



Accelerated somatic mutation calling for whole-genome and whole-exome sequencing data from heterogenous tumor samples

Shuangxi Ji, Wenyi Wang, Tong Zhu, et al.

Genome Res. published online April 8, 2024

Access the most recent version at doi:[10.1101/gr.278456.123](https://doi.org/10.1101/gr.278456.123)

P<P	Published online April 8, 2024 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at http://creativecommons.org/licenses/by/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 Accelerated somatic mutation calling for whole-genome and whole-exome sequencing data
2 from heterogenous tumor samples

3

4 Shuangxi Ji¹, Tong Zhu², Ankit Sethia², Wenyi Wang^{1*}

5 ¹Department of Bioinformatics and Computational Biology, The University of Texas MD
6 Anderson Cancer Center, Houston, TX, USA.

7 ²NVIDIA Corporation, Santa Clara, CA, USA.

8

9 *Correspondence to: Wenyi Wang, wwang7@mdanderson.org

10

11 Running title: MuSE with multi-step parallelization

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27 **Abstract**

28 Accurate detection of somatic mutations in DNA sequencing data is a fundamental prerequisite
29 for cancer research. Previous analytical challenge was overcome by consensus mutation calling
30 from four to five popular callers. This, however, increases the already nontrivial computing time
31 from individual callers. Here, we launch MuSE 2, powered by multi-step parallelization and
32 efficient memory allocation, to resolve the computing time bottleneck. MuSE 2 speeds up 50
33 times than MuSE 1 and 8-80 times than other popular callers. Our benchmark study suggests
34 combining MuSE 2 and the recently accelerated Strelka2 achieves high efficiency and accuracy
35 in analyzing large cancer genomic datasets.

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53 **Introduction**

54 Cancer arises and evolves by accumulating various types of genetic alterations, such as single
55 nucleotide variation (SNV), copy number alteration (CNA) and structural variation (SV). The
56 high-throughput sequencing (HTS) technology has revolutionized the way we look at many
57 human diseases, particularly cancer. With its constantly improved capacity and reduced cost,
58 HTS is enabling investigations of genetic alterations within large human patient cohorts, hence
59 advancing both basic and translational cancer research. Many computational tools have been
60 developed for calling somatic variants (Xu 2018), which typically require, as input, whole-
61 genome sequencing (WGS) or whole-exome sequencing (WES) data from the tumor tissue, as
62 well as from the blood of the patient to serve as the germline control. WGS provides the most
63 comprehensive coverage to sequence both protein-coding and non-coding regions across the
64 entire genome; whereas WES provides an efficient alternative to WGS by targeting only protein-
65 coding regions that accounts for 1-2% of the genome (Alfares et al. 2018), hence achieving both
66 higher read depth (Sun et al. 2021; Barbitoff et al. 2020) and lower sequencing cost.

67 We previously launched MuSE 1 (Fan et al. 2016), a statistical approach for somatic
68 mutation calling, where we introduced a combination of nucleotide base-specific Markov
69 substitution model for molecular evolution and a tumor sample-specific error model to estimate
70 tier-based cutoffs for selecting SNVs. Due to its high sensitivity and specificity, especially for
71 calling subclonal SNVs, MuSE 1 was adopted in multiple pipelines, including as a major
72 contributing caller to reach final consensus calls by the Cancer Genome Atlas (TCGA)
73 PanCanAtlas project (Ellrott et al. 2018) across ~13,000 tumor samples, and the International
74 Cancer Genome Consortium Pan-Cancer Analysis of Whole Genomes (ICGC-PCAWG)
75 initiative (The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium 2020) across
76 ~2,700 tumor samples.

77 One major limitation of MuSE 1, like many other mutation callers (Larson et al. 2012; Cibulskis
78 et al. 2013; Koboldt et al. 2012), is the computational speed. It takes 2-3 days to finish running

79 the WGS data of a tumor-normal pair on a typical Linux server with an Intel Xeon processor and
80 more than 100 gigabytes (GB) random access memory (RAM), which explains the commonly
81 seen long wait-times for completing mutation calling before any downstream analysis in large
82 patient cohort studies. Here, we present MuSE 2, which maintains the same input, output and
83 mathematical model as MuSE 1, but accelerates significantly for both WES and WGS data by
84 adopting a new algorithmic programming backbone. MuSE 2 employs a multithreaded producer-
85 consumer model and the OpenMP library for parallel computing, including parsing and
86 uncompressing reads from binary sequence alignment/map formatted (BAM) files, detecting and
87 filtering variants, and writing output. It is also optimized by adopting a more efficient memory
88 allocator. In this paper, we have benchmarked the accuracy of MuSE 2 against three somatic
89 mutation callers, i.e., MuTect2 (Cibulskis et al. 2013), SomaticSniper (Larson et al. 2012) and
90 VarScan2 (Koboldt et al. 2012), which are the other highlighted somatic mutation callers in the
91 National Cancer Institute Genomic Data Commons (GDC) DNA-seq analysis pipeline:
92 https://docs.gdc.cancer.gov/Data/Bioinformatics_Pipelines/DNA_Seq_Variant_Calling_Pipeline/,
93 as well as a recently accelerated mutation caller Strelka2 (Kim et al. 2018). We use the
94 consensus mutation calls generated by previous consortial studies with 3-5 un-accelerated
95 callers (Ellrott et al. 2018; The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes
96 Consortium 2020; Craig et al. 2016). Here we demonstrate the improved utility of our new caller
97 using WES data generated from 5 tumor-normal pairs and WGS data generated from 7 tumor-
98 normal pairs.

