

1 **Large-Scale Uniform Analysis of Cancer Whole Genomes in Multiple** 2 **Computing Environments**

3
4 Christina K. Yung^{1,*}, Brian D. O'Connor^{1,2,*}, Sergei Yakneen^{1,3,*}, Junjun Zhang^{1,*}, Kyle Ellrott⁴,
5 Kortine Kleinheinz^{5,6}, Naoki Miyoshi⁷, Keiran M. Raine⁸, Romina Royo⁹, Gordon B. Saksena¹⁰,
6 Matthias Schlesner⁵, Solomon I. Shorser¹, Miguel Vazquez¹¹, Joachim Weischenfeldt^{3,12}, Denis
7 Yuen¹, Adam P. Butler⁸, Brandi N. Davis-Dusenbery¹³, Roland Eils^{14,6}, Vincent Ferretti¹, Robert L.
8 Grossman¹⁵, Olivier Harismendy^{16,17}, Youngwook Kim¹⁸, Hidewaki Nakagawa¹⁹, Steven J.
9 Newhouse²⁰, David Torrents^{9,21}, Lincoln D. Stein^{1,22,‡} on behalf of the PCAWG Technical Working
10 Group²³ and the PCAWG Network

11
12 * *These authors contributed equally to this work.*

13 ‡ *Corresponding author: lincoln.stein@gmail.com*

14
15 ¹Informatics and Biocomputing Program, Ontario Institute for Cancer Research, Toronto, Ontario, M5G 0A3, Canada.
16 ²UC Santa Cruz Genomics Institute, University of California Santa Cruz, Santa Cruz, California, 95065, USA.
17 ³Genome Biology Unit, European Molecular Biology Laboratory, Heidelberg, Baden-Württemberg, 69120, Germany.
18 ⁴Department of Computational Biology, Oregon Health and Science University, Portland, Oregon, 97239, USA.
19 ⁵Division of Theoretical Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg, Baden-Württemberg,
20 69120, Germany. ⁶Department for Bioinformatics and Functional Genomics, Institute for Pharmacy and Molecular
21 Biotechnology and BioQuant, Heidelberg University, Heidelberg, Baden-Württemberg, 69120, Germany. ⁷Human
22 Genome Center, Institute of Medical Science, University of Tokyo, Tokyo, 108-8639, Japan. ⁸Cancer Ageing and
23 Somatic Mutation Programme, Wellcome Trust Sanger Institute, Hinxton, Cambridgeshire, CB10 1SA, United
24 Kingdom. ⁹Department of Life Sciences, Barcelona Supercomputing Center, Barcelona, Catalunya, 8034, Spain.
25 ¹⁰Cancer Program, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, 02142, USA. ¹¹Structural
26 Computational Biology Group, Centro Nacional de Investigaciones Oncológicas, Madrid, Madrid, 28029, Spain.
27 ¹²BRIC/Finsen Laboratory, Rigshospitalet, Copenhagen, 2200, Denmark. ¹³Seven Bridges, Cambridge,
28 Massachusetts, 02142, USA. ¹⁴Theoretical Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg,
29 Baden-Württemberg, 69120, Germany. ¹⁵Center for Data Intensive Science, University of Chicago, Chicago, Illinois,
30 60637, USA. ¹⁶Department of Medicine, University of California San Diego, San Diego, California, 92093, USA.
31 ¹⁷Moore's Cancer Center, Department of Medicine, Division of Biomedical Informatics, University of California San
32 Diego, San Diego, California, 92093, USA. ¹⁸Samsung Advanced Institute of Health Science and Technology,
33 Sungkyunkwan University, School of Medicine, Seoul, 135-710, South Korea. ¹⁹Laboratory for Genome Sequencing
34 Analysis, RIKEN Center for Integrative Medical Sciences, Tokyo, 108-8639, Japan. ²⁰Technical Services Cluster,
35 European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton, Cambridge, CB10 1SD, United
36 Kingdom. ²¹Institució Catalana de Recerca i Estudis Avançats, Barcelona, Catalunya, 8010, Spain. ²²Department of
37 Molecular Genetics, University of Toronto, Toronto, Ontario, M5S 1A1, Canada. ²³Full lists of members and
38 affiliations appear at the end of the paper.

39

40 **Abstract**

41 The International Cancer Genome Consortium (ICGC)'s Pan-Cancer Analysis of Whole Genomes
42 (PCAWG) project aimed to categorize somatic and germline variations in both coding and non-
43 coding regions in over 2,800 cancer patients. To provide this dataset to the research working
44 groups for downstream analysis, the PCAWG Technical Working Group marshalled ~800TB of
45 sequencing data from distributed geographical locations; developed portable software for uniform
46 alignment, variant calling, artifact filtering and variant merging; performed the analysis in a
47 geographically and technologically disparate collection of compute environments; and
48 disseminated high-quality validated consensus variants to the working groups. The PCAWG
49 dataset has been mirrored to multiple repositories and can be located using the ICGC Data Portal.
50 The PCAWG workflows are also available as Docker images through Dockstore enabling
51 researchers to replicate our analysis on their own data.

52 **Introduction**

53 The International Cancer Genome Consortium (ICGC)/The Cancer Genome Atlas (TCGA) Pan-
54 Cancer Analysis of Whole Genomes (PCAWG) study has characterized the pattern of mutations
55 in over 2,800 cancer whole genomes. Extending TCGA Pan-Cancer analysis project, which
56 focused on molecular aberrations in protein coding regions only¹, PCAWG undertook the study of
57 whole genomes, allowing for the discovery of driver mutations in cis-regulatory sites and non-
58 coding RNAs, examination of the patterns of large-scale structural rearrangements, identification
59 of signatures of exposure, and elucidation of interactions between somatic mutations and germline
60 polymorphisms.

61 The PCAWG dataset comprises a total of 5,789 whole genomes of tumors and matched normal
62 tissue spanning 39 tumor types. The tumor/normal pairs came from a total of 2,834 donors

63 collected and sequenced by 48 sequencing projects across 14 jurisdictions (Supplementary Fig. 1).
64 In addition, RNA-Seq profiles were obtained from a subset of 1,284 of the donors². While the
65 individual sequencing projects contributing to PCAWG had previously identified genomic variants
66 within their individual cancer cohorts, each project had used their own preferred methods for read
67 alignment, variant calling and artifact filtering. During initial evaluation of the data set, we found
68 that the different analysis pipelines contributed high levels of technical variation, hindering
69 comparisons across multiple cancer types³. To eliminate the variations arising from non-uniform
70 analysis, we reanalyzed all samples starting with the raw sequencing reads and using a
71 standardized set of alignment, variant calling and filtering methods. These “core” workflows
72 yielded uniformly analyzed genomic variants for downstream analyses by various PCAWG
73 working groups. A subset of these variants were validated through targeted deep sequencing to
74 estimate the accuracy of our approach⁴.

75 To create this uniform analysis set, multiple logistic and technical challenges had to be overcome.
76 First, projects participating in the PCAWG study employed their own metadata conventions for
77 describing their raw sequencing data sets. Hence, we had to establish a PCAWG metadata standard
78 suitable for all the participating projects. Second, and more significantly, the data was large in size
79 -- 800TB of raw sequencing reads -- and distributed geographically across the world. During
80 realignment, the data transiently doubled in size, and after final variant calling and other
81 downstream analysis, the full data set reached nearly 1PB. Furthermore, the compute necessary to
82 fully harmonize the data was estimated at more than 30 million core-hours. Both the storage and
83 compute requirements made it impractical to complete the analysis at any single research institute.
84 In addition, legal constraints across the various jurisdictions imposed restrictions as to where
85 personal data could be stored, analyzed and redistributed⁵. Hence, we needed a protocol to spread

86 the compute and storage resources across multiple commercial and academic compute centers.
87 This requirement, in turn, necessitated the development of analysis pipelines that would be
88 portable to different compute environments and yield consistent analysis results independent of
89 platform. With multiple analysis pipelines running simultaneously in multiple compute
90 environments, the assignment of workload, tracking of progress, quality checking of data and
91 dissemination of results all required sophisticated and flexible planning.

92 Our approach to tackling these challenges was unique and substantially different from previous
93 large-scale genome analysis endeavors. First, as a collaborative effort among a wide range of
94 institutions not backed by a centralized funding source, a high degree of coordination among a
95 large task force of volunteer software engineers, bioinformaticians and computer scientists was
96 required. Second, the project fully embraced the use of both public and private cloud compute
97 technologies while leveraging established high-performance computing (HPC) infrastructures to
98 fully utilize the compute resources contributed by the partner organizations. The cloud technology
99 platforms we utilized included both Infrastructure as a Service (IaaS): OpenStack, Amazon Web
100 Services and Microsoft Azure; and Platform as a Service (PaaS): Seven Bridges (SB). Lastly, the
101 project made heavy use of Docker, a new lightweight virtualization technology that ensured
102 workflows, tools and infrastructure would work identically across the large number of compute
103 environments utilized by the project.

