

# 1 **High performance of a GPU-accelerated variant calling tool** 2 **in genome data analysis**

3 **Qian Zhang<sup>1†</sup>, Hao Liu<sup>2†</sup>, Fengxiao Bu<sup>1\*</sup>**

4 <sup>1</sup>Institute of Rare Diseases, West China Hospital, Sichuan University, Chengdu, China

5 <sup>2</sup>Department of Neurosurgery, West China Hospital, Sichuan University, Chengdu,  
6 China

7 \* **Correspondence:** Fengxiao Bu, [bufengxiao@wchscu.cn](mailto:bufengxiao@wchscu.cn)

8 † These authors have contributed equally to this work

9 **Keywords: bioinformatics tool, WGS, variant calling, acceleration tool, GPU-**  
10 **based acceleration, GATK pipeline**

## 11 **Abstract**

12 Rapid advances in next-generation sequencing (NGS) have facilitated ultralarge  
13 population and cohort studies that utilized whole-genome sequencing (WGS) to  
14 identify DNA variants that may impact gene function. Massive sequencing data  
15 require highly efficient bioinformatics tools to complete read alignment and variant  
16 calling as the fundamental analysis. Multiple software and hardware acceleration  
17 strategies have been developed to boost the analysis speed. This study  
18 comprehensively evaluated the germline variant calling of a GPU-based acceleration  
19 tool, BaseNumber, using WGS datasets from several sources, including gold-standard  
20 samples from the Genome in a Bottle (GIAB) project and the Golden Standard of  
21 China Genome (GSCG) project, resequenced GSCG samples, and 100 in-house  
22 samples from the China Deafness Genetics Consortium (CDGC) project. Sequencing  
23 data were analyzed on the GPU server using BaseNumber, the variant calling outputs  
24 of which were compared to the reference VCF or the results generated by the  
25 Burrows-Wheeler Aligner (BWA) + Genome Analysis Toolkit (GATK) pipeline on a  
26 generic CPU server. BaseNumber demonstrated high precision (99.32%) and recall  
27 (99.86%) rates in variant calls compared to the standard reference. The variant calling  
28 outputs of the BaseNumber and GATK pipelines were very similar, with a mean F1 of  
29 99.69%. Additionally, BaseNumber took only 23 minutes on average to analyze a 48X  
30 WGS sample, which was 215.33 times shorter than the GATK workflow. The GPU-  
31 based BaseNumber provides a highly accurate and ultrafast variant calling capability,  
32 significantly improving the WGS analysis efficiency and facilitating time-sensitive  
33 tests, such as clinical WGS genetic diagnosis, and sheds light on the GPU-based  
34 acceleration of other omics data analyses.

35 (3254 words)

## 36 **1 Introduction**

37 Genomic DNA variation represents a critical genetic source of variation that alters  
38 protein function and expression, affects human diseases and phenotypes, and can  
39 provide markers associated with various populational traits. Next-generation  
40 sequencing (NGS) is commonly used to identify DNA variants in a high-throughput  
41 manner. Compared to other enrichment-based NGS approaches, whole-genome  
42 sequencing (WGS) covers up to 98% of the human genome, providing comprehensive  
43 and unbiased variant detection in individuals <sup>1</sup>. With the development of NGS at a  
44 pace far exceeding that predicted by Moore's Law in recent decades, the cost of  
45 sequencing is rapidly decreasing, facilitating the broad application of WGS in clinical  
46 and research genetic testing <sup>2,3</sup>. Ultralarge WGS projects have been initiated, such as  
47 Genomics England <sup>4</sup>, All of Us <sup>5</sup>, and the China Metabolic Analytics Project <sup>6</sup>, aiming  
48 for in-depth understanding of the underlying molecular mechanisms of human  
49 diseases and traits, where accurate identification and interpretation of DNA variants  
50 are the foundation. However, due to the enormous amount of sequencing data, the  
51 efficiency of variant calling has become a significant bottleneck, with technical  
52 complexities remaining.

53 Conventional variant calling pipelines are often based on CPU servers and open-  
54 source software, such as Genome Analysis Toolkit (GATK) <sup>7</sup> and VarScan <sup>8</sup>. GATK  
55 HaplotypeCaller calls germline variants through the local de novo assembly of  
56 haplotypes in a region showing signs of variation. The pipeline based on GATK Best  
57 Practice reportedly completed WGS variant calling on a sample in about 24 hours <sup>9</sup>.  
58 Significant improvements are required to meet the massive needs of WGS studies. A  
59 widely used solution is the high-performance cluster (HPC)-based speedup, which  
60 reduces the total computation time by using multiple CPU servers to analyze a set of  
61 data simultaneously. The use of large clusters implies an increase in the cost of initial  
62 equipment procurement, maintenance, and energy consumption. Therefore, it is  
63 essential to maximize the utilization of computing resources on a single server to  
64 achieve satisfying efficiency. Multiple algorithms and tools for speeding up variant  
65 callers have been developed <sup>10-12</sup>, such as Sentieon DNaseq, a reimplement of  
66 the GATK Best Practice workflow <sup>13</sup>. By optimizing and recompiling the variant  
67 calling algorithms, DNaseq can achieve a 10-fold increase in processing speed while  
68 providing results nearly identical to those of the GATK pipeline <sup>13,14</sup>. However, these  
69 CPU computing-based solutions have two major drawbacks: 1) limited parallelism in  
70 a CPU environment and 2) the growing gap between CPU computing power and  
71 sequencing data throughput. As a promising solution, heterogeneous computing-based  
72 approaches have emerged to use different architectures of computing units, such as  
73 graphics processing units (GPUs) or field-programmable gate arrays (FPGAs), which  
74 are more parallelable than CPUs and more powerful under certain conditions. As an  
75 example, Illumina Dragen integrated an FPGA card to boost the variant caller. It also  
76 adopted its own bioinformatics algorithms to make good use of the integrated circuits