99

100 **Results**

101 **Overview of approach.** MuSE 2 takes as input the indexed reference genome FASTA file, the
102 BAM format sequencing data from a pair of tumor-normal tissues (**Supplemental Fig. S1**) and
103 the dbSNP (Sherry et al. 2001) variant call format (VCF) file, which is bgzip compressed, tabix
104 (Li 2011) indexed using the same reference genome. Unlike MuSE 1, which can only utilize one

105 core, MuSE 2 takes advantage of the multi-core resources in a modern computer or a
106 computing node for somatic SNV calling from WES/WGS data (**Fig. 1A-B**).

107 Since our benchmarking study requires a large amount of computational resources to cover
108 multiple callers and scenarios, we only include results for WES data from 5 tumor-normal pairs
109 and WGS data from 5 tumor-normal pairs, which are randomly selected and downloaded from
110 the GDC (Genomic Data Commons) data portal and the ICGC data portal, respectively. The
111 sequencing depths from these samples reflect the wide ranges presented in both datasets
112 (**Supplemental Table S1, Supplemental Fig. S2**). We further include WGS data from 2 tumor-
113 normal pairs to evaluate mutation calling performances on newer sequencing platforms.

114 We compare the SNV entries in the output VCF files generated by MuSE 2 with those by
115 MuSE 1 for each patient sample with the same or different number of CPU cores. Since each
116 SNV entry is denoted by one line of string in a VCF file, we compare the strings from both
117 methods line by line. The result shows that all the entries from the two methods are identical
118 with the same number or different number of cores (**Supplemental Fig. S3**).

119 **Accuracy benchmarking for real tumor samples.** We evaluate the performance of MuSE
120 2 and compare to other callers using the consensus SNV calls from TCGA (for WES data) and
121 PCAWG (for WGS data) as truth sets. The truth sets include 168-2,553 somatic SNVs
122 (mean=1,394, median=932) for WES and 3,813-19,081 somatic SNVs (mean=10,146,
123 median=8,073) for WGS data. We first compared among callers whose predecessors
124 contributed to the consensus call: MuSE 2, MuTect2, Somatic Sniper, VarScan2 for TCGA WES,
125 MuSE 2 and MuTect2 for PCAWG WGS. We divided the mutation positions into multiple bins
126 defined by variant allele frequency (VAF), or sequencing read depth, by classes of variant
127 effects or clonality (**Methods**), and calculated precision, recall and F1 score, i.e., the harmonic
128 mean of precision and recall for each bin. For both WES and WGS data, MuSE 2 achieves a
129 higher precision at a similar or higher recall, hence a higher F1 score (**Supplemental Fig. S4A**
130 for TCGA, **Supplemental Fig. S4B** for PCAWG) across all bins of VAF and read depth, and

131 across variant classes. It also achieves a higher recall for subclonal consensus SNV calls from
132 PCAWG WGS (Dentro et al. 2021, **Methods**). We then compared the performance between
133 MuSE 2 and Strelka2 (Kim et al. 2018), which used machine learning and curated data to train a
134 position-specific error score, and was developed after the release of PCAWG consensus calls.
135 Note that MuSE 2 utilizes the same Bayesian Markov model as MuSE 1 to explicitly define an
136 evolutionary process and estimate model-based parameters based on data. Overall, Strelka2
137 performs well in both WES and WGS data. Compared to MuSE 2, its performance is lower in
138 WES (**Fig. 2**) across all bins of VAF, read depth and different variant class. However, in the
139 case of WGS data, it either matches or surpasses MuSE 2 in precision, at a similar recall rate
140 (**Fig. 3A**). On the other hand, Strelka2 still shows a slightly lower recall than MuSE 2 at
141 positions with low VAF (<0.2) or low read depth ($<20\times$).

142 **Accuracy benchmarking for tumor cell line data.** As the TCGA and PCAWG data used
143 above were generated by Illumina HiSeq 2000 more than a decade ago, we further obtained a
144 set of tumor cell line sequencing data that were generated using newer sequencing platforms to
145 compare the caller performances. A cell line tumor-normal pair COLO829/COLO829BL was
146 used to generate WGS data on HiSeq X Ten, called COLO829 Illumina, and NovaSeq, called
147 COLO829 10X, and the consensus calls from three mutation callers were generated through a
148 multi-institutional effort (Craig et al. 2016, **Methods**). Although the sequencing data is newer,
149 the consensus call effort was relatively old and the three callers did not entirely match what was
150 established by the TCGA and PCAWG projects. They include MuTect2's predecessor MuTect
151 (Cibulskis et al. 2013), Strelka2's predecessor Strelka (Saunders et al. 2012) and Seurat
152 (Christoforides et al. 2013). We therefore implemented two strategies to benchmark the
153 accuracy of MuSE 2 and Strelka2. First, we took the consensus calls as the truth set, which
154 includes 35,543 SNVs from two pairs of samples, COLO829 Illumina and COLO829 10X.
155 Second, we put aside the consensus calls and instead used the mutation calls made in the
156 tumor cell line (100%) as the truth set. This truth set then includes 45,853 SNVs from MuSE 2

157 and 44,257 SNVs from Strelka2, respectively called from COLO829 Illumina. We then evaluated
158 how many of these initial calls were recovered in a *in silico* diluted dataset where the tumor cell
159 proportion decreases to 75%, 50%, 25%, 20% and 10% (**Methods**). MuSE 2 and Strelka2 did
160 equally well with both truth sets. MuSE 2 presented a slightly higher recall than Strelka2 at
161 positions with low VAF (<0.2) or in samples with low tumor cell proportion ($\leq 10\%$) (**Fig. 3B, C**).
162 We note that this slight advantage in MuSE 2 is consistently observed across all WGS data,
163 whereas an advantage of higher precision in MuSE 2 is consistently observed in WES data.