104 Utilizing the compute capacity contributed by academic HPC, academic clouds and commercial
105 clouds (Table 1), we were able to complete a uniform analysis of the entire set of 5,789 whole
106 genomes in just over 23 months (Figure 1). Figure 3 illustrates the three broad phases of the project:
107 (1) Marshalling and upload of the data into data analysis centres (3 months); (2) Alignment and
108 variant calling (18 months); and (3) Quality filtering, merging, synchronization and distribution of

109 the variant calls to downstream research groups (2 months). A fourth phase of the project, in which
110 PCAWG working groups used the uniform variant calls for downstream analysis, such as cancer
111 driver discovery, began in the summer of 2016 and continued through the first two quarters of
112 2017.

113 The following sections will describe the technical solutions used to accomplish each of the phases
114 of the project.

115 **Phase 1: Data Marshalling and Upload**

116 A significant challenge for the project was that at its inception, a large portion of the raw read
117 sequencing data had yet to be submitted to a read archive and thus had no standard retrieval
118 mechanism. In addition, the metadata standards for describing the raw data varied considerably
119 from project to project. For this reason, we asked the participating projects to prepare and upload
120 the 774 TB of raw whole genome sequencing (WGS) data and 27 TB raw RNA-seq data into a
121 series of geographically distributed data repositories, each running a uniform system for registering
122 the data set, accepting and validating the raw read data and standardized metadata.

123 We utilized seven geographically distributed data repositories located at: (1) Barcelona
124 Supercomputing Centre (BSC), (2) European Bioinformatics Institute (EMBL-EBI) in the UK, (3)
125 German Cancer Research Center (DKFZ) in Germany; (4) the University of Tokyo in Japan; (5)
126 Electronics and Telecommunications Research Institute (ETRI) in South Korea; (6) the Cancer
127 Genome Hub (CGHub) and (7) the Bionimbus Protected Data Cloud (PDC) in the USA (Figure 2
128 and Suppl Table 1).

129 To accept and validate sequence set uploads, each data repository ran a commercial software
130 system, GNOS (Annai Systems). We chose GNOS because of the heavy testing it had previously

131 received as the engine powering TCGA CGHub, and its support for validation of metadata
132 according to the Sequence Read Archive (SRA) standard and file submission, strong user
133 authentication and encryption, as well as its highly optimized data transfer protocol⁶. Each of the
134 seven data centers initially allocated several hundred terabytes of storage to accept raw sequencing
135 data from submitters within the region. The data centers also provided co-located compute
136 resources to perform alignment and variant calling on the uploaded data.

137 Genomic data uploaded to the GNOS repositories was accompanied with detailed and accurate
138 metadata to describe the cancer type, sample type, sequencing type and other attributes for
139 managing and searching the files. We required that identifiers for project, donor, sample follow a
140 standardized convention such that validation and auditing tools could be implemented. Most of the
141 naming conventions in PCAWG were adopted from the well established ICGC data dictionary
142 (<http://docs.icgc.org/dictionary/about/>).

143 Since most member projects at the time of upload already had sequencing reads aligned and
144 annotated using their own metadata standards, a non-trivial effort was required to prepare the
145 sequencing data for submission to GNOS. Each member project had to (1) prepare lane-level
146 unaligned reads in BAM format, (2) reheader the BAM files with metadata following the PCAWG
147 conventions, (3) generate metadata XML files, and (4) upload the BAM files along with the
148 metadata XML files to GNOS. To facilitate this process, we developed the *PCAP-core* tool
149 (<https://github.com/ICGC-TCGA-PanCancer/PCAP-core>) to extract the metadata from the BAM
150 headers, validate the metadata, transform the metadata into the XML files conforming to the SRA
151 specifications, and submitting the BAM files along with the metadata XML files to GNOS.

152

153 **Phase 2: Sequence Alignment and Variant Calling**

154 We began the process of sequence alignment about two months after the uploading process had
155 begun. Both tumor and matched normal reads were subjected to uniform sequence alignment using
156 BWA-MEM⁷ on top of a common GRCh37-based reference genome that was enhanced with decoy
157 sequences, viral sequences, and the revised Cambridge reference genome for the mitochondria.

158 Efforts by the project QC group demonstrated that employing multiple variant callers in ensemble
159 fashion improved calling sensitivity³, thus the aligned tumor/normal pairs were subjected to
160 somatic variant calling using three “best practice” software pipelines. These pipelines were
161 developed by the Sanger Institute⁸⁻¹¹; jointly by DKFZ¹² and the European Molecular Biology
162 Laboratory (EMBL)¹³; and the Broad Institute¹⁴ with contribution from MD Anderson Cancer
163 Center-Baylor College of Medicine¹⁵. Each pipeline represents the best practices from the
164 authoring organizations and include the current versions of each institute’s flagship tools. Each
165 pipeline consists of multiple software tools for calling of single and multiple nucleotide variants
166 (SNVs and MNVs), small insertions/deletions (indels), structural variants (SVs) and somatic copy
167 number alterations (SCNAs). The minimum compute requirements, median runtime and the
168 analytical algorithms for each pipeline are shown in Table 2.

169 When possible, both the alignment and variant calling pipelines were executed in the same regional
170 compute centers to which the data sets were uploaded. As the project progressed, we utilized
171 additional compute resources from AWS, Azure, iDASH, the Ontario Institute for Cancer
172 Research (OICR), the Sanger Institute, and Seven Bridges (Figure 2). These centers computed on
173 data sets located in the same region to optimize data transfer. Over the course of the project, some
174 centers outpaced others and we rebalanced data sets as needed to use resources as efficiently as

175 possible. Figure 1 shows the progress of the analytic pipelines with more details shown in
176 Supplementary Figures 2-6.

177 **Phase 3: Variant merging, filtering, and synchronization**

178 Following the completion of the three variant calling workflows, variants were passed to an
179 additional pipeline referred as the “OxoG workflow”. This pipeline filtered out oxidative artifacts
180 in SNVs using the OxoG algorithm¹⁶, normalized indels using the bcftools “norm” function,
181 annotated genomic features for downstream merging of variants, and generated one “minibam”
182 per specimen using the VariantBam algorithm¹⁷. Minibams are a novel format for representing the
183 evidence that underlies genomic variant calls. Read pairs spanning a variant within a specified
184 window were extracted from the whole genome BAM to generate the minibam. The windows we
185 chose were +/- 10 base pairs (bp) for SNVs, +/- 200 bp for indels, and +/- 500 bp for SV
186 breakpoints. The resulting minibams are about 0.5% of the size of whole genome BAMs, totalling
187 to about four terabytes for all PCAWG specimens, making it much easier to download and store
188 for the purpose of inspecting variants and their underlying read evidence.

189 Following filtering, we applied a series of merge algorithms to merge variants from the multiple
190 variant calling pipelines into consensus call sets with higher accuracies than the individual
191 pipelines alone. The SNV and indel merge algorithms were developed on the basis of experimental
192 validation of the individual variant calling pipelines using deep targeted sequencing, a process
193 detailed in the PCAWG-1 marker paper⁴. The algorithm for consensus SVs is described in the
194 PCAWG-6 marker paper¹⁸. The consensus SCNAs were built upon the base-pair breakpoint
195 results from the consensus SVs using a multi-tiered bespoke approach combining results from 6
196 SCNA algorithms¹⁹.

197 Following merging, the SNV, indel, SV and SCNA consensus call sets were subjected to intensive
198 examination by multiple groups in order to identify anomalies and artefacts, including uneven
199 coverage of the genome, strand and orientation bias, contamination with reads from non-human
200 species, contamination of the library with DNA from an unrelated donor, and high rates of common
201 germline polymorphisms among the somatic variant calls^{4,11}. In keeping with our mission to
202 provide a high-quality and uniformly annotated data set, we developed a series of filters to annotate
203 and/or remove these artefacts. Tumor variant call sets that were deemed too problematic to use for
204 downstream analysis were placed on an “exclusion list” (353 specimens, 176 donors). In addition,
205 we established a “grey list” (150 specimens, 75 donors), of call sets that had failed some tests but
206 not others and could be used, with caution, for certain types of downstream analysis. The criteria
207 for classifying callsets into exclusion and grey list are described in more detail in the PCAWG-1
208 paper¹⁰.

209 Following the filtering steps, we used GNOS to synchronize the aligned reads and variant call sets
210 among a small number of download sites for use by PCAWG downstream analysis working groups
211 (Suppl Table 2). We also provided login credentials to members of PCAWG working groups for
212 compute cloud-based access to the aligned read data across several of the regional data analysis
213 centers, which avoided the overhead of downloading the data.

214 **Software and Protocols**

215 This section describes the software and protocols developed for this project in more detail. All the
216 software that we created for this project is available for use by any research group to conduct
217 similar cloud-based cancer genome analyses economically and at scale.

