	98.57%	99.78%	99.82%	98.49%	99.79%	99.85%	98.58%	99.80%	99.85%
	InDels	True positive	398324	458096	466055	398912	458728	466804	396100	457507	466171	429199	475007	477924	428632	474850	477933	429485	474906	477897
		False positive	35000	21111	18550	34291	20094	17163	35069	20892	17960	19737	2741	1688	19793	2775	1719	19661	2782	1650
		False negative	82941	23168	15210	82352	22536	14462	85162	23757	15094	52066	6258	3343	52632	6417	3334	51783	6361	3370
		Precision	91.92%	95.59%	96.17%	92.08%	95.80%	96.45%	91.87%	95.63%	96.29%	95.60%	99.43%	99.65%	95.59%	99.42%	99.64%	95.62%	99.42%	99.66%
		Sensitivity	82.77%	95.19%	96.84%	82.89%	95.32%	96.99%	82.30%	95.06%	96.86%	89.18%	98.70%	99.31%	89.06%	98.67%	99.31%	89.24%	98.68%	99.30%
		F-score	87.10%	95.39%	96.50%	87.24%	95.56%	96.72%	86.82%	95.35%	96.58%	92.28%	99.06%	99.48%	92.21%	99.04%	99.47%	92.32%	99.05%	99.48%
DeepVariant	SNPs	True positive	3175227	3206052	3207074	3174871	3205898	3206874	3174231	3206023	3207115	3185071	3207185	3207765	3184356	3207120	3207784	3184844	3207140	3207757
		False positive	12381	2879	1875	13175	2917	1886	13301	2942	1851	7973	2119	1700	7830	2138	1698	7709	2077	1675
		False negative	35030	4205	3182	35386	4359	3383	36026	4234	3142	25185	3072	2492	25900	3137	2473	25413	3117	2500
		Precision	99.61%	99.91%	99.94%	99.59%	99.91%	99.94%	99.58%	99.91%	99.94%	99.75%	99.93%	99.95%	99.75%	99.93%	99.95%	99.76%	99.94%	99.95%
		Sensitivity	98.91%	99.87%	99.90%	98.90%	99.86%	99.89%	98.88%	99.87%	99.90%	99.22%	99.90%	99.92%	99.19%	99.90%	99.92%	99.21%	99.90%	99.92%
		F-score	99.26%	99.89%	99.92%	99.24%	99.89%	99.92%	99.23%	99.89%	99.92%	99.48%	99.92%	99.93%	99.47%	99.92%	99.94%	99.48%	99.92%	99.93%
	InDels	True positive	436833	468201	474042	437513	468409	474219	436066	468101	474127	470823	477670	478584	470466	477597	478546	470756	477582	478538
		False positive	23254	8085	4458	22359	7786	4257	23105	8152	4453	6233	2111	1614	6514	2120	1614	6182	2141	1616
		False negative	44515	13120	7283	43834	12915	7093	45290	13212	7179	10553	3641	2719	10915	3711	2756	10617	3718	2760
		Precision	94.95%	98.30%	99.07%	95.14%	98.36%	99.11%	94.97%	98.29%	99.07%	98.69%	99.56%	99.66%	98.63%	99.56%	99.66%	98.70%	99.55%	99.66%
		Sensitivity	90.75%	97.27%	98.49%	90.89%	97.32%	98.53%	90.59%	97.25%	98.51%	97.81%	99.24%	99.44%	97.73%	99.23%	99.43%	97.79%	99.23%	99.43%
		F-score	92.80%	97.79%	98.78%	92.97%	97.84%	98.82%	92.73%	97.77%	98.79%	98.25%	99.40%	99.55%	98.18%	99.39%	99.55%	98.25%	99.39%	99.54%
DNAScope	SNPs	True positive	3172347	3205898	3207358	3172082	3205807	3207242	3170755	3206001	3207353	3181144	3207111	3207908	3180305	3206997	3207905	3181142	3207056	3207920
		False positive	8291	1522	996	8325	1577	1046	8438	1595	1021	6050	1269	927	5969	1237	888	6083	1255	889
		False negative	37909	4359	2899	38175	4450	3015	39502	4256	2904	29113	3146	2349	29951	3260	2352	29115	3201	2337
		Precision	99.74%	99.95%	99.97%	99.74%	99.95%	99.97%	99.73%	99.95%	99.97%	99.81%	99.96%	99.97%	99.81%	99.96%	99.97%	99.81%	99.96%	99.97%
		Sensitivity	98.82%	99.86%	99.91%	98.81%	99.86%	99.91%	98.77%	99.87%	99.91%	99.09%	99.90%	99.93%	99.07%	99.90%	99.93%	99.09%	99.90%	99.93%
		F-score	99.28%	99.91%	99.94%	99.27%	99.91%	99.94%	99.25%	99.91%	99.94%	99.45%	99.93%	99.95%	99.44%	99.93%	99.95%	99.45%	99.93%	99.95%
	InDels	True positive	443499	466856	470799	444329	467200	471196	443141	466836	471115	467078	477999	479118	466450	477903	479072	467027	478004	479118
		False positive	21902	8738	6002	21250	8438	5540	21807	8721	5777	5615	1617	1194	5775	1648	1217	5460	1595	1195
		False negative	37767	14409	10466	36938	14066	10070	38121	14429	10150	14187	3268	2149	14815	3365	2194	14238	3265	2149
		Precision	95.29%	98.16%	98.74%	95.44%	98.23%	98.84%	95.31%	98.17%	98.79%	98.81%	99.66%	99.75%	98.78%	99.66%	99.75%	98.84%	99.67%	99.75%
		Sensitivity	92.15%	97.01%	97.83%	92.32%	97.08%	97.91%	92.08%	97.00%	97.89%	97.05%	99.32%	99.55%	96.92%	99.30%	99.54%	97.04%	99.32%	99.55%
		F-score	93.70%	97.58%	98.28%	93.85%	97.65%	98.37%	93.67%	97.58%	98.34%	97.92%	99.49%	99.65%	97.84%	99.48%	99.65%	97.93%	99.49%	99.65%

**Table1S. Variant calling performance of PCR and PCR-free with three variant callers** Variant called from each library and variant caller was evaluated by Vcfeval tool in RTGtools against NIST truth set at high confident regions.



**Figure 6S. False negative detection in DNBSEQ™ PCR libraries that may cause mis-diagnostic in clinical utilities related to AKT1 gene.** Clinical Phenotype: Cancer, Cowden syndrome 6; All PCR libraries failed to detect: Chr14: 105262025 T to TC insertion at Exon1 UTR, and Chr14: 105262041 CG to GC SNP at Exon1 UTR.



**Figure 7S. False negative detection in DNBSEQ™ PCR libraries that may cause mis-diagnostic in clinical utilities related to GNAS gene.** Clinical Phenotype: Pseudohypoparathyroidism; All PCR libraries failed to detect: Chr20: 57466734 C to CCG insertion at Exon1 UTR.



**Figure 8S. False negative detection in 3 ILMN libraries and 2 DNBSEQ™ PCR libraries that may cause mis-diagnostic in clinical utilities related to MAF gene.** Clinical Phenotypes: Ayme-Gripp syndrome; Cataract 21, multiple types; All ILMN libraries failed to detect: Chr16: 79632062 C to CT insertion at Exon1, which was picked up by all DNBSEQ™ PCR free libraries at even 15x depth.