164
165 **Speed benchmarking.** We compare the speed of running MuSE 2 against MuSE 1,
166 MuTect2, SomaticSniper, VarScan2 and Strelka2 on a computing cluster. Each method is
167 tested with the number of CPU cores at 1, 5, 10, 20, 28, 40 and 80. All methods are assigned
168 with the same RAM of 50GB for WES data and 150GB for WGS data. The time cost of each
169 method for each pair of data is shown in **Fig. 4A**. Except for COLO829 10X, both MuSE 2 and
170 Strelka2 continue to gain computational advantages with increasing number of CPU cores,
171 while the other four methods do not. We examine the overall speed performances of these
172 methods with MuSE 2 at 80 cores, Strelka2 at 80 scores, and the average time cost across
173 multiple runs over the different numbers of cores except for core=1 (where the computing
174 resource is too limited) for the other methods (**Fig. 4B**). Both MuSE 2 and Strelka2 achieve
175 much faster SNV calling compared to all other methods. For WES data, MuSE 2 accelerates 28-
176 58 times (mean=44) than MuSE 1, 68-83 times (mean=77) than MuTect2, 5-8 times (mean=8)
177 than SomaticSniper, 33-39 times (mean=36) than VarScan2. Similarly, for WGS data, it
178 accelerates 48-59 times (mean=57) than MuSE 1, 33-44 times (mean=41) than MuTect2, 7-8
179 times (mean=8) than SomaticSniper, 33-43 times (mean=37) than VarScan2 for WGS data. On
180 the other hand, Strelka2 is faster than MuSE 2 for all the WES data and the WGS data except
181 for COLO829 10X: it is about two-fold speedup on average compared to MuSE 2. For COLO829

182 10X, however, Strelka2 stopped gaining speed after 5 CPU cores, while MuSE 2 continued
183 accelerating; its computing time is 20 times of MuSE 2 at the number of cores of 80.

184 Overall, MuSE 2 and Strelka2 are the two top accelerated methods compared to others in the
185 above benchmarking. We further examine the difference between the SNV calls reported by
186 them for the same patient sample (**Fig. 4C**). For WES data, 46-78% (mean=66%) of the calls
187 are identified by both; 2-16% (mean=7%) of the calls are unique to MuSE 2, and 13-51%
188 (mean=27%) of the calls are unique to Strelka2. For WGS data, 41-77% (mean=62%) of the
189 calls are identified by both; 6-18% (mean=11%) of the calls are unique to MuSE 2, and 13-45%
190 (mean=26%) of the calls are unique to Strelka2.

191 We further investigate the feasibility of using the intersect calls from these two methods to
192 reproduce the consensus calls for these data generated by previous studies (Ellrott et al. 2018;
193 The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium 2020; Craig et al. 2016).
194 For WES data, these calls notably improve the F1 scores of Strelka2 calls, while maintaining a
195 comparable F1 score with the MuSE 2 calls (**Fig. 4D** and **Supplemental Table S2**). This
196 suggests that running MuSE 2 alone might be sufficient for WES data when the computing
197 resource is limited. For WGS data, on the contrary, Strelka2 calls (0.76-0.92) have higher F1
198 scores than MuSE 2 calls (0.63-0.88) in all the patient samples; while the intersection calls
199 outperform the two individual callers and reach the highest F1 scores (0.91-0.96). Also, the
200 intersect calls achieve the highest precision values (0.92-0.96 for WES, 0.87-0.95 for WGS) for
201 all the data benchmarked, despite the differences of read depths and sequencing platform they
202 were generated (**Supplemental Table S1**). The intersect calls maintain good recall values at
203 0.74-0.89 (median =0.86) for WES data, and at 0.96-0.98 (median=0.96) for WGS data. The
204 recall values of either the intersect call sets, or the individual call sets from the two methods are
205 consistently higher for WGS data than WES (**Fig. 4D**, all results from the WGS data fall in the
206 top rectangle). Also, MuSE 2 call sets achieve higher precisions and F1 scores, lower recalls for
207 WES, but higher recalls, lower precisions and F1 scores for WGS data, when compared to the

208 calls from Strelka2. Since the F1 scores of intersect calls are lower in WES than in WGS, we
209 further investigated whether the inclusion of a third caller could improve the accuracy for WES.
210 This led to an interesting observation that adding VarScan2 or SomaticSniper as a third caller did not
211 improve precision, recall or F1, while adding MuTect2 improved recall but decreased precision and
212 F1 (**Supplemental Table S3**). In the project of Multi-Center Mutation Calling in Multiple Cancers
213 (MC3) organized by TCGA, the latter strategy was further expanded to include five callers,
214 maximizing recall, and followed by a post-filtering pipeline to further improve precision,
215 independent of any callers (Ellrott et al. 2018). This post-filtering pipeline removes potentially
216 false positive variants that are caused by germline contamination, sequence artifacts, low read
217 depth in the normal sample, as well as non-exonic variants. In summary, combining mutation
218 calls from the two accelerated callers MuSE 2 and Strelka2, e.g., by simply taking an
219 intersection of the calls, is promising to achieve optimized mutation calling in a significantly
220 shorter wait-time. This strategy is particularly useful for WGS data and for analysis of large
221 patient cohorts. With WES data, running MuSE 2 alone can be a cost-effective strategy to obtain
222 mutation calls with high precision and a reasonable F1 score.

223 **A Snakemake pipeline for somatic SNV calling.** Finally, we introduce a fully automated
224 mutation calling pipeline for general users who do not have the time or expertise to learn about
225 the nuances in optimizing mutation calling accuracy, using the Snakemake workflow
226 management system (Köster and Rahmann 2012) (**Supplemental Fig. S5**). This user-friendly
227 pipeline allows for running all preprocessing, postprocessing, as well as MuSE 2 and Strelka2
228 together, in the background, without manual curation. It is compatible with typical Linux systems
229 and computing clusters, and optimizes memory and CPU utilization by parallelizing independent
230 tasks.