218

219 Centralized Metadata Management System

220 The metadata describing the donors, specimens, raw sequencing reads, WGS and RNA-Seq
221 alignments, variant calls from the three pipelines, OxoG-filtered variants, and mini-BAMs were
222 collected from globally distributed GNOS repositories, consolidated and indexed nightly using
223 ElasticSearch (<https://www.elastic.co>) in a specially designed object graph model. This centrally
224 managed metadata index was a key component of our operations and data provenance tracking.
225 First, the metadata index was critical for tracking the status of each sequencing read set and for
226 scheduling the next analytic step. The index also tracked the current location of each BAM and
227 variant call set, allowing the pipelines to access the needed input data efficiently. Second, the
228 metadata index provided the basis for a dashboard (<http://pancancer.info>) for all stakeholders to
229 track day-to-day progress of each pipeline at each compute site. By reviewing the throughput of
230 each compute site on a daily basis, we were able to identify issues early and to assign work
231 accordingly to keep our compute resources productive. Third, the metadata index was also used
232 by the ICGC Data Coordination Centre (DCC) to transfer PCAWG core datasets to long-term
233 genomic data archive systems. Finally, the metadata index was imported into the ICGC Data Portal
234 (<https://dcc.icgc.org>) to create a faceted search for PCAWG data allowing users to quickly locate
235 data based on queries about the donor, cancer type, data type or data repositories.

236 Docker Containers & Consonance

237 Given that the compute resources donated to the PCAWG project were a mix of cloud and HPC
238 environments, we required a mechanism to encapsulate the analytical workflows to allow them to
239 run smoothly across a wide variety of compute sites. The approaches we used evolved over time
240 to incorporate better ways of abstracting and packaging tools to facilitate this portability. Initially,
241 we used SeqWare workflow execution engine²⁰ for bundling software and executing workflows,

242 but this system required extensive and time consuming setup for the worker virtual machines
243 (VMs). Later, we adopted Docker (<http://www.docker.com>) as a key enabling technology for
244 running workflows in an infrastructure-independent manner. As a lightweight, infrastructure-
245 agnostic containerization technology, Docker allowed PCAWG pipeline authors to fully
246 encapsulate tools and system dependencies into a portable image. This included the fleet of VMs
247 on commercial and academic clouds, as well as the project's HPC clusters that were modified to
248 support Docker containers. Each of our major pipelines was encapsulated in a single Docker
249 image, along with a suitable workflow execution engine, reference data sets, and software libraries
250 (Table 2) .

251 Another key component of the PCAWG software infrastructure stack was cloud-agnostic
252 technology to provision virtual machines on both academic and commercial clouds. Our initial
253 attempts to scale the analytic pipelines across multiple cloud systems were complicated by
254 transient failures in many of the academic cloud environments, subtle differences between
255 seemingly identical clouds, and misconfigured services within the clouds. Initially, we attempted
256 to replicate within the clouds standard components of conventional HPC environments, including
257 shared file systems and cluster load balancing systems. However, we quickly learned that these
258 perform poorly in the dynamic environments of the cloud. After several design iterations, we
259 developed Consonance (<https://github.com/consonance>), a cloud-agnostic provisioning and
260 queueing platform. For each of the cloud platforms in use in PCAWG, including OpenStack,
261 VMWare, AWS, and Azure, Consonance provided a queue where work scheduling was decoupled
262 from the worker nodes. As the fleet of working nodes shrank or expanded, each queue queried the
263 centralized metadata index to obtain the next batch of tasks to execute. Consonance then created
264 and maintained a fleet of worker VMs, launched new pipeline jobs, detected and relaunched failed

265 VMs, and reran workflows as needed. Consonance allowed us to dynamically allocate cloud
266 resources depending on the workload at hand, and even interacted with the AWS spot marketplace
267 to minimize our commercial cloud costs.

268 The Operations: whitelist, work queue, cloud shepherds

269 For the duration of the project, several personnel were required to operate the Docker images,
270 Consonance and the metadata index effectively (Figure 4). Each compute environment was
271 managed by a “cloud shepherd” responsible for completing the workflows on a set of pre-assigned
272 donors or specimens. All the HPC environments (BSC, DKFZ, UTokyo, UCSC, Sanger) were
273 shepherded by personnel local to the institute who were already familiar with the specific file
274 systems and work schedulers, and obtained technical support from their local system
275 administrators. The majority of the cloud environments (AWS, Azure, DKFZ, EMBL-EBI, ETRI,
276 OICR, PDC) granted tenancy to OICR whose personnel acted as cloud shepherds. The other clouds
277 (iDASH, SB), newly launched at the time, assigned their own cloud shepherds who also tested and
278 fine tuned their environments in the process.

279 A project manager acted as the point of contact for all the cloud shepherds to report any technical
280 issues and progress, such that the overall availability of compute resources and throughput at any
281 time point could be estimated. Combining this knowledge with the information from the
282 centralized metadata index, the project manager assigned donors and workflows to compute
283 environments in the form of “whitelists” on a weekly basis. Cloud shepherds then added the
284 whitelist of donors to their workflow queue for execution. This approach allowed us to be agile in
285 responding to data availability disruptions, planned or unplanned downtime while optimizing data
286 transfer and operations throughput.

287 While quotas shifted throughout the duration of the analysis, as demands and workloads on the
288 individual centers changed, the overall peak commitment received was on the order of the 15,000
289 cores, approximately 60TB of RAM, and a peak usage of ~630 virtual machines.

290 Software Distribution through Dockstore

291 The workflows used during PCAWG production include several PCAWG-specific elements that
292 may limit their usability by researchers outside of the project. To facilitate the long term usage of
293 these workflows by a broad range of cancer genomic researchers, we have simplified the tools to
294 make most workflows standalone (Suppl Table 4). These Docker-packaged workflows have been
295 extensively tested for their reproducibility and are registered on the Dockstore²¹
296 (<http://dockstore.org>), a service compliant with Global Alliance for Genomics and Health
297 (GA4GH) standards to provide computational tools and workflows through Docker and described
298 with Common Workflow Language²² (CWL). This enables other researchers to run the workflows
299 on their own data, extend their utility, and replicate the work we have done in any CWL-compliant
300 environment. By running the identical PCAWG workflows on their own data, researchers will be
301 able to make direct comparisons and add to the existing PCAWG dataset.

302 The Docker-packaged BAM alignment and variant calling workflows were tested in different
303 cloud environments and found to be easy to enact by third parties. Some discrepancies with the
304 official data were observed and attributed to improvements in the underlying software (Sanger,
305 Delly) or to the stochastic nature of the software, and deemed to have a low overall impact. Despite
306 not achieving a completely identical results, the reproducibility of the process is satisfactory,
307 especially considering that it involves software developed independently by different teams.

308

309 Data Distribution / Data Portal

310 While GNOS was used for the core pipelines, Synapse²³ was used to provide an interface to the
311 files generated by the working groups and other intermediate results created throughout the project.
312 Unlike GNOS which is focused on archival storage, Synapse allowed for collective editing in the
313 form of a wiki, provenance tracking and versioning of results through a web interface as well as
314 programmatic APIs. While Synapse provided an interface that allowed analyses to be shared
315 rapidly across the consortia, the controlled access data was stored on a secure SFTP server
316 provided by the National Cancer Institute (NCI). When the working groups complete their
317 analysis, the metadata is retained in Synapse while the final version of the results is transferred to
318 the ICGC Data Portal for archival.

319 In addition to GNOS-based repositories, the PCAWG dataset has been mirrored to multiple
320 locations: the European Genome-phenome Archive (EGA,
321 <https://www.ebi.ac.uk/ega/studies/EGAS00001001692>), AWS Simple Storage Service (S3,
322 <https://dcc.icgc.org/icgc-in-the-cloud/aws>), and the Cancer Genome Collaboratory
323 (<http://cancercollaboratory.org>). The data holdings at each repository at the time of publication are
324 summarized in Suppl Table 2. To help researchers locate the PCAWG data, the ICGC Data Portal
325 (<https://dcc.icgc.org>) provides a faceted search interface to query about donor, cancer type, data
326 type or data repositories. Users can browse the collection of released PCAWG data and generate
327 a manifest that facilitates downloading of the selected files.

328 The data repositories hosted at AWS S3 and the Collaboratory are powered by an open source
329 object-based ICGC Storage System (<https://github.com/icgc-dcc/dcc-storage>) that enables fast,
330 secure and multi-part downloads of files. Since AWS and the Collaboratory also have compute
331 power co-located with the PCAWG data, they serve as effective cloud resources for researchers

332 wishing to conduct further analyses on the PCAWG data without having to provision local
333 compute resources and to download terabytes of data to their local compute environment.