77 (ICs) on board. The read mapper of Dragen was based on a hash algorithm instead of  
78 the Burrows–Wheeler transform (BWT) algorithm used by Burrows-Wheeler Aligner  
79 (BWA). In addition, the variant calling algorithm was based on a hidden Markov  
80 model, whose running speed can be improved significantly because of its parallelable  
81 nature<sup>15</sup>. However, the concerns of FPGA server-based variant callers include the  
82 limited reusability of the dedicated FPGA server for other analytic tasks and  
83 hardware-associated upgrading processes. Given their powerful parallel capabilities  
84 and wide range of usability, GPU server-based variant calling solutions are gaining  
85 more attention<sup>16–18</sup>. Compared to the GATK pipeline, NVIDIA Clara Parabricks, a  
86 GPU-accelerated computational genomics application framework, can reduce the  
87 running time of 30X WGS germline analysis on an 8\*A100 GPU server by 60-fold<sup>19</sup>.  
88 With a large number of processing units and high memory bandwidth, GPUs can also  
89 significantly increase deep learning-based variant callers' training and inference speed  
90 <sup>20,21</sup>.

91 The present study aimed to comprehensively evaluate a GPU-based variant caller,  
92 BaseNumber (SaileGene Inc, Beijing). Using gold standard samples and an in-house  
93 WGS dataset, we compared the BaseNumber to the GATK-based pipeline in terms of  
94 efficiency, accuracy, reproducibility, scalability, and energy consumption in germline  
95 variant calling on human genome data, providing an overall assessment of this high-  
96 throughput variant identification tool.

## 97 **2 Methods**

### 98 **2.1 WGS Data Preparation**

99 The WGS data for this study were from four sources: 1) seven standards (HG001 ~  
100 HG007) from the Genome in a Bottle (GIAB) project hosted by the National Institute  
101 of Standards and Technology<sup>22</sup>, 2) four standards (D5, D6, F7, and M8) provided by  
102 the Golden Standard of China Genome (GSCG), 3) 24 additional WGS resequenced  
103 on the GSCG cell lines in six different DNA sequencing facilities using different  
104 sequencing platforms (named the Retested\_GSCG samples), and 4) 100 in-house  
105 samples from the China Deafness Genetics Consortium (CDGC) project. Reference  
106 BAM and VCF files were obtained from GIAB and GSCG. For the in-house data, 100  
107 samples were randomly selected from 1085 subjects of the CDGC project that  
108 underwent WGS, representing realistic experimental conditions and outputs. As  
109 shown in Table 1 and Table S1, the four GSCG standards had high coverage of  
110 144.2X to 146.5X, while all seven GIAB standards had an ultra-high data volume of  
111 247.1X. The average coverage of the CDGC and Retested\_GSCG samples was  
112 48.2X, ranging from 31.0X to 118.0X.

### 113 **2.2 Variant Calling Pipelines and Testing Environment**

114 The GATK pipeline followed the widely implemented protocol of best practice that  
115 includes Fastp (v0.20.1) for FASTQ preprocessing<sup>23</sup>, BWA (v0.7.17) for read  
116 alignment<sup>24</sup>, Samtools (v1.9) for BAM sorting<sup>25</sup>, Picard (v2.23.5) for duplication

117 removal<sup>26</sup>, and GATK (v4.1.7.0) for base quality score recalibration (BQSR) and  
118 variant calling<sup>27</sup>. The BaseNumber pipeline included Saile-Aligner (SLA, v1.0.3) for  
119 FASTQ and BAM processing and Saile-Caller (SLC, v1.0.3) for BQSR and variant  
120 calling. Docker (v19.03) was used to pack and implement the pipelines on each  
121 server. Genome Reference Consortium Human Build 37 (GRCh37) was used as the  
122 reference genome, and corresponding reference files were downloaded from the  
123 GATK Resource Bundle (<ftp.broadinstitute.org/bundle/>).

124 Two GPU servers and ten high-performance CPU servers were utilized for this study.  
125 The detailed configurations of the servers are shown in Table 2. In brief, the  
126 GPU\_Generic server was configured with 8X NVIDIA Tesla V100 cards and more  
127 CPU cores, aiming for common GPU-related computing, such as training of deep  
128 learning models. The configuration of the GPU\_BaseNumber server was optimized  
129 for the BaseNumber pipeline, including improved RAM size and the cooling system.  
130 The high-performance CPU (CPU\_Generic) servers were designed for generic CPU-  
131 related analytic tasks, containing a cost-effectively balanced hardware configuration  
132 and no GPU acceleration. UNI-T UT230A-II power sockets were used to monitor the  
133 instantaneous power consumption and total power consumption. All testing data were  
134 stored on the network-attached storage system, connected with the GPU and CPU  
135 servers through the InfiniBand network.

### 136 **2.3 Comprehensive Assessments**

137 The complete evaluation consisted of seven assessments, as illustrated in Figure 1.  
138 Gold standard data were applied in the accuracy and scalability assessments, while  
139 locally yielded data were used to evaluate the efficiency, consistency, reproducibility,  
140 and energy consumption. The accuracy of BaseNumber was assessed using the  
141 sequencing data of the Retested\_GSCG samples and GIAB standards on the  
142 GPU\_BaseNumber server. The variant calling output of BaseNumber was compared  
143 to the reference VCF files using hap.py from Haplotype Comparison Tools to  
144 calculate the precision, recall, and F1 scores of the whole genome and to give  
145 genomic regions<sup>28</sup>. Seqtk was used to randomly select a given proportion of the reads  
146 from the original FASTQ files to constitute new FASTQ files with targeted depth<sup>29</sup>.  
147 Raw and downsampled WGS data of GSCG and GIAB standards were used to  
148 evaluate the correlation between accuracy and sequencing depth. In addition,  
149 downsampled WGS data were applied to assess the correlation between the analysis  
150 time and sequencing depth. WGS data of 100 CDGC samples were analyzed on the  
151 GPU\_BaseNumber and CPU\_Generic servers using the BaseNumber and GATK  
152 pipelines, respectively, to compare the efficiency. The wall clock time was recorded  
153 for each sample and used to calculate the acceleration ratio. BAM and VCF outputs  
154 were compared to assess the reproducibility of BaseNumber. Variant calling results of  
155 BaseNumber and GATK were compared using HAP for consistency. The analysis  
156 times of raw and downsampled HG001 by GPU\_BaseNumber and GPU\_Generic  
157 servers were compared to assess the hardware impact. The GPU\_Generic server was  
158 configured to assess the influence of the GPU card number on the analysis time.