231

232 **Discussion**

233 Precision medicine and personalized cancer treatments have advanced remarkably in the last
234 decade, which greatly benefited from the accurate identification of genetic variations in the
235 tumor tissue using HTS data. An efficient and accurate somatic mutation caller is crucial to the
236 scientific studies of all cancers and their clinical management. Previously the accuracy and
237 utility of MuSE 1, either alone (Fan et al. 2016) or as a member of a multi-caller consensus
238 calling strategy has been validated by multiple consortial projects (Ellrott et al. 2018; The
239 ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium 2020). This study further
240 develops MuSE 2, in order to fully utilize resources on a high-performance computing machine,
241 including both the CPU cores and memory allocation. The producer-consumer model behind the
242 parallelization implemented in the step of ‘MuSE call’ gives MuSE 2 the ability to manage
243 multiple processes (workers) at the same time: they run independently at their own rates without
244 being affected by the computing load of other processes. Since the calculation in the step of
245 ‘MuSE sump’ is more straightforward – the computing speed bottlenecks only reside in several
246 for-loop iterations, we therefore use the OpenMP library, with which the parallelization is
247 relatively trivial. The speed-up by MuSE 2 becomes evident when it is run on at least 4-5 cores
248 to take advantage of the multi-step parallelization. In summary, MuSE 2 improves the mutation
249 calling utility of MuSE 1 by accelerating its computing speed by up to 50-60 times for both WES
250 and WGS data. MuSE 2 reduces the computational time cost of a somatic mutation calling
251 project from ~40 hours to < 1 hour for WGS data, and from 2-4 hours to ~5 minutes for WES
252 data, from each pair of tumor-normal samples.

253 MuSE 2 is much faster than the other three benchmarked callers, i.e., MuTect2,
254 SomaticSniper and VarScan2. It is slightly slower than Strelka2 for the sequencing data
255 generated by HiSeq 2000 and HiSeq X Ten, but much faster than the latter with NovaSeq.
256 Since we only include one pair of tumor-normal WGS data from HiSeq X Ten and NovaSeq,
257 respectively, more data are needed to validate this result in the future study. The intersection of
258 MuSE 2 and Strelka2 calls can substantially improve precisions without much loss in recalls,

259 hence improving the overall F1 scores for both the WES and WGS data benchmarked. We
260 therefore demonstrate the utility of the intersection calls from these two fast callers, as
261 compared to using each caller individually or using un-accelerated callers.

262 In contrast to the consensus calls of TCGA and PCAWG and the cell line study, running
263 MuSE 2 and Strelka2 to generate intersect calls may greatly improve the efficiency of genomic
264 data analysis for large patient cohorts, especially for those with WGS data. We also found
265 running MuSE 2 alone is a cost-effective solution for mutation calling in WES data, as it would
266 otherwise require 4-5 callers plus post-filters to achieve much higher recall and precision. Finally,
267 we provide a Snakemake workflow pipeline that automatically runs preprocessing, intersect
268 mutation calling using the two accelerated callers, and postprocessing without human
269 intervention, in order to improve accessibility by general users. We note that the hg19 genome
270 assembly was used throughout the study due to the fact that all consensus calls were based on
271 hg19. Given the underlying models of MuSE 2 and Strelka2, we expect the performance of
272 variant calling of both methods to be insensitive to genome assemblies. As the switch of
273 assembly build from hg19 to hg38 can impact preprocessing and read mapping to generate the
274 input data, some difference in variant calls could be observed, which should not substantially
275 affect the conclusions (Gao et al., 2019). Future development of MuSE 2 includes indel calling,
276 SNV calling from formalin-fixed paraffin-embedded (FFPE) samples and from tumor samples
277 only, all of which require an advanced error model construction as well as further benchmarking.
278 In summary, we expect the proposed MuSE 2 to significantly accelerate the variant calling
279 process and benefit the cancer research and clinical communities.

280

281 **Methods**

282 *Sample selection*

283 The consensus mutation calls of the TCGA portion of the PCAWG project were downloaded
284 from the ICGC data portal (https://dcc.icgc.org/releases/PCAWG/consensus_snv_indel). The
285 consensus mutation calls of the TCGA MC3 project were downloaded from GDC
286 (<https://gdc.cancer.gov/about-data/publications/mc3-2017>). We randomly selected 5 patient
287 samples from each of the two repositories, and downloaded the BAM files from
288 <https://dcc.icgc.org/repositories> and <https://portal.gdc.cancer.gov/repository>, respectively. The
289 sequencing data of these 10 patient samples were generated by Illumina HiSeq 2000. We
290 further downloaded BAM files of WGS data from the metastatic melanoma cell line COLO829
291 and the matched normal lymphoblastoid line COLO829BL from European Nucleotide Archive
292 (ENA) (accession code PRJEB27698, <https://www.ebi.ac.uk/ena/browser/view/PRJEB27698>).
293 We downloaded the latest WGS dataset (Espejo Valle-Inclan et al. 2022). The sequencing
294 libraries were either prepared with Illumina TruSeq Nano reagent kit and sequenced on the
295 HiSeq X Ten platform (COLO829 Illumina), or prepared on the 10X Chromium platform and
296 sequenced on the NovaSeq platform (COLO829 10X). For COLO829 Illumina, we also
297 downloaded the BAM files of the cell line with mixed tumor cell proportions of 75%, 50%, 25%,
298 20% and 10%. These data were simulated by *in silico* mixing of reads from COLO829 (100%
299 tumor) and COLO829BL (normal sample) with different ratios.

300 *BAM pre-processing*

301 MuSE 2 adopts the same preprocessing steps for the unaligned sequencing reads of the tumor-
302 normal pair as MuSE 1, which include trimming poor-quality bases, removing adapters, marking
303 duplicate reads, performing local indel realignment for the paired tumor-normal BAM files jointly,
304 and recalibrating base quality scores (**Supplemental Fig. S1**). In this study, the sequencing
305 reads are aligned against the hg19 reference genome build using BWA-MEM (Li 2013).