334 **Discussion: Replicating PCAWG Analysis on Your Own Data**

335 This project provided us with a rare opportunity to directly compare three categories of compute
336 environment: traditional HPC, academic compute clouds and commercial clouds. In terms of
337 stability and first time setup effort, we found that the traditional HPC environment routinely
338 outperformed academic cloud systems, and often outperformed the commercial clouds. However,
339 most of the academic cloud systems we worked with had been recently installed and some of the
340 stability issues resulted from the shake-down period. The major benefit of the commercial clouds
341 was the ability to scale compute resources up or down as needed, the ease of replicating the setup
342 in different regions, and the availability of cloud-based data centers in different geographic
343 regions, which allowed us to minimize data transfer overhead. For groups interested in replicating
344 PCAWG results, or using the analytic pipelines for their own data, we are comfortable
345 recommending running the analysis on a commercial cloud.

346 In terms of cost, we have summarized in Figure 5 the costs of computing on AWS and the tradeoff
347 in accuracy if running a subset of the variant calling pipelines. The cost of aligning one normal
348 specimen and one tumor specimen, and running three variant calling workflows followed by the
349 OxoG workflow is about \$100 per donor. This is based on a mean WGS coverage of 30X for
350 normal specimens, and a bimodal coverage distribution with maxima at 38X and 60X for tumor
351 specimens²⁴. In addition, the hourly rate of the VMs are approximated from the spot instance
352 pricing we experienced during production runs. With three variant calling workflows, we achieved
353 an F1 score of 0.92. If one is willing to sacrifice some accuracy in order to reduce costs, then

354 running only one variant calling workflow may be an option. Despite the higher costs, running two
355 workflows does not result in increased accuracy. Unfortunately, we were not able to directly
356 compare the analysis costs among commercial clouds, academic clouds and HPC due to the
357 difficulty in assessing the fully loaded cost of provisioning and running an academic compute
358 cluster.

359 In terms of time, the major benefit of operating on commercial clouds is the availability of ample
360 resources for simultaneous parallel runs. For example, in a scenario to analyze a total of 100
361 donors, one runs 200 VMs each aligning one tumor or normal specimen, followed by 300 VMs
362 each running one of the three variant calling workflows on one donor, and 100 VMs to run OxoG
363 workflow, the analysis will in principle take under 9 days to complete. In practice, additional time
364 must be allowed for testing, scaling up, and the inevitability of failed jobs. A more realistic
365 estimate of the time taken to run 100 donors through the complete PCAWG analysis on a
366 commercial cloud is a few weeks.

367 Another issue when planning a large-scale genome analysis project is the variance in execution
368 time from donor to donor. The variant calling pipelines took between 40 and 65 hours of wall time
369 to complete a tumor/genome pair, with the EMBL/DKFZ pipeline running the quickest and the
370 Broad and Sanger pipelines taking somewhat longer. In addition to the variant calling step, the
371 Broad pipeline was preceded by a GATK co-cleaning process taking an additional 24 hours. For
372 each pipeline there was significant variation in the runtime taken for each genome, and some
373 tumor/normal pairs required an excessive amount of time to complete. Because long-running jobs
374 can have economic and logistic impacts, we investigated the cause of this variation by applying
375 linear regression to a number of features describing the raw sequencing sets, including coverage,
376 read quality and mapping scores, number of mismatched end pairs and others (data not shown).

377 We found that a single factor, genomic coverage, explained the variation in wall clock time which
378 increased roughly linearly with coverage.

379 In conclusion, we tackled the challenge of performing uniform analysis on a large dataset across a
380 geographically and technologically disparate collection of compute resources by developing
381 technologies that realized the efficiencies of moving algorithms to the data. This is becoming a
382 necessity as genomic datasets continue to increase in size and are geographically distributed with
383 some jurisdictions restricting the geographical storage and computing of specific datasets. Our
384 approach serves as a model for large scale collaborative efforts that engage many organizations
385 and spread the computation work around the globe.

386 Our effort resulted in three key deliverables. First and foremost, we produced a high-quality,
387 validated consensus variant and alignment dataset of 2,834 cancer donors. To date, this is the
388 largest whole genome cancer dataset analyzed in a consistent and uniform way. The dataset formed
389 the basis for the research by the PCAWG working groups, and will continue to provide value to
390 the research community for many years into the future. Second, we produced a series of best-
391 practice analytical workflows that are portable through the use of Docker and are available on the
392 Dockstore. These workflows are usable in a multitude of compute environments giving researchers
393 the ability to replicate our analysis on their own data. Finally, the infrastructure we built to
394 coordinate analyses between cloud and HPC environments will be helpful for other projects
395 requiring the same distributed approaches.

396 **Acknowledgements**

397 The authors would like to acknowledge the donation of the following compute resources: the
398 PRACE Research Infrastructure resource MareNostrum3 at Barcelona Supercomputing Center

399 with technical expertise provided by the Red Española de Supercomputación and funding support
400 by the Spanish Ministry of Health, ISCIII, in the project Instituto Nacional de Bioinformática
401 (PRB2: PT13/0001/0028); the Cancer Genome Collaboratory, jointly funded by the Natural
402 Sciences and Engineering Research Council of Canada, the Canadian Institutes of Health
403 Research, Genome Canada, and the Canada Foundation for Innovation, and with in-kind support
404 from the Ontario Research Fund of the Ministry of Research, Innovation and Science through the
405 Discovery Frontiers: Advancing Big Data Science in Genomics Research program (grant no.
406 RGPGR/448167-2013); the EMBL-EBI Embassy Cloud supported by UK's (BBSRC) Large
407 Facilities Capital Fund and Cancer Research UK's EMBL-EBI Bioinformatics Resource (grant
408 no. C32939/A20952); sFTP server provided by the Center for Biomedical Informatics &
409 Information Technology (CBIIT) at National Cancer Institute; infrastructure at the Ontario
410 Institute for Cancer Research funded by the Government of Ontario and the Canada Foundation
411 for Innovation (Project #21039); ETRI's OpenStack supported by Institute for Information &
412 communications Technology Promotion with funding from the Korea government (MSIP)
413 (No.B0101-15-0104, The Development of Supercomputing System for the Genome Analysis),
414 Ministry of Health & Welfare, Republic of Korea (grant no: HI14C0072), Korean national research
415 foundation (grant no NRF-2017R1A2B2012796, NRF-2016R1D1A1B03934110), and generous
416 support from Wan Choi and Kwang-Sung; 'Shirokane'_provided by Human Genome Center, the
417 Institute of Medical Science, the University of Tokyo along with technical assistance from Hitachi,
418 Ltd.; Microsoft Azure contributed through a grant to the UC Santa Cruz Genomics Institute and
419 supported by the National Human Genome Research Institute of the National Institutes of Health
420 (grant no U54HG007990) and NCI ITCR (grant no 1R01CA180778); iDASH HIPAA cloud which

421 is a member of the NIH/NHLBI National Centers for Biomedical Computing (U54HL108460) to
422 UC San Diego Health Sciences, Department of Biomedical Informatics.

423 In addition, the Broad team was supported by G.G. funds at MGH and Broad Institute. The DKFZ
424 team was supported by the BMBF-funded Heidelberg Center for Human Bioinformatics (HD-
425 HuB) within the German Network for Bioinformatics Infrastructure (de.NBI) (#031A537A,
426 #031A537C) and the BMBF-funded grants ICGC PedBrain (01KU1201A, 01KU1201B), ICGC
427 EOPC (01KU1001A), ICGC MMML-seq (01KU1002B), and ICGC DE-MINING (01KU1505E).
428 Variant calling with the DKFZ/EMBL pipeline made use of the Roddy framework, and provision
429 of data and metadata of the German ICGC projects was assisted by the One Touch Pipeline (OTP).
430 The OICR team was funded by the Government of Ontario and the Canada Foundation for
431 Innovation (Project #21039). The Sanger team was supported by the Wellcome Trust grant
432 (098051) with contributions by Shriram G Bhosle, David R Jones, Andrew Menzies, Lucy
433 Stebbings, Jon W Teague.

434

435 **References**

- 436 1. Network, T.C.G.A.R. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nature*
437 *Genetics* **45**, 1113-1120 (2013).
- 438 2. PCAWG-3. Pan-Cancer Study of Recurrent and Heterogeneous RNA Aberrations and
439 Association with Whole-Genome Variants. (in preparation).
- 440 3. Alioto, T.S. *et al.* A comprehensive assessment of somatic mutation detection in cancer
441 using whole-genome sequencing. *Nat Commun* **6**, 10001 (2015).
- 442 4. PCAWG-1. Consistent Detection of Short Somatic Mutations in 2,778 Cancer Whole
443 Genomes. (in preparation).