## 159 **3 Results**

### 160 **3.1 Accurate Variant Identification using BaseNumber on GSCG and GIAB** 161 **Standards**

162 We first evaluated the accuracy of the BaseNumber germline variant detection using  
163 24 Retested\_GSCG samples and seven GIAB standards (Figure 2A and Table S2).  
164 The GIAB standards were downsampled to 30X depth since the raw coverage was  
165 excessively high. The reference VCF files of the GSCG and GIAB standards were  
166 compared to the BaseNumber variant calling results. For the 31 tested samples, the  
167 mean precision rate of all variants was 99.32% (with a standard deviation [SD] of  
168  $\pm 0.21\%$ ), the mean recall rate was 99.86% ( $\pm 0.08\%$ ), and the mean F1 was 99.59%  
169 ( $\pm 0.10\%$ ). Single nucleotide variant (SNV) calling (mean F1 99.63 $\pm 0.09\%$ )  
170 performed slightly better than indel calling (mean F1 99.05 $\pm 0.09\%$ ). We compared the  
171 data of the Retested\_GSCG samples generated by DNBSEQ-T7 to the same samples  
172 generated by NovaSeq 6000. The F1 values were 99.65 $\pm 0.03\%$  for the DNBSEQ-T7  
173 dataset and 99.56 $\pm 0.13\%$  for the NovaSeq 6000 dataset, which were comparable. The  
174 precision (99.26 $\pm 0.13\%$ ) and recall (99.80 $\pm 0.09\%$ ) rates of the downsampled GIAB  
175 samples were also satisfactory.

176 The correlation between the WGS depth and the variant calling accuracy by  
177 BaseNumber was measured with the GSCG standards and GIAB HG001. The reads  
178 were randomly retrieved from the FASTQ files with a given proportion to construct a  
179 set of simulated WGS samples with gradationally downsampled depth from 10X to  
180 300X (Table S3). The mean F1 was 97.10% for 10X simulated GSCG samples and  
181 94.71% for 10X HG001, but the performance quickly improved with increasing depth  
182 (Figure 2B). When the coverage was over 30X, the analysis accuracy reached a stable  
183 level even with extremely high coverage of 300X.

### 184 **3.2 Boosted Variant Calling Efficiency of BaseNumber**

185 The wall clock time spent on each step of variant calling on the data of 100 CDGC  
186 samples was recorded for the BaseNumber pipeline on the GPU\_BaseNumber server  
187 and the GATK pipeline on the CPU\_Generic servers. As shown in Figure 3A, the  
188 mean total analysis time of BaseNumber (from FASTQ to VCF) for each CDGC  
189 sample was 23.35 $\pm 4.75$  minutes (ranging from 19.92 to 50.4 minutes). More  
190 specifically, the steps of read alignment and polymerase chain reaction (PCR)  
191 duplication removal (carried out by SLA) required 10.64 $\pm 2.68$  minutes, the time for  
192 generating recalibrated BAM files (carried out by SLC) was 11.59 $\pm 2.12$  minutes, and  
193 the variant calling step that outputs gVCF and VCF files (carried out by SLC)  
194 required only 1.11 $\pm 0.28$  minutes. As a comparison, a mean of 5018.35 $\pm 1330.97$   
195 minutes (ranging from 3598 to 10179 minutes) was required for GATK to generate  
196 final variant calling results from raw sequencing read data of the same sample set.  
197 BaseNumber greatly (215.33 $\pm 37.45$  times) improved the computational efficiency of  
198 germline variant calling. The analysis time was linearly correlated with the coverage



199 of the WGS data, which ranged from 31X to 118X. The  $R^2$  of the correlations was  
200 0.863 for BaseNumber and 0.719 for GATK.

201 We further explored the effect of data volume on the analysis time using the GSCG  
202 and GIAB standards. The FASTQ files of four GSCG standards were randomly  
203 downsampled from 10X to 120X. The GIAB HG001 FASTQ was sampled from 10X  
204 to 300X. As shown in Figure 3B, a linear correlation between the depth and analysis  
205 time was observed even when the coverage was ultrahigh (300X). Interestingly, the  $R^2$   
206 for SLA was greater than that for SLC, suggesting a stronger linear correlation  
207 between the amount of input data and SLA analysis time.

### 208 **3.3 Reproducibility of BaseNumber Outputs**

209 We repeatedly analyzed the CDGC samples to evaluate the reproducibility of  
210 BaseNumber. The message-digest algorithm 5 (MD5) values of output files  
211 corresponding to the same samples were identical for all three rounds of analysis  
212 (Table S4), indicating complete, reproducible, and robust outcomes with  
213 BaseNumber. Next, we compared the variant calling results for the CDGC samples  
214 from the BaseNumber and GATK pipelines. The precision, recall, and F1 scores were  
215 calculated using the corresponding GATK result as the reference for each sample  
216 (Figure 3C). The F1 of all variants was on average  $99.69 \pm 0.19\%$ , ranging from  
217  $98.62\%$  to  $99.80\%$ . The mean F1 scores were  $99.71 \pm 0.19\%$  for SNVs and  
218  $99.55 \pm 0.19\%$  for indels, representing highly similar variant calling outcomes between  
219 BaseNumber and GATK on the same sequencing data. We also compared the F1  
220 scores of different sequencing facilities and platforms, and no significant differences  
221 were observed.