306 *Sequencing depth*

307 The sequencing depth of each BAM file after pre-processing is estimated by SAMtools with the
308 'depth' command. For WGS data, the overall depth was calculated as the average of the read

309 depths of all genomic locations. For WES data, the overall depth was calculated as the average
310 of the read depths of the genomic locations in the exon regions defined by the exome capture
311 kit downloaded from GDC (<https://gdc.cancer.gov/about-data/publications/mc3-2017>).

312 *Parallel computing implementation for MuSE*

313 **MuSE call.** We implement a multithreaded producer-consumer model which deploys threads for
314 parsing and uncompressing reads from BAM files, variant filtering and detection, writing outputs
315 and monitoring the whole process. The model connects all the threads concurrently by thread-
316 safe queues and atomic variables. We also adopt a faster and more efficient memory allocator
317 (i.e., TCMalloc: <https://github.com/google/tcmalloc>) rather than use the default malloc in C and
318 new in C++ in this step. The parallelization model starts with creating of 6 threads, 3 for the
319 BAM of tumor sample and the other 3 for the BAM of normal sample: 1 of the 3 threads parses
320 the compressed binary data and sends its reference to two queues, namely ChunkReadQueue
321 and ChunkUnzipQueue; the other two threads take the data from the ChunkUnzipQueue,
322 decompress it and label it as 'processed'. This change is also effective for the data in
323 ChunkReadQueue, since these two queues in fact store the same data. Another thread (i.e.,
324 read) is then created, which takes uncompressed data from ChunkReadQueue and recover it to
325 read format for both the BAM tumor sample and the BAM of normal sample, and pushes them
326 to the same queue, ReadQueue. A new thread named processReads is created; it parses the
327 reads from ReadQueue and sends them to the queue, processQ. n threads named workers are
328 created to take the reads from processQ and process them following the same pre-filtering and
329 the evolutionary model as MuSE 1. The last thread is named as 'monitor', which prints the sizes
330 of the queues every second. Here, users can specify n according to the number of cores
331 available in the input of MuSE 2 (**Supplemental Fig. S6**).

332 **MuSE sump.** We use the OpenMP library to parallelize the three most time-consuming parts in
333 MuSE sump. The first is the loading of candidate variants, the corresponding estimates of
334 equilibrium frequencies for all four alleles (A, C, T, G) for each variant from MuSE call, and

335 filtering out the variants whose ratio between the variant allele frequency from the normal
336 sample and the variant allele frequency from the tumor sample above a predefined cutoff (Fan
337 et al. 2016) (0.05). The second is scanning for the remaining variants in the dbSNP and marked
338 as 'true' or 'false' if they appear in the database or not. For WGS data, MuSE 1 fits a two-
339 component Gaussian mixture model to the allele frequencies of the post-filtered variants to
340 separate true mutations from background noise. The parameters (e.g., mean, standard
341 deviation and proportion) of the two components are estimated using the expectation-
342 maximization algorithm which are repeated 50 times with random initializations. For the three
343 parts, we parallel the for-loop iterations using the 'omp parallel for' clause from OpenMP in
344 MuSE 2 to deploy the computation on multiple cores.

345 *Speed benchmarking settings*

346 For all the benchmarked methods, if the number of cores requested lies in {1, 5, 10, 20, 28}, the
347 processor is Intel(R) Xeon(R) Gold 6132 CPU @ 2.60GHz; if the number of cores requested lies
348 in {40, 80}, the processor is Intel(R) Xeon(R) Platinum 8380 CPU @ 2.30GHz. We run each
349 method by submitting the LSF (Load Sharing Facility) job script using the bsub command, with
350 which we can easily control the random-access memory (RAM) and the number of cores
351 specified for each method. The options for the 6 callers can be found in **Supplemental Table**
352 **S4**.

353 *Precision and recall*

354 For the samples from TCGA and PCAWG, we used the consensus SNV calls published
355 previously (Ellrott et al. 2018; The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes
356 Consortium 2020) as a truth set; for the cell line COLO829/COLO829BL, we used the SNV calls
357 downloaded from (Craig et al. 2016) as a truth set. The version information of the benchmarked
358 callers is listed in **Supplemental Table S5**. Indel calls were removed from the call set before
359 any comparison. For WGS data, we took calls from all tiers in MuSE 2 (Fan et al. 2016), and

360 only calls from the PASS category from the other callers for each patient sample. We filtered out
 361 low quality SNVs from the consensus calls from PCAWG WGS data that are labeled as
 362 'LOWSUPPORT', 'OXOGFAIL', 'bSeq', 'bPcr', 'GERM1000G', 'GERMOVLP', 'NORMALPANEL'
 363 or 'REMAPFAIL' (https://dcc.icgc.org/releases/PCAWG/consensus_snv_indel). Consistently, we
 364 used consensus calls for the cell line COLO829/COLO829BL which had already gone through a
 365 similar post-filtering process (Craig et al. 2016). For WES data, we selected calls from all the
 366 categories except for 'Tier5' from MuSE 2, and only calls from the PASS category from the other
 367 callers for each patient sample. We also used the consensus calls of TCGA WES that had
 368 already gone through post-filtering (Ellrott et al. 2018). The intersection between any two sets
 369 from the same patient sample were identified by matching the SNV ids, which combined the
 370 columns of CHROM, POS, REF and ALT from the two VCF files. For WES data, we removed
 371 the SNVs from the intersection calls outside the regions defined by the exome capture kit of
 372 TCGA.