- 444 5. Phillips, M. & Knoppers, B. Building an International Code of Conduct for Genomic Cloud
445 Research. (in preparation).
- 446 6. Wilks, C. *et al.* The Cancer Genomics Hub (CGHub): overcoming cancer through the
447 power of torrential data. *Database (Oxford)* **2014**(2014).
- 448 7. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
449 (2013).
- 450 8. Jones, D. *et al.* cgpCaVEManWrapper: Simple Execution of CaVEMan in Order to Detect
451 Somatic Single Nucleotide Variants in NGS Data. *Curr Protoc Bioinformatics* **56**, 15.10.1-
452 15.10.18 (2016).
- 453 9. Raine, K.M. *et al.* cgpPindel: Identifying Somatically Acquired Insertion and Deletion
454 Events from Paired End Sequencing. *Curr Protoc Bioinformatics* **52**, 15.7.1-12 (2015).
- 455 10. Raine, K.M. *et al.* ascatNgs: Identifying Somatically Acquired Copy-Number Alterations
456 from Whole-Genome Sequencing Data. *Curr Protoc Bioinformatics* **56**, 15.9.1-15.9.17 (2016).
- 457 11. BRASS. (<https://github.com/cancerit/BRASS>).
- 458 12. Rimmer, A. *et al.* Integrating mapping-, assembly- and haplotype-based approaches for
459 calling variants in clinical sequencing applications. *Nat Genet* **46**, 912-8 (2014).
- 460 13. Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end and split-
461 read analysis. *Bioinformatics* **28**, i333-i339 (2012).
- 462 14. Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and
463 heterogeneous cancer samples. *Nat Biotechnol* **31**, 213-9 (2013).
- 464 15. Fan, Y. *et al.* MuSE: accounting for tumor heterogeneity using a sample-specific error
465 model improves sensitivity and specificity in mutation calling from sequencing data. *Genome Biol*
466 **17**, 178 (2016).

- 467 16. Costello, M. *et al.* Discovery and characterization of artifactual mutations in deep coverage
468 targeted capture sequencing data due to oxidative DNA damage during sample preparation.
469 *Nucleic Acids Res* **41**, e67 (2013).
- 470 17. Wala, J., Zhang, C.Z., Meyerson, M. & Beroukhi, R. VariantBam: filtering and profiling
471 of next-generation sequencing data using region-specific rules. *Bioinformatics* **32**, 2029-31
472 (2016).
- 473 18. PCAWG-6. PCAWG-6 paper. (in preparation).
- 474 19. PCAWG-11. PCAWG-11 paper. (in preparation).
- 475 20. O'Connor, B.D., Merriman, B. & Nelson, S.F. SeqWare Query Engine: storing and
476 searching sequence data in the cloud. *BMC Bioinformatics* **11 Suppl 12**, S2 (2010).
- 477 21. O'Connor, B.D. *et al.* The Dockstore: enabling modular, community-focused sharing of
478 Docker-based genomics tools and workflows. *FI000Res* **6**, 52 (2017).
- 479 22. Amstutz, P. *et al.* Common Workflow Language, v1.0. *figshare* (2016).
- 480 23. Omberg, L. *et al.* Enabling transparent and collaborative computational analysis of 12
481 tumor types within The Cancer Genome Atlas. *Nat Genet* **45**, 1121-6 (2013).
- 482 24. PCAWG-QC. Framework for quality assessment of whole genome, cancer sequences. (in
483 preparation).

484

485 **Additional Members of the PCAWG Technical Working Group**

486 Javier Bartolomé Rodríguez¹, Keith A. Boroevich², Rich Boyce³, Angela N. Brooks⁴, Alex
487 Buchanan⁵, Ivo Buchhalter^{6,7}, Niall J. Byrne⁸, Andy Cafferkey⁹, Peter J. Campbell¹⁰, Zhaohong
488 Chen¹¹, Sunghoon Cho¹², Wan Choi¹³, Peter Clapham¹⁴, Francisco M. De La Vega^{15,16}, Jonas
489 Demeulemeester^{17,18}, Michelle T. Dow¹⁹, Lewis J. Dursi^{8,20}, Juergen Eils²¹, Claudiu Farcas²²,
490 Francesco Favero²³, Nodirjon Fayzullaev⁸, Paul Flicek³, Nuno A. Fonseca³, Josep L.I. Gelpi^{24,25},
491 Gad Getz^{26,27}, Bob Gibson⁸, Michael C. Heinold^{7,6}, Julian M. Hess²⁶, Oliver Hofmann²⁸, Jongwhi
492 H. Hong²⁹, Thomas J. Hudson^{30,31}, Daniel Huebschmann^{6,7}, Barbara Hutter^{32,33}, Carolyn M.

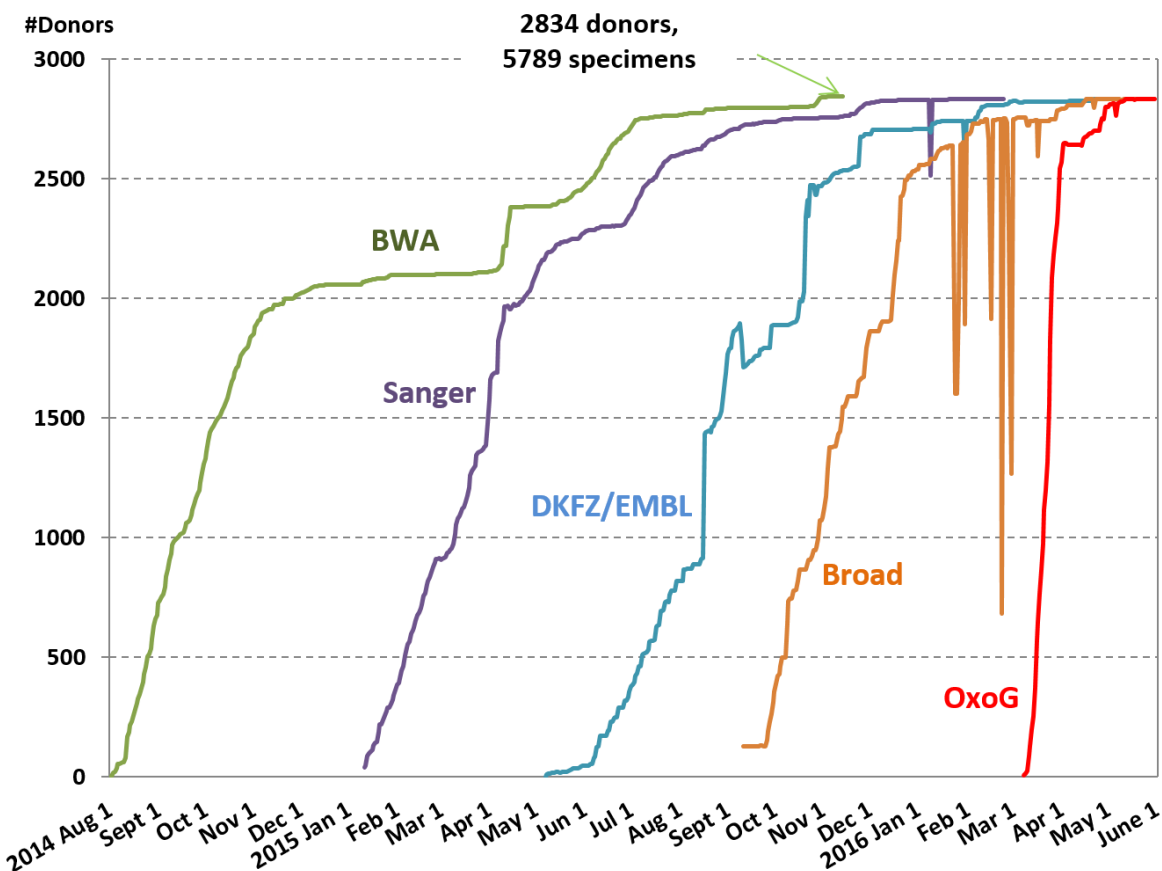
493 Hutter³⁴, Seiya Imoto³⁵, Sinisa Ivkovic³⁶, Seung-Hyup Jeon¹³, Wei Jiao⁸, Jongsun Jung³⁷, Rolf
494 Kabbe⁶, Andre Kahles^{38,39}, Jules Kerssemakers⁴⁰, Hyunghwan Kim¹³, Hyung-Lae Kim^{41,42},
495 Jihoon Kim¹¹, Jan O. Korbel^{43,3}, Michael Koscher⁴⁰, Antonios Koures¹¹, Milena Kovacevic³⁶,
496 Chris Lawrenz⁶, Ignaty Leshchiner²⁶, Dimitri G. Livitz²⁶, George L. Mihaiescu⁸, Sanja
497 Mijalkovic³⁶, Ana Mijalkovic Lazic³⁶, Satoru Miyano⁴⁴, Hardeep K. Nahal⁸, Mia Nastic³⁶,
498 Jonathan Nicholson¹⁴, David Ocana³, Kazuhiro Ohi⁴⁴, Lucila Ohno-Machado²², Larsson
499 Omberg⁴⁵, B.F. Francis Ouellette^{8,46}, Nagarajan Paramasivam^{6,47}, Marc D. Perry⁸, Todd D. Pihl⁴⁸,
500 Manuel Prinz⁶, Montserrat Puiggròs²⁴, Petar Radovic³⁶, Esther Rheinbay^{26,49}, Mara W.
501 Rosenberg^{26,49}, Charles Short³, Heidi J. Sofia⁵⁰, Jonathan Spring⁵¹, Adam J. Struck⁵, Grace
502 Tiao²⁶, Nebojsa Tijanic³⁶, Peter Van Loo^{17,18}, David Vicente¹, Jeremiah A. Wala^{26,52}, Zhining
503 Wang⁵³, Johannes Werner⁶, Ashley Williams¹¹, Youngchoon Woo¹³, Adam J. Wright⁸, Qian
504 Xiang⁸