### 222 **3.4 Impact of the GPU Configuration on the BaseNumber Performance**

223 The configuration of the GPU cards directly impacted the speed of the BaseNumber  
224 algorithms. We assessed the running time of the BaseNumber variant caller with  
225 different numbers of GPU cards. Downsampled HG001 data (from 10X to 180X)  
226 were used for the analysis on the GPU\_Generic server configured with eight NVIDIA  
227 Tesla V100 cards. As expected, the total analysis time increased as the data volume  
228 increased or the number of GPU cards decreased (Figure 4, Table S5). We further  
229 dissected the total time into the times for SLA and SLC. The SLA analysis time,  
230 which was responsible for alignment and variant calling, continued to decrease as the  
231 number of GPU cards increased, while the analysis time of SLC, which included the  
232 BQSR process, insignificantly decreased from four GPU cards to eight cards and even  
233 increased for HG001\_10X. This may be due to IO bottlenecks triggered by SLC when  
234 outputting recalibrated BAM files.

### 235 **3.5 Energy Consumption**

236 A comparison of the energy consumption of the BaseNumber and GATK pipelines  
237 was conducted. Electricity usage was recorded for the variant calling processes of 100

238 CDGC samples; this was repeated three times. In total, 224 kilowatt-hours were  
239 consumed in analyzing these samples three times using BaseNumber on the  
240 GPU\_BaseNumber server. For each CDGC sample, the energy consumption of  
241 BaseNumber was estimated to be 0.746 kilowatt-hours. Similarly, we estimated that  
242 the average power consumption per CDGC sample was 33.5 kilowatt-hours for the  
243 GATK pipeline on CPU\_Generic servers. As a result, the "GPU sever + BaseNumber"  
244 solution used 44.9 times less power than the "CPU server + GATK".

#### 245 **4 Discussion**

246 This study comprehensively evaluated the performance of the GPU-accelerated  
247 BaseNumber pipeline in germline variant calling. BaseNumber achieved an average  
248 F1 value of 99.59% on the gold standards from the GSCG and GIAB projects, in  
249 which the accuracy was similar to that of the best-performing variant callers in the  
250 PrecisionFDA Consistency Challenge  
251 (<https://precision.fda.gov/challenges/consistency/results>). Additionally, we compared  
252 the results of BaseNumber and GATK on CDGC WGS samples with a mean coverage  
253 of 48X, and the consistency was remarkable, with an average F1 value of 99.68%  
254 using the GATK results as the references. More importantly, we observed an average  
255 of 23 minutes taken to analyze the 48X WGS sample using the GPU server, which  
256 was more than 200 times faster than the BWA + GATK pipeline. These results show  
257 that the BaseNumber pipeline is an attractive alternative to the commonly used BWA  
258 + GATK pipeline.

259 With the emergence of million-sample-sized WGS projects, ultrafast and highly  
260 accurate variant calling is essential for further genome analysis. For example, the  
261 Genome Aggregation Database (gnomAD) (v3.1) integrated 76,156 WGS samples  
262 worldwide, providing invaluable population genetic information to support a wide  
263 range of research and clinical applications<sup>30</sup>. If the average depth of gnomAD  
264 samples is 30X, BaseNumber requires only 90 days to complete the variant calling  
265 from raw sequencing data using ten servers configured with GPU. Clearly, leveraging  
266 the capabilities of BaseNumber can result in significant savings in terms of  
267 investment in server hardware, room space, support facilities, and staffing. The  
268 reduced analysis time also significantly decreases energy consumption and computing  
269 costs. Based on the price of GPU servers and power consumption per sample, we  
270 estimate that if 100,000 30X WGS samples are processed in four years, the analysis  
271 cost per sample can be less than \$0.40. The analysis strategy can be flexibly and  
272 finely formulated, facilitated by the significant improvement in efficiency and cost  
273 control. Reference genomes, read mapping, variant calling, and quality control  
274 algorithms can be tuned at a relatively low sunk cost even as a large cohort study  
275 progresses. Moreover, the agile and easy implementation of GPU-based software  
276 ensures the operability and timeliness of the pipeline adjustments. During idle time,  
277 the GPU server can be used for other high-performance parallel computing tasks, such  
278 as deep learning model training, image recognition, and natural language processing,  
279 to maximize the usage of the server resources.

280 The principle of BaseNumber's germline variant calling algorithm is similar to that of  
281 GATK HaplotypeCaller. It reassembles the aligned reads in the active regions where  
282 variants may be present to capture SNVs and indels accurately. However, it is  
283 challenging to accelerate variant calling on the GPU-based platform. Simple  
284 transplantation of HaplotypeCaller to the GPU was reportedly inefficient, with only a  
285 2.3 times speedup<sup>16,31</sup>. The alignment and variant calling processes were relatively  
286 coarse-grained; therefore, the parallelism of corresponding algorithms required  
287 reconstruction to support the concurrent operation of thousands of GPU cores. In  
288 addition, considering the ultrahigh bandwidth of graphics memory, a dedicated  
289 management system needs to be developed to modulate the computing flow. I/O  
290 operations have to be optimized to eliminate bottlenecks of the system. With these  
291 improvements, BaseNumber achieves high-throughput parallel acceleration and  
292 analysis efficiency. Similar performance was claimed for the NVIDIA Clara  
293 Parabricks; reportedly, its GPU caller can complete variant analysis of a 30X human  
294 WGS in 22 minutes on a DGX A100 server, but Parabricks has not yet revealed  
295 detailed information of the evaluation<sup>19</sup>.

296 This study focused on the WGS germline variant calling scenario, while BaseNumber  
297 is also a proper solution for analyzing ultrahigh coverage sequencing data. Given the  
298 high computational capabilities, BaseNumber had a shorter running time, requiring no  
299 downsampling processes such as that implemented in GATK to handle excessive  
300 coverage regions. As a result, BaseNumber yielded excellent reproducibility in that  
301 the outputs for the same data were identical. With the support of the high processing  
302 speed of GPUs, some other time-consuming methods, such as graph alignment<sup>32</sup> and  
303 genotype imputation<sup>33</sup>, could be applied to boost the accuracy of read mapping and  
304 variant calling. BaseNumber might benefit from the adoption of advanced storage  
305 systems, such as parallel file systems and flash solid-state drive (SSD) network-  
306 attached storage (NAS), and enhanced I/O capabilities to further boost the  
307 performance, which became a relative bottleneck in this evaluation.