373 We considered any calls reported by the consensus, but not by the intersection calls as false
 374 negatives, any calls reported by the intersection calls, but not by the consensus as false
 375 positives. We calculated precision, recall and F1 score to evaluate the accuracy of a call set
 376 against a truth set.

$$F_1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

377 *Definitions of variant allele frequency bins, sequencing depth bins, clonality and variant effect*
 378 *annotation for SNVs*

379 To resolve the issue that the read depth (including the number of reads supporting the reference
 380 allele and the alternate allele) can be different for the same SNVs from MuSE 2, Strelka2 and
 381 the consensus calls, we used alleleCount (<https://github.com/cancerit/alleleCount>) to recalculate
 382 the read depth and VAF for all the unique SNVs from MuSE 2, Strelka2 and the consensus calls.
 383 We finally generated the bins of VAFs (i.e., 0-0.2, 0.2-0.3, 0.3-0.4, 0.4-0.5 and >0.5) and read

384 depths (<40x, 40-80x, 80-120x, 120-160x and >160x for WES data, <20x, 20-40x, 40-60x, 60-
385 80x and >80x for WGS data) for the calls of each method for detailed comparisons.

386 We downloaded the consensus subclonal reconstruction results (Dentro et al. 2021) for the
387 consensus calls of PCAWG WGS data from ICGC data portal:

388 https://dcc.icgc.org/releases/PCAWG/subclonal_reconstruction/. Each SNV in the consensus
389 calls is defined as either clonal or subclonal when such information is available. To compare the
390 performance between MuSE 2 and other callers, we restricted the SNVs from each caller
391 overlapping with the ones from the consensus calls such that they can be annotated as clone or
392 subclone. Therefore, only recall is evaluated.

393 We used Ensembl VEP (McLaren et al. 2016) (v101) to predict the effect of a SNV. For
394 simplicity, we merged nonsense and missense variants into nonsynonymous variants; variants
395 in 3'UTR (untranslated region), 5' UTR, 3' Flank and 5' Flank into untranslated region variants;
396 variants in splice region, translation start site and RNA variants into others category. We also
397 renamed silent variants to synonymous variants. We have four classes for the SNVs from TCGA
398 WES data: nonsynonymous, synonymous, untranslated region and others, and six classes for
399 the SNVs from PCAWG or cell line COLO829/COLO829BL WGS data: nonsynonymous,
400 synonymous, intergenic region, intron, untranslated region and others.

401 *Software availability*

402 MuSE 2 is implemented in C++ and is available at GitHub <https://github.com/wwwylab/MuSE> with
403 the GPL-2.0 license. A Dockerfile is included in the repository for building MuSE 2 into a Docker
404 container running on Linux machines. A Snakemake pipeline for somatic SNV calling,
405 MuSE.Snakemake 1.0, is also available on the GitHub repository. The source code of the
406 repository can also be found in the **Supplemental Material**.

407

408 **Competing interest statement**

409 Tong Zhu and Ankit Sethia are employees of NVIDIA Corp. and own NVIDIA stock as part of the
410 standard compensation package.

411

412 **Acknowledgements**

413 S.J. is supported by the MD Anderson Colorectal Cancer Moon Shot Program and National
414 Institutes of Health (NIH) R01CA268380. W.W. is supported by DoD PC210079, NIH
415 R01CA268380, P30CA016672. S.J. and W.W are also supported by The MD Anderson Cancer
416 Center SPORE in Gastrointestinal Cancer Grant P50 CA221707. The authors thank the
417 Biomedical Visualization team at Anderson Cancer Center for their assistance in creating Figure
418 1. Copyright used with the permission of The Board of Regents of the University of Texas
419 System through The University of Texas MD Anderson Cancer Center.

420

421 **Author contributions**

422 S.J. implemented the parallelization for the 'MuSE sump' step, performed the speed
423 benchmarking, analyzed the results, and wrote the manuscript in collaboration with other
424 authors. T.Z. and A.S. implemented the parallelization for the 'MuSE call' step. W.W. conceived
425 the project, planned, and supervised the work, wrote the manuscript, in collaboration with all
426 other authors. All authors commented on and approved the final manuscript.