505

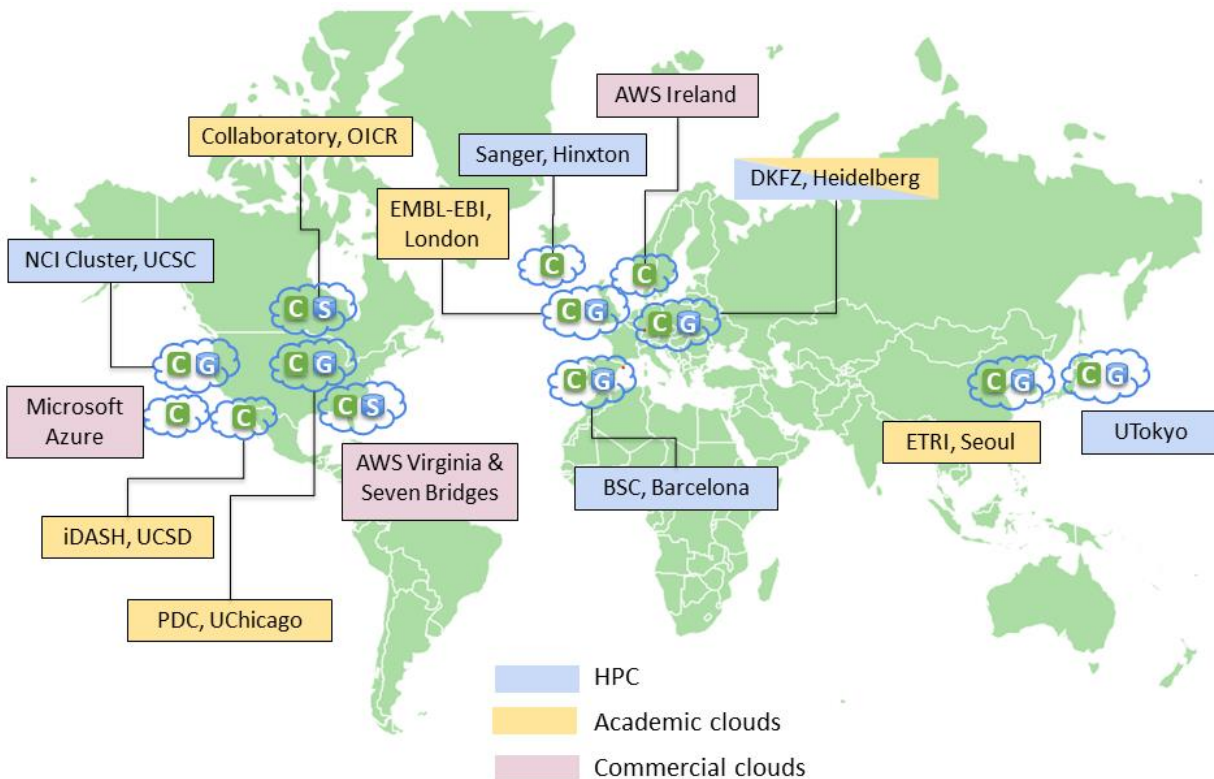
506 ¹Department of Operations, Barcelona Supercomputing Center, Barcelona, Catalunya, 8034, Spain. ²Laboratory for
507 Medical Science Mathematics, RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa, 230-0045,
508 Japan. ³European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton, Cambridge, CB10
509 1SD, United Kingdom. ⁴Biomolecular Engineering, University of California Santa Cruz, Santa Cruz, California,
510 95065, USA. ⁵Department of Computational Biology, Oregon Health and Science University, Portland, Oregon,
511 97239, USA. ⁶Division of Theoretical Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg,
512 Baden-Württemberg, 69120, Germany. ⁷Department for Bioinformatics and Functional Genomics, Institute for
513 Pharmacy and Molecular Biotechnology and BioQuant, Heidelberg University, Heidelberg, Baden-Württemberg,
514 69120, Germany. ⁸Informatics and Biocomputing Program, Ontario Institute for Cancer Research, Toronto, Ontario,
515 M5G 0A3, Canada. ⁹Technical Services Cluster, European Molecular Biology Laboratory, European Bioinformatics
516 Institute, Hinxton, Cambridge, CB10 1SD, United Kingdom. ¹⁰Cancer Genome Project, Wellcome Trust Sanger
517 Institute, Hinxton, Cambridgeshire, CB10 1SA, United Kingdom ¹¹Department of Medicine, University of
518 California San Diego, San Diego, California, 92093, USA. ¹²PDXen Biosystems Inc., Seoul, 4900, South Korea.
519 ¹³Electronics and Telecommunications Research Institute, Daejeon, 34129, South Korea. ¹⁴Informatics Support
520 Group, Wellcome Trust Sanger Institute, Hinxton, Cambridgeshire, CB10 1SA, United Kingdom. ¹⁵Department of
521 Biomedical Data Science, Stanford University School of Medicine, Stanford, California, 94305, USA. ¹⁶Annai
522 Systems, Inc., Carlsbad, California, 92011, USA. ¹⁷The Francis Crick Institute, London, NW1 1AT, United
523 Kingdom. ¹⁸Department of Human Genetics, University of Leuven, B-3000 Leuven, Belgium ¹⁹Biomedical
524 Informatics, University of California San Diego, San Diego, California, 92093, USA. ²⁰The Centre for
525 Computational Medicine, The Hospital for Sick Children, Toronto, Ontario, M5G 0A4, Canada. ²¹Theoretical
526 Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg, Baden-Württemberg, 69120, Germany.
527 ²²Health System Department of Biomedical Informatics, University of California San Diego, La Jolla, California,
528 92093, USA. ²³BRIC/Finsen Laboratory, Rigshospitalet, Copenhagen, 2200, Denmark. ²⁴Department of Life
529 Sciences, Barcelona Supercomputing Center, Barcelona, Catalunya, 8034, Spain. ²⁵Department of Biochemistry and
530 Molecular Biomedicine, University of Barcelona, Barcelona, Catalunya, 8028, Spain. ²⁶Cancer Program, Broad
531 Institute of MIT and Harvard, Cambridge, Massachusetts, 02142, USA. ²⁷Cancer Center and Department of
532 Pathology, Massachusetts General Hospital, Boston, Massachusetts, 02114, USA. ²⁸Center for Cancer Research,
533 University of Melbourne, Melbourne, VIC 3001, Australia. ²⁹Genome Data Integration Center, Syntekabio Inc.,
534 Daejeon, 34025, South Korea. ³⁰Genomics Program, Ontario Institute for Cancer Research, Toronto, Ontario, M5G
535 0A3, Canada. ³¹Oncology Discovery and Early Development, AbbVie, Redwood City, California, 94063, USA.
536 ³²Division of Applied Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg, Baden-Württemberg,
537 69120, Germany. ³³Division of Applied Bioinformatics, National Center for Tumor Diseases, Heidelberg, Baden-
538 Württemberg, 69120, Germany. ³⁴Division of Genomic Medicine, National Human Genome Research Institute,

539 Bethesda, Maryland, 20852, USA. ³⁵Health Intelligence Center, Institute of Medical Science, University of Tokyo,
540 Tokyo, 108-8639, Japan. ³⁶Seven Bridges, Cambridge, Massachusetts, 02142, USA. ³⁷Genome Data Integration
541 Center, Syntekabio Inc., Daejeon, 34025, South Korea ³⁸Department of Computer Science, ETH Zurich, Zurich,
542 Zurich, 8092, Switzerland. ³⁹Computational Biology Center, Memorial Sloan Kettering Cancer Center, New York,
543 New York, 10065, USA. ⁴⁰German Cancer Research Center (DKFZ), Heidelberg, Baden-Württemberg, 69120,
544 Germany. ⁴¹Department of Biochemistry, Ewha Womans University, Seoul, 07985, South Korea. ⁴²PGM21, Seoul,
545 07985, South Korea. ⁴³Genome Biology Unit, European Molecular Biology Laboratory, Heidelberg, Baden-
546 Württemberg, 69120, Germany. ⁴⁴Human Genome Center, Institute of Medical Science, University of Tokyo,
547 Tokyo, 108-8639, Japan. ⁴⁵Systems Biology, Sage Bionetworks, Seattle, Washington, 98112, USA. ⁴⁶Department of
548 Cell and Systems Biology, University of Toronto, Toronto, Ontario, M5S 3G5, Canada. ⁴⁷Medical Faculty
549 Heidelberg, Heidelberg University, Heidelberg, Baden-Württemberg, 69120, Germany. ⁴⁸CSRA Incorporated,
550 Fairfax, Virginia, 22042, USA. ⁴⁹Cancer Center, Massachusetts General Hospital, Boston, Massachusetts, 02114,
551 USA. ⁵⁰National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland, 20892-
552 9305, USA. ⁵¹Center for Data Intensive Science, University of Chicago, Chicago, Illinois, 60637, USA.
553 ⁵²Department of Cancer Biology, Dana-Farber Cancer Institute, Boston, Massachusetts, 02115, USA. ⁵³TCGA
554 Program Office, National Cancer Institute, Bethesda, Maryland, 20892, USA.
555