308 The results should be interpreted considering several limitations. First, we were  
309 unable to compare GATK and BaseNumber on the same GPU server due to the large  
310 size of samples for the evaluation and slow process of the GATK pipeline. The  
311 GPU\_BaseNumber server was configured with better I/O components (such as more  
312 RAM and a larger SSD) than the CPU\_Generic servers. The performance of GATK  
313 might be slightly improved on the GPU\_BaseNumber server, but it is unlikely to  
314 change the main results. Second, due to the nature of short-read sequencing<sup>24,34,35</sup>,  
315 BaseNumber, like other variant callers based on NGS data, suffered a significant  
316 decrease in variant calling accuracy in low mappability regions. Such regions include  
317 segmental duplications (segdups)<sup>36</sup>, tandem repeats<sup>37</sup>, variable, diversity and joining  
318 (VDJ) recombination regions<sup>38</sup>, and other genomic regions enriched with repetitive  
319 sequence. Variant identification and interpretation in these regions depend on the  
320 advancement of long-read sequencing techniques.



321 In conclusion, the GPU-based BaseNumber provides a high accuracy and ultrafast  
322 variant calling. It can maximize the efficiency of WGS analysis in large population  
323 studies, improve the utilization of hardware resources, and meet the requirements of  
324 time-sensitive tests, such as clinical WGS genetic diagnosis. Moreover, BaseNumber  
325 shed light on the use of GPUs to accelerate other bioinformatic pipelines based on a  
326 similar concept and design ideas to satisfy the increasing demands of multi-omics data  
327 analysis, providing powerful support for future individual precision medicine.

## 328 **5 Conflict of Interest**

329 The authors declare that the research was conducted in the absence of any commercial  
330 or financial relationships that could be construed as a potential conflict of interest.

## 331 **6 Author Contributions**

332 F. B., H. L., and Q. Z. conceived the study. H. L. wrote the manuscript, with the  
333 contributions by F. B.. Q. Z. organized and performed the evaluation, with  
334 contributions by F. B. and H. L.

## 335 **7 Funding**

336 This work was supported by the 1·3·5 project for disciplines of excellence, West  
337 China Hospital, Sichuan University, and the project No. 82171836 of National  
338 Natural Science Foundation of China.

## 339 **8 Acknowledgments**

340 We are grateful to Prof. Leming Shi provided the GSCG gold standard data.

## 341 **9 References**

- 342 1. Waterston, R. H., Lander, E. S. & Sulston, J. E. On the sequencing of the human  
343 genome. *Proceedings of the National Academy of Sciences* **99**, 3712–3716 (2002).
- 344 2. Check Hayden, E. Technology: the \$1,000 genome. *Nature News* **507**, 294 (2014).
- 345 3. Lappalainen, T., Scott, A. J., Brandt, M. & Hall, I. M. Genomic Analysis in the  
346 Age of Human Genome Sequencing. *Cell* **177**, 70–84 (2019).
- 347 4. Siva, N. UK gears up to decode 100,000 genomes from NHS patients. *Lancet* **385**,  
348 103–104 (2015).
- 349 5. All of Us Research Program Investigators *et al.* The 'All of Us' Research Program.  
350 *N Engl J Med* **381**, 668–676 (2019).
- 351 6. Cao, Y. *et al.* The ChinaMAP analytics of deep whole genome sequences in  
352 10,588 individuals. *Cell research* **30**, 717–731 (2020).

- 353 7. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using  
354 next-generation DNA sequencing data. *Nature genetics* **43**, 491–498 (2011).
- 355 8. Koboldt, D. C. *et al.* VarScan 2: somatic mutation and copy number alteration  
356 discovery in cancer by exome sequencing. *Genome research* **22**, 568–576 (2012).
- 357 9. Goyal, A. *et al.* Ultra-fast next generation human genome sequencing data  
358 processing using DRAGENTM bio-IT processor for precision medicine. *Open*  
359 *Journal of Genetics* **7**, 9–19 (2017).
- 360 10. Herzeel, C., Costanza, P., Decap, D., Fostier, J. & Verachtert, W. elPrep 4: A  
361 multithreaded framework for sequence analysis. *PLoS One* **14**, e0209523 (2019).
- 362 11. Kelly, B. J. *et al.* Churchill: an ultra-fast, deterministic, highly scalable and  
363 balanced parallelization strategy for the discovery of human genetic variation in  
364 clinical and population-scale genomics. *Genome biology* **16**, 1–14 (2015).
- 365 12. Kathiresan, N. *et al.* Accelerating next generation sequencing data analysis with  
366 system level optimizations. *Scientific reports* **7**, 1–11 (2017).
- 367 13. Freed, D., Aldana, R., Weber, J. A. & Edwards, J. S. The Sentieon Genomics  
368 Tools-A fast and accurate solution to variant calling from next-generation  
369 sequence data. *BioRxiv* 115717 (2017).
- 370 14. Kendig, K. I. *et al.* Sentieon DNaseq Variant Calling Workflow Demonstrates  
371 Strong Computational Performance and Accuracy. *Front Genet* **10**, 736 (2019).
- 372 15. *Illumina DRAGEN Bio-IT Platform 3.7 User Guide*. (Illumina, 2020).
- 373 16. Ren, S., Ahmed, N., Bertels, K. & Al-Ars, Z. GPU accelerated sequence  
374 alignment with traceback for GATK HaplotypeCaller. *BMC genomics* **20**, 103–  
375 116 (2019).
- 376 17. Ren, S., Bertels, K. & Al-Ars, Z. Efficient acceleration of the pair-hmms forward  
377 algorithm for gatk haplotypecaller on graphics processing units. *Evolutionary*  
378 *Bioinformatics* **14**, 1176934318760543 (2018).
- 379 18. Wang, J., Xie, X. & Cong, J. Communication optimization on GPU: A case study  
380 of sequence alignment algorithms. in *2017 IEEE International Parallel and*  
381 *Distributed Processing Symposium (IPDPS)* 72–81 (IEEE, 2017).
- 382 19. Braunstein, V. & Burnett, G. GPU-Accelerated Tools Added to NVIDIA Clara  
383 Parabricks v3.6 for Cancer and Germline Analyses. *GPU-Accelerated Tools*  
384 *Added to NVIDIA Clara Parabricks v3.6 for Cancer and Germline Analyses*  
385 [https://developer.nvidia.com/blog/gpu-accelerated-tools-added-to-nvidia-clara-](https://developer.nvidia.com/blog/gpu-accelerated-tools-added-to-nvidia-clara-parabricks-v3-6-for-cancer-and-germline-analyses/)  
386 [parabricks-v3-6-for-cancer-and-germline-analyses/](https://developer.nvidia.com/blog/gpu-accelerated-tools-added-to-nvidia-clara-parabricks-v3-6-for-cancer-and-germline-analyses/) (2021).