427 **References**

- 428 Alfares A, Aloraini T, subaie L Al, Alissa A, Qudsi A Al, Alahmad A, Mutairi F Al, Alswaid A,
429 Alothaim A, Eyaid W, et al. 2018. Whole-genome sequencing offers additional but limited
430 clinical utility compared with reanalysis of whole-exome sequencing. *Genetics in Medicine*
431 **20**: 1328–1333.
- 432 Barbitoff YA, Polev DE, Glotov AS, Serebryakova EA, Shcherbakova I V, Kiselev AM, Kostareva
433 AA, Glotov OS, Predeus A V. 2020. Systematic dissection of biases in whole-exome and
434 whole-genome sequencing reveals major determinants of coding sequence coverage. *Sci*
435 *Rep* **10**: 1–13.
- 436 Christoforides A, Carpten JD, Weiss GJ, Demeure MJ, Von Hoff DD, Craig DW. 2013.
437 Identification of somatic mutations in cancer through Bayesian-based analysis of
438 sequenced genome pairs. *BMC Genomics* **4**: 302.
- 439 Cibulskis K, Lawrence MS, Carter SL, Sivachenko A, Jaffe D, Sougnez C, Gabriel S, Meyerson
440 M, Lander ES, Getz G. 2013. Sensitive detection of somatic point mutations in impure and
441 heterogeneous cancer samples. *Nat Biotechnol* **31**: 213–219.
- 442 Craig DW, Nasser S, Corbett R, Chan SK, Murray L, Legendre C, Tembe W, Adkins J, Kim N,
443 Wong S, et al. 2016. A somatic reference standard for cancer genome sequencing. *Sci*
444 *Rep* **6**: 24607.
- 445 Dentre SC, Leshchiner I, Haase K, Tarabichi M, Wintersinger J, Deshwar AG, Yu K, Rubanova
446 Y, Macintyre G, Demeulemeester J, et al. 2021. Characterizing genetic intra-tumor
447 heterogeneity across 2,658 human cancer genomes. *Cell* **184**: 2239-2254.e39.
- 448 Ellrott K, Bailey MH, Saksena G, Covington KR, Kandoth C, Stewart C, Hess J, Ma S, Chiotti
449 KE, McLellan MD, et al. 2018. Scalable Open Science Approach for Mutation Calling of
450 Tumor Exomes Using Multiple Genomic Pipelines. *Cell Syst* **6**: 271-281.e7.
- 451 Espejo Valle-Inclan J, Besselink NJM, de Bruijn E, Cameron DL, Ebler J, Kutzera J, van
452 Lieshout S, Marschall T, Nelen M, Priestley P, et al. 2022. A multi-platform reference for
453 somatic structural variation detection. *Cell Genomics* **2**: 100139.
- 454 Fan Y, Xi L, Hughes DST, Zhang J, Zhang J, Futreal PA, Wheeler DA, Wang W. 2016. MuSE:
455 accounting for tumor heterogeneity using a sample-specific error model improves
456 sensitivity and specificity in mutation calling from sequencing data. *Genome Biol* **17**: 178.
- 457 Gao GF, Parker JS, Reynolds SM, Silva TC, Wang LB, Zhou W, Akbani R, Bailey M, Balu S,
458 Berman BP, et al. 2019. Before and After: Comparison of Legacy and Harmonized TCGA
459 Genomic Data Commons' Data. *Cell Syst* **9**: 24-34.e10.
- 460 Kim S, Scheffler K, Halpern AL, Bekritsky MA, Noh E, Källberg M, Chen X, Kim Y, Beyter D,
461 Krusche P, et al. 2018. Strelka2: fast and accurate calling of germline and somatic variants.
462 *Nat Methods* **15**: 591–594.
- 463 Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, Miller CA, Mardis ER, Ding L,
464 Wilson RK. 2012. VarScan 2: Somatic mutation and copy number alteration discovery in
465 cancer by exome sequencing. *Genome Res* **22**: 568–576.
- 466 Köster J, Rahmann S. 2012. Snakemake—a scalable bioinformatics workflow engine.
467 *Bioinformatics* **28**: 2520–2522.
- 468 Larson DE, Harris CC, Chen K, Koboldt DC, Abbott TE, Dooling DJ, Ley TJ, Mardis ER, Wilson
469 RK, Ding L. 2012a. SomaticSniper: identification of somatic point mutations in whole
470 genome sequencing data. *Bioinformatics* **28**: 311–317.
- 471 Larson DE, Harris CC, Chen K, Koboldt DC, Abbott TE, Dooling DJ, Ley TJ, Mardis ER, Wilson
472 RK, Ding L. 2012b. SomaticSniper: identification of somatic point mutations in whole
473 genome sequencing data. *Bioinformatics* **28**: 311–317.
- 474 Li, H. 2011. Tabix: Fast retrieval of sequence features from generic TAB-delimited files.
475 *Bioinformatics* **27**: 718–719.

- 476 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
477 *arXiv:13033997*.
- 478 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R.
479 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079.
- 480 McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, Flicek P, Cunningham F. 2016.
481 The Ensembl Variant Effect Predictor. *Genome Biol* **17**: 1–14.
- 482 Saunders CT, Wong WSW, Swamy S, Becq J, Murray LJ, Cheetham RK. 2012. Strelka:
483 Accurate somatic small-variant calling from sequenced tumor-normal sample pairs.
484 *Bioinformatics* **28**: 1811–1817.
- 485 Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K. 2001. dbSNP:
486 The NCBI database of genetic variation. *Nucleic Acids Res* **29**: 308–311.
- 487 Sun Y, Liu F, Fan C, Wang Y, Song L, Fang Z, Han R, Wang Z, Wang X, Yang Z, et al. 2021.
488 Characterizing sensitivity and coverage of clinical WGS as a diagnostic test for genetic
489 disorders. *BMC Med Genomics* **14**: 1–13.
- 490 The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. 2020. Pan-cancer
491 analysis of whole genomes. *Nature* **578**: 82–93.
- 492 Xu C. 2018. A review of somatic single nucleotide variant calling algorithms for next-generation
493 sequencing data. *Comput Struct Biotechnol J* **16**: 15–24.
- 494

495

496 **Figure legends**

497 **Figure 1. Assembly line illustration of the multi-step parallelization implemented in MuSE 2.**

498 (A) 'MuSE call'. Workers (threads) keep fetching chunks from the input BAM files from the tumor and
499 normal samples and unzipping them to the text format of reads. Downstream workers combine the
500 reads from tumor and normal samples and send to a queue; from there, other workers detect
501 candidate variants. (B) 'MuSE sump'. Multiple workers are used to take the candidate variants and
502 their corresponding estimated summary statistic π 's and scan them against dbSNP database,
503 labeling those appearing in the database. For candidate variants from the WGS data, we fit two-
504 component Gaussian Mixture Models (GMMs) with multiple initializations, distributed to multiple
505 workers, in order to separate true variants from background noise; for candidate variants from the
506 WES data, no parallelization is implemented due to computational simplicity as we simply fit a beta
507 distribution to π 's.

508

509 **Figure 2. Comparisons of F1 score, precision and recall between MuSE 2 and Strelka2**
510 **within each bin of variant allele frequency (VAF, top) or sequencing read depth (middle),**

511 **variant class (bottom) for TCGA WES data.** The calls of each method, and the consensus
512 calls, which are used as the truth set, are pooled from the WES data of 5 patient samples in
513 TCGA.