556 **Figures**



557
558
559 Figure 1: Progress of the 5 workflows over time. The “flat line” of the BWA workflow was due to
560 two major tranches of sequencing data submissions, with a first tranche of ~2000 donors and a
561 second tranche of ~800 donors that were uploaded later. The staggered start of the three
562 variant calling pipelines was dictated more by the time required to develop and package the
563 workflows, and less by the availability of compute power. The “dips” on the plots resulted from
564 quality issues with some sets of variant calls that were withdrawn, reprocessed and resubmitted.
565 In the case of the Broad workflow, the variant calls were withdrawn for post-processing before
566 being considered complete. If all workflows and data would have been in place at the beginning
567 of the project, we estimate the computation across the full set of 5,789 genomes could have
568 been completed in under 6 months.

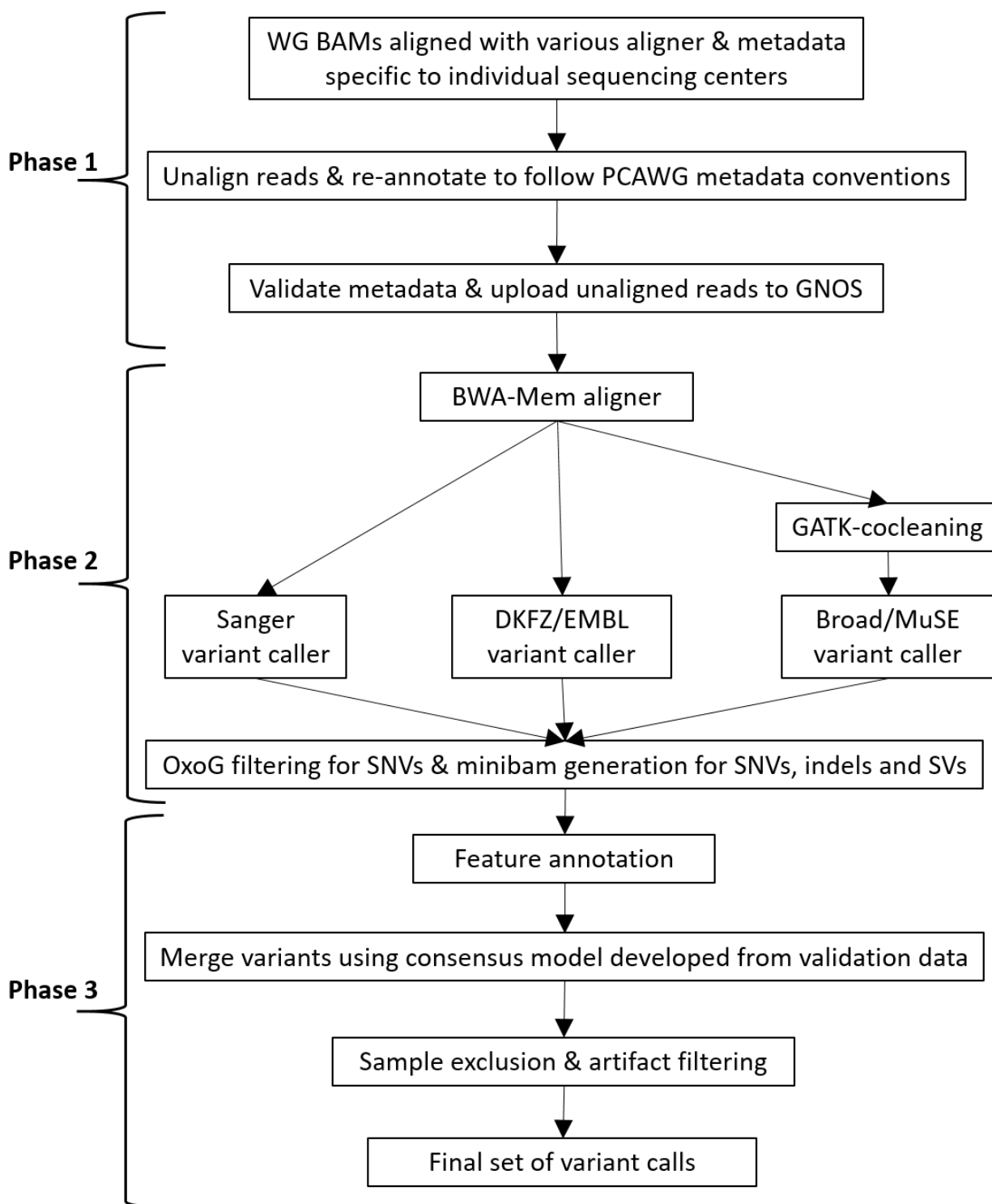


569

570 Figure 2: Geographical distribution of compute centers (C), GNOS servers (G), and

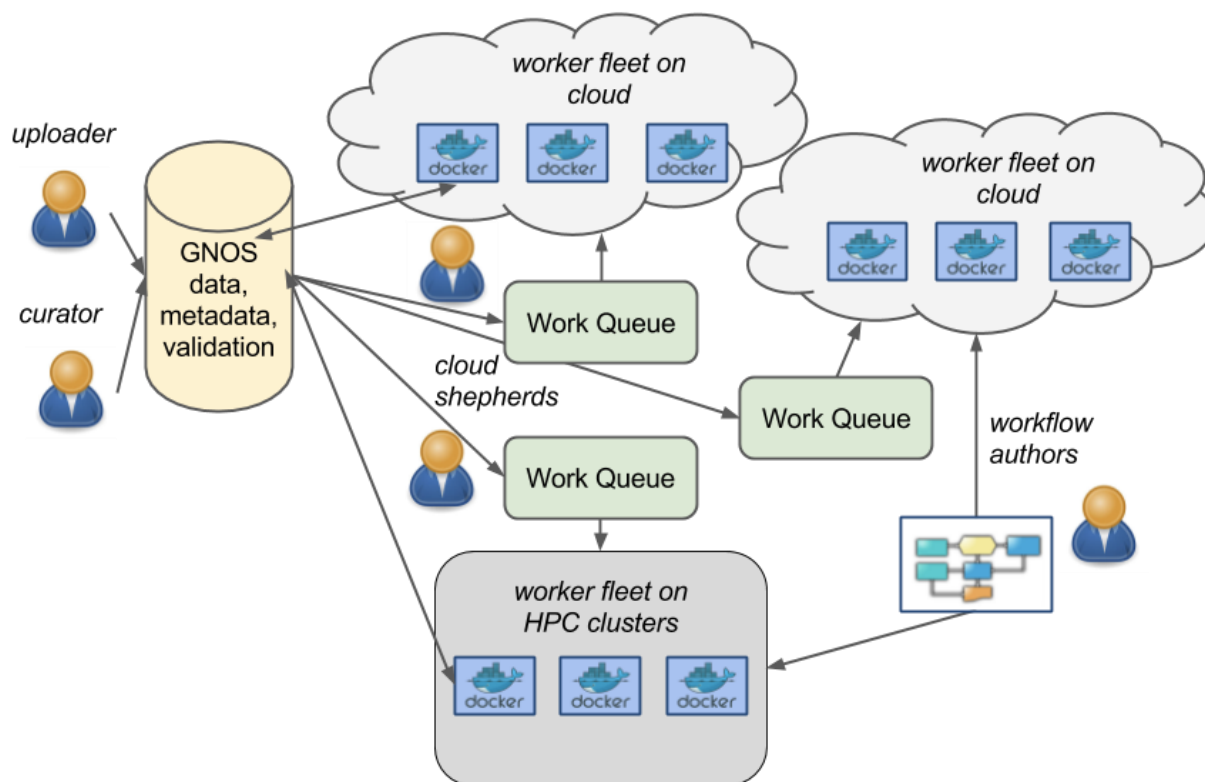
571 S3-compatible data storage (S).