- 387 20. Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural  
388 networks. *Nature biotechnology* **36**, 983–987 (2018).
- 389 21. Luo, R., Sedlazeck, F. J., Lam, T.-W. & Schatz, M. C. A multi-task convolutional  
390 deep neural network for variant calling in single molecule sequencing. *Nature*  
391 *communications* **10**, 1–11 (2019).
- 392 22. Zook, J. M. *et al.* An open resource for accurately benchmarking small variant  
393 and reference calls. *Nature biotechnology* **37**, 561–566 (2019).
- 394 23. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ  
395 preprocessor. *Bioinformatics* **34**, i884–i890 (2018).
- 396 24. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-  
397 MEM. *arXiv preprint arXiv:1303.3997* (2013).
- 398 25. Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics*  
399 **25**, 2078–2079 (2009).
- 400 26. Picard toolkit. *Broad Institute, GitHub repository* (2019).
- 401 27. Van der Auwera, G. A. *et al.* From FastQ data to high-confidence variant calls: the  
402 genome analysis toolkit best practices pipeline. *Current protocols in*  
403 *bioinformatics* **43**, 11–10 (2013).
- 404 28. Krusche, P. *et al.* Best practices for benchmarking germline small-variant calls in  
405 human genomes. *Nature biotechnology* **37**, 555–560 (2019).
- 406 29. Li, H. & others. Seqtk: a fast and lightweight tool for processing FASTA or  
407 FASTQ sequences. (2013).
- 408 30. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from  
409 variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
- 410 31. Franke, K. R. & Crowgey, E. L. Accelerating next generation sequencing data  
411 analysis: an evaluation of optimized best practices for Genome Analysis Toolkit  
412 algorithms. *Genomics & informatics* **18**, (2020).
- 413 32. Garrison, E. *et al.* Variation graph toolkit improves read mapping by representing  
414 genetic variation in the reference. *Nature biotechnology* **36**, 875–879 (2018).
- 415 33. Browning, B. L., Zhou, Y. & Browning, S. R. A one-penny imputed genome from  
416 next-generation reference panels. *The American Journal of Human Genetics* **103**,  
417 338–348 (2018).
- 418 34. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2.  
419 *Nature methods* **9**, 357–359 (2012).

- 420 35. Marco-Sola, S., Sammeth, M., Guigó, R. & Ribeca, P. The GEM mapper: fast,  
421 accurate and versatile alignment by filtration. *Nature methods* **9**, 1185–1188  
422 (2012).
- 423 36. Bailey, J. A. *et al.* Recent segmental duplications in the human genome. *Science*  
424 **297**, 1003–1007 (2002).
- 425 37. Gemayel, R., Vincens, M. D., Legendre, M. & Verstrepen, K. J. Variable tandem  
426 repeats accelerate evolution of coding and regulatory sequences. *Annual review of*  
427 *genetics* **44**, 445–477 (2010).
- 428 38. Schatz, D. G. & Swanson, P. C. V (D) J recombination: mechanisms of initiation.  
429 *Annual review of genetics* **45**, 167–202 (2011).

430

431

432

433

434 **Table 1.** Summary of WGS samples for the evaluation

| Category      | Sample size | Coverage | Evaluation content   |
|---------------|-------------|----------|--|
| CDGC          | 100         | 48.89X   | Efficiency, consistency, reproducibility, energy consumption |
| GIAB          | 7           | 247.14X  | Accuracy, scalability  |
| GSCG          | 4           | 145.05X  | Accuracy, scalability  |
| Retested_GSCG | 24          | 45.56X   | Accuracy   |