514

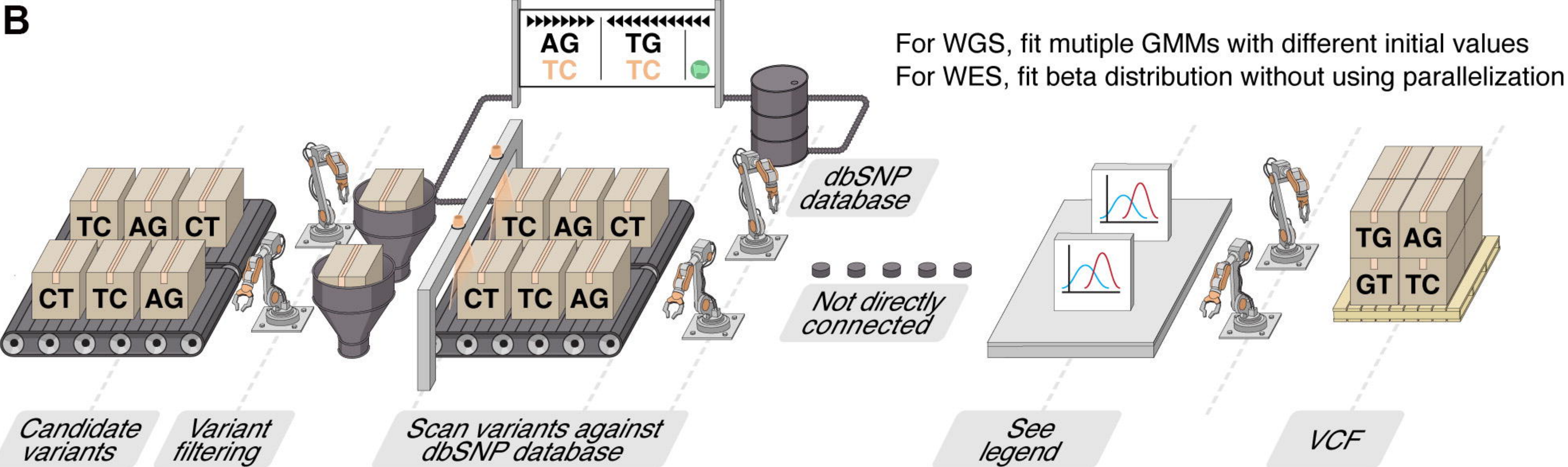
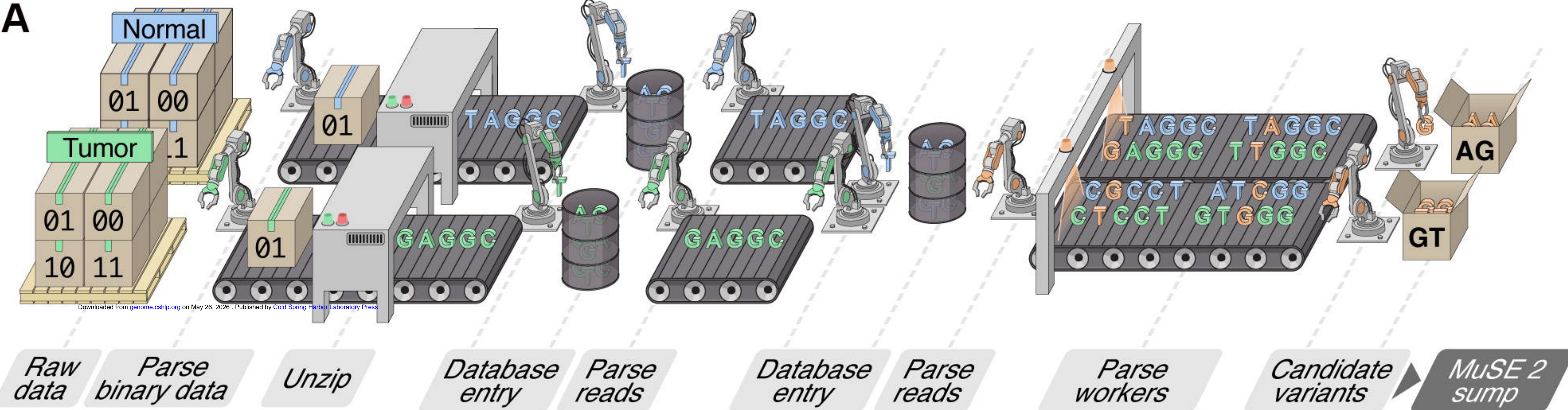
515 **Figure 3. Accuracy benchmarking of MuSE 2 and Strelka2** within each bin of VAF (first row) or
516 sequencing read depth (second row), class of variant effects (third row) for PCAWG WGS (A) and
517 cell line WGS data (B). Comparison of recall between the two methods within different clonality for
518 PCAWG WGS data is shown in the last row of (A). The calls of each method, and the truth set, are
519 pooled from the WGS data of the selected 5 patient samples from PCAWG (A), or the WGS data of
520 the cell line COLO829/COLO829BL generated from two platforms (B). The number of consensus

521 calls for a variant class is included in the x-axis labels. (C) F1 scores of MuSE 2 and Strelka2 with
522 varied tumor cell proportions. Calls from the tumor purity of 100% were used as the truth set.

523

524 **Figure 4. Benchmarking the speed and usability of MuSE 2.** (A) The runtime of MuSE 2 against
525 MuSE 1 and the other four methods for both WES and WGS data across different numbers of cores.
526 The numbers in the plot represent the fold speedup of MuSE 2 (with 80 cores) relative to the other
527 methods whose time cost is averaged across different numbers of cores (excluding No. of core=1).
528 For Strelka2, only the time cost with 80 cores is considered. (B) A simplified version of (A) in which
529 the time cost of each method is averaged across all samples (excluding COLO829 10X as an outlier,
530 see (A)) and different numbers of cores. (C) Venn diagrams showing the unique and shared SNV
531 calls of MuSE 2 and Strelka2. (D) Scatter plot of the precision, recall for the intersection calls from
532 MuSE 2 (in red) and Strelka2 (in blue), the calls from each of the two methods (in purple) against the
533 previously reported consensus calls which are considered as the benchmark. For both WES (circle)
534 and WGS (triangle) data, the median F1 scores of the intersection calls, calls from each individual
535 method are shown. Two shaded rectangles highlight the difference of the performance metrics
536 between WES and WGS data. Results from the WGS data are located in the top rectangle.

537



TCGA WES