435

436



437 **Table 2.** Hardware configuration of testing servers

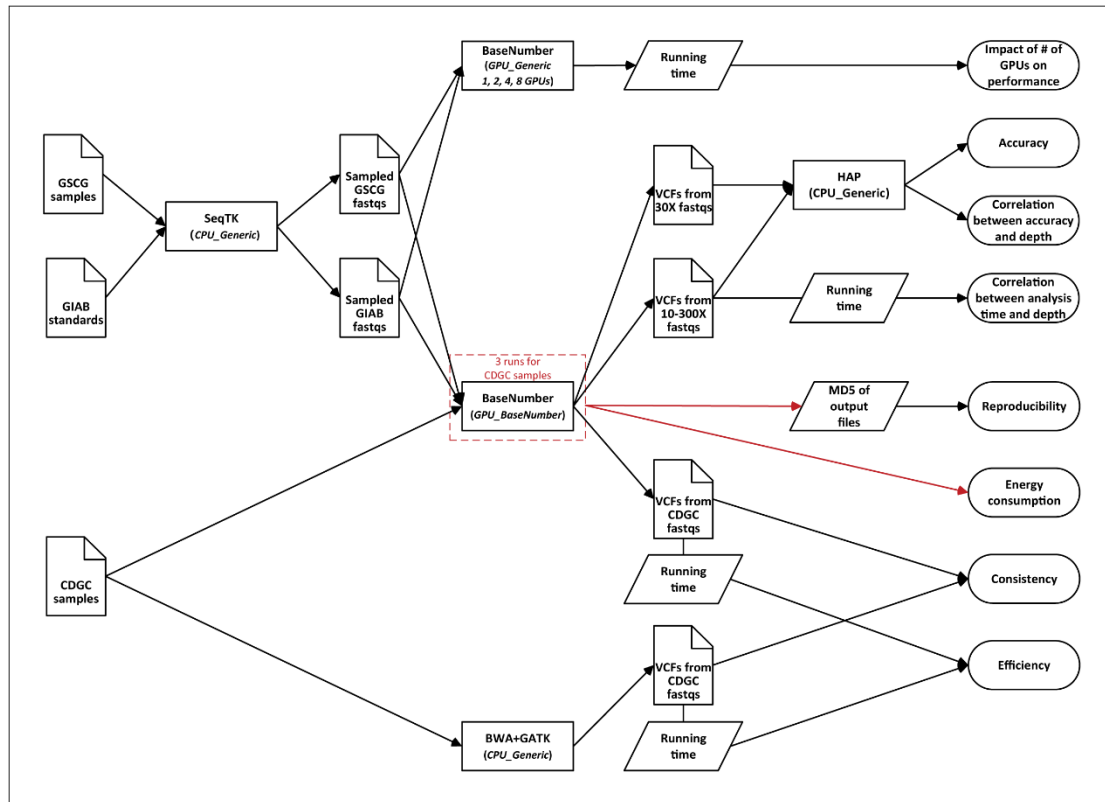
| Server name    | CPU   | Memory    | GPU                                | Hard disk                             |
|----------------|---|-----------|------------------------------------|---------------------------------------|
| GPU_BaseNumber | 2X Intel 6226R processor (16 cores, 32 threads) | DDR4 1T   | 4* NVIDIA Geforce RTX 2080 Ti 11GB | 2X 3.2TB NVMe; 2X 8T 7200 RPM HDDs    |
| CPU_Generic    | 2X Intel 6248 processor (20 cores, 40 threads)  | DDR4 125G | N/A                                | 250T NAS                              |
| GPU_Generic    | 2X Intel 6240R processor (48 cores, 96 threads) | DDR4 256G | 8 * NVIDIA Tesla V100 32GB         | 2X 3.2TB NVMe 3T; 2X 8T 7200 RPM HDDs |

438

439

440

441

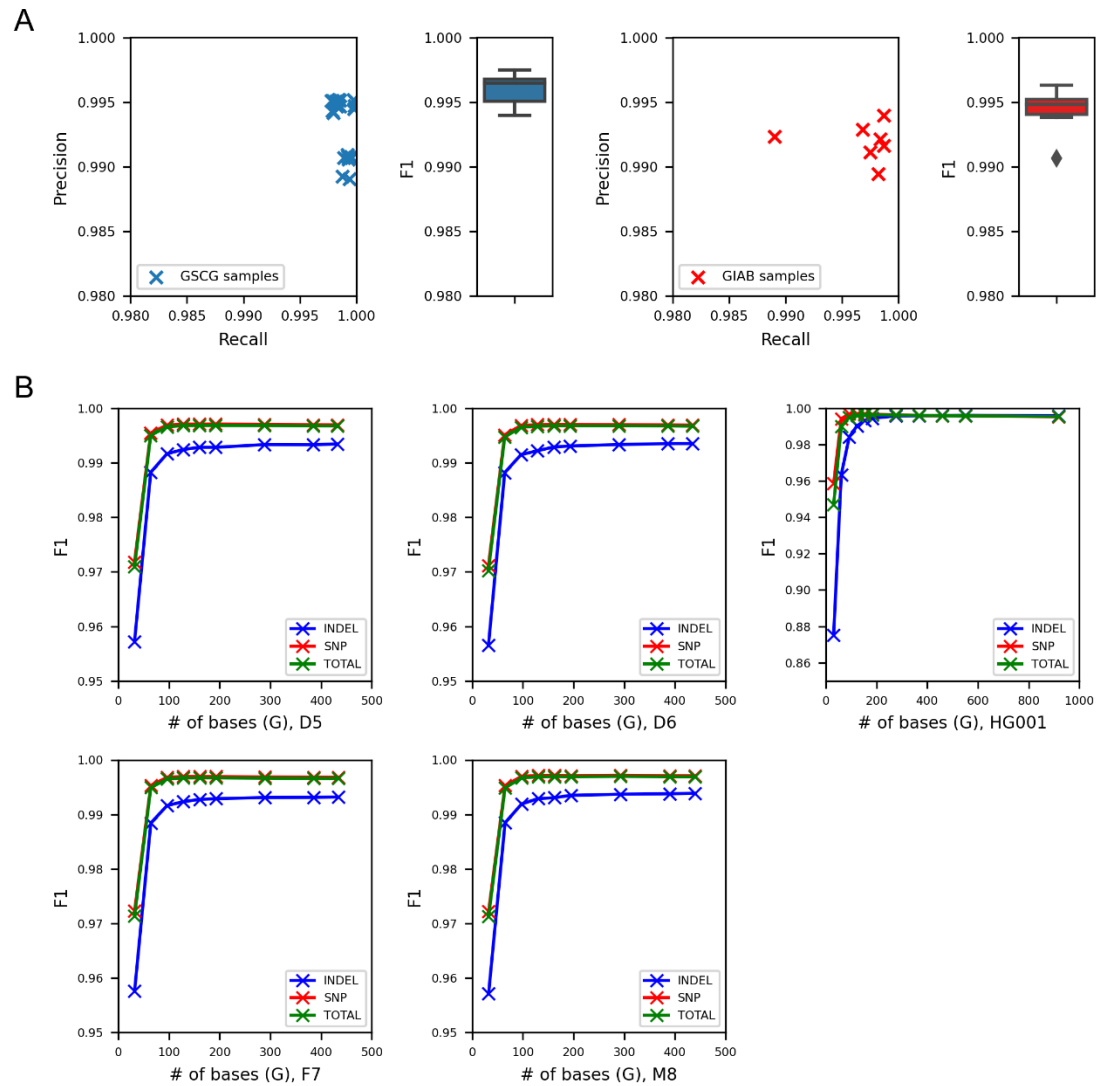


442

443

**Figure 1.** Schema of the study design.

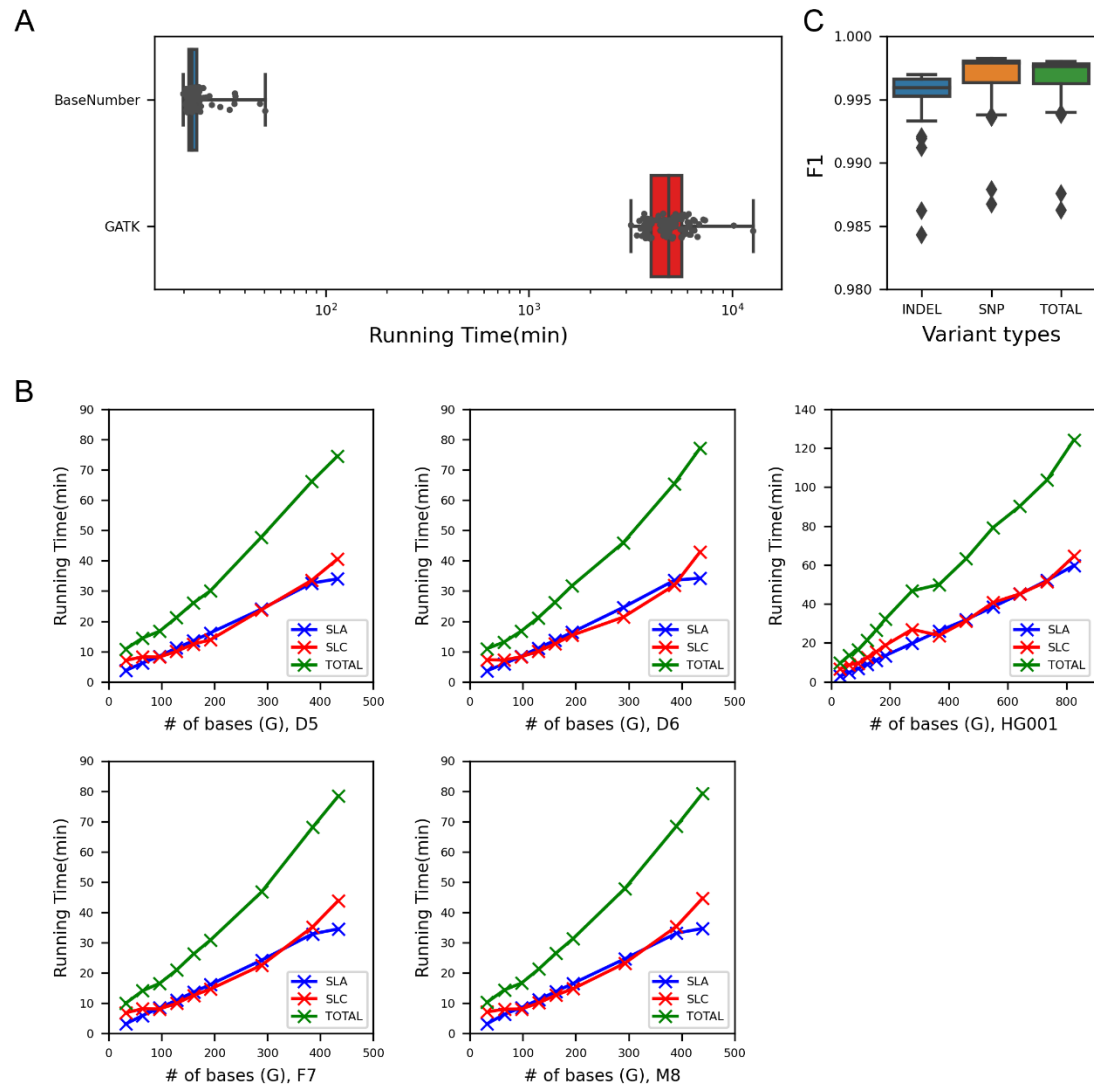
444



445

446 **Figure 2.** Accuracy evaluation of the BaseNumber using gold standard samples. A)  
447 Precision, recall, and F1 score of the BaseNumber on 24 resequenced GSCG and  
448 seven GIAB samples. B) F1 score of the BaseNumber rapidly increased along with  
449 the sequencing depth.

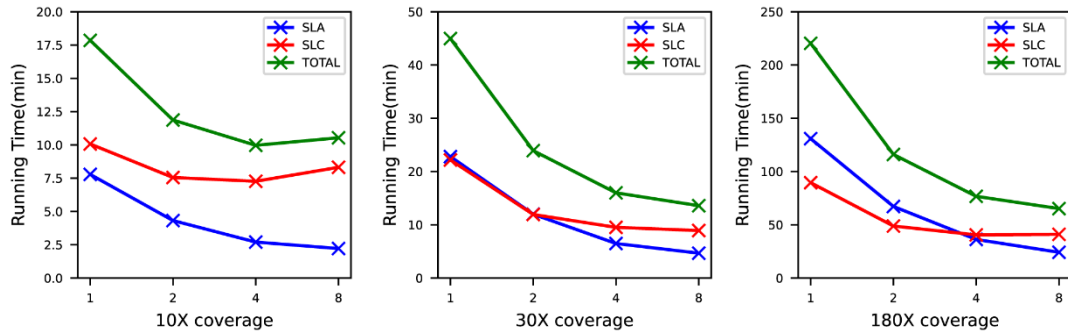
450



451

452 **Figure 3.** Efficiency of the BaseNumber variant calling process. A) Comparison of  
453 analysis time between the BaseNumber and GATK process using 100 CDGC WGS  
454 samples. B) Correlation between the BaseNumber analysis time and sequencing  
455 depth. C) Variant calling results of the BaseNumber and GATK process were highly  
456 consistent

457



458

459 **Figure 4.** Analysis time of BaseNumber was correlated with GPU configuration and  
460 sequencing depth

461

462

463



464 **Supplementary Appendix**

465 Table S1. Detailed information of the testing samples

466 Table S2. Recall, precision, and F1 scores of BaseNumber in the Retested\_GSCG and  
467 GIAB samples.

468 Table S3. Correlation between the accuracy of BaseNumber and sequencing depth

469 Table S4. MD5 values of BaseNumber output files from three rounds of analysis

470 Table S5. Effects of GPU card counts and sequencing depth on the analysis time