572

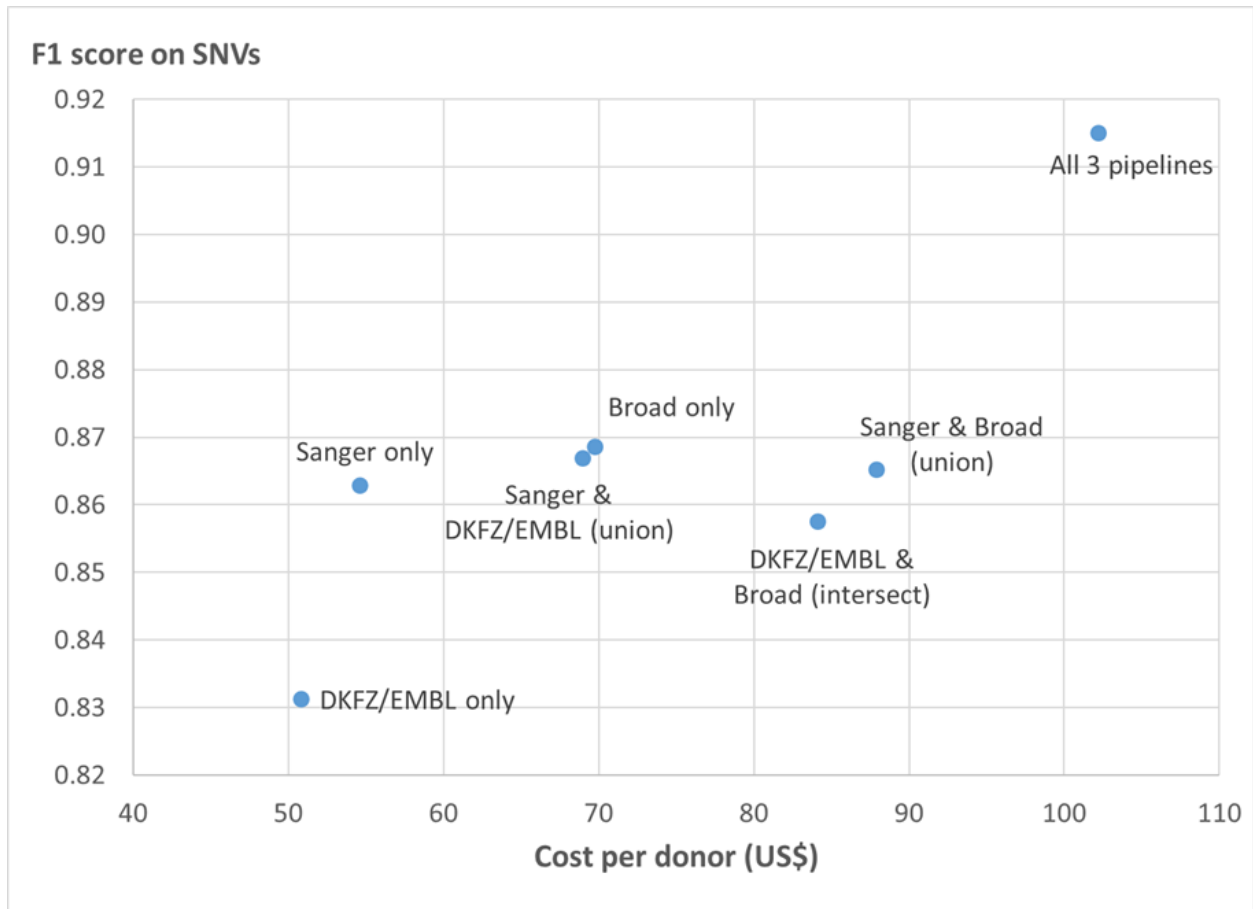


573
574
575
576
577
578
579

Figure 3: The uniform analysis of whole genomes involves three broad phases. Phase 1: Data marshalling and upload. Phase 2: Sequence alignment and variant calling. Phase 3: Variant merging and filtering. The algorithms for merging SNVs and indels are described in the PCAWG-1 paper, SVs in the PCAWG-6 paper, and CNVs in the PCAWG-11 paper.



580
581 Figure 4: Infrastructure used on cloud and HPC compute environments for core analysis.
582



583
584
585
586
587
588
589
590
591

Figure 5: Costs for analyzing a tumor/normal pair through BWA-Mem, different combinations of variant calling pipelines, and OxoG filtering. The cost is calculated based on AWS instances at average spot pricing we experienced during the project, and includes egress costs to transfer the result files. PCAWG ran all 3 variant calling pipelines and achieved an F1 score of 0.9151 for SNVs. If running only one or two pipelines, there will be savings in cost but sacrifice in accuracy. Detailed cost analysis is shown in Suppl Table 3.

592 **Tables**

593

594 Table 1. Compute resources. * Shared between environments. ** Transient storage used for
595 local data processing.

596

	Type	Allocated CPU/Cores	Allocated memory	Data Co-location Repository	Local Storage Amount
AWS	Cloud	variable	variable	Y	420TB
Azure	Cloud	variable	variable	N	-
BSC	HPC	1000	7.75TB	Y	300TB
Collaboratory	Cloud	350	3.2TB	Y	132TB
DKFZ	HPC	800	3.5TB	Y	1.7PB*
DKFZ	Cloud	1024	4TB	Y	1.7PB*
EMBL-EBI	Cloud	1000	4TB	Y	1PB
ETRI	Cloud	800	2TB	Y	750TB
iDASH	Cloud	304	2.8TB	N	9TB**
PDC	Cloud	108	324GB	Y	732TB
Sanger	HPC	1500	12TB	N	750TB**
SBG	Cloud	variable	variable	Y	-
UCSC	HPC	4000	33TB	Y	300TB
UTokyo	HPC	2496	2.5TB	Y	400TB

597

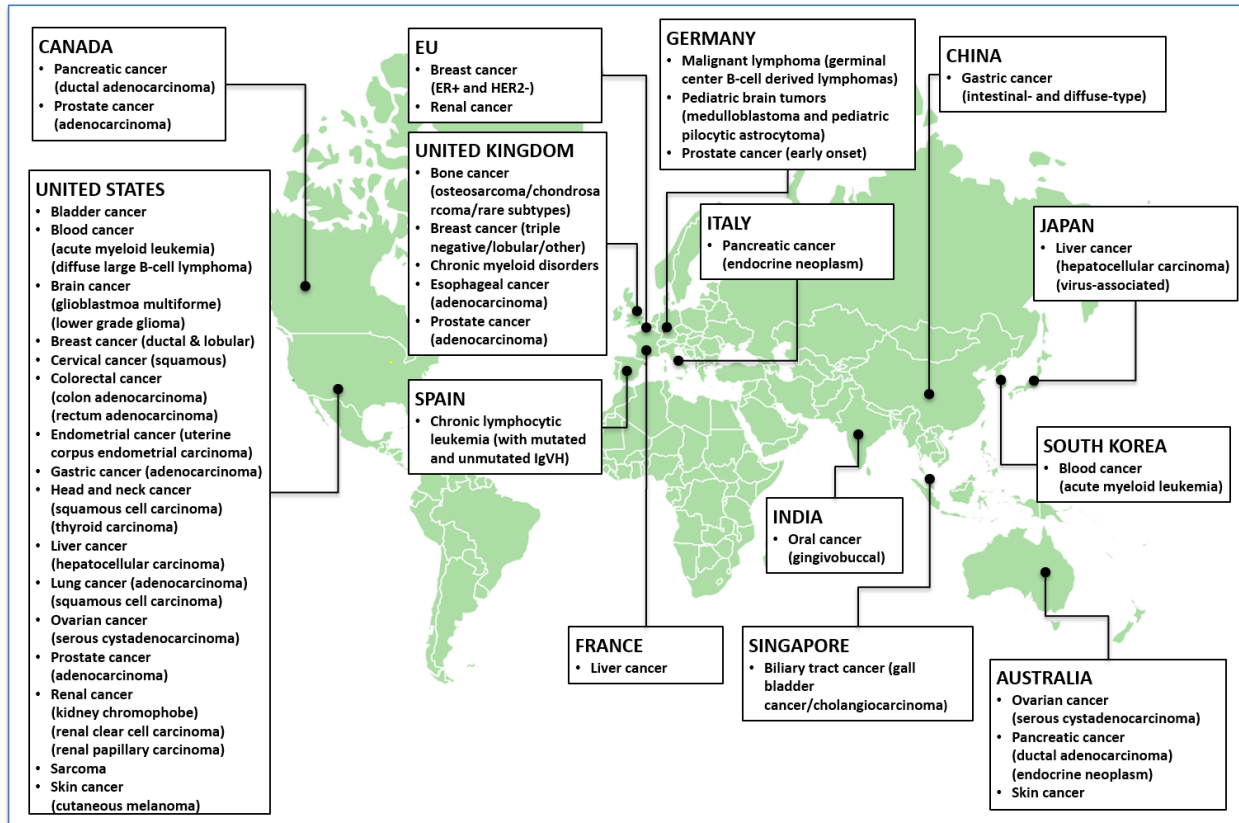
598 Table 2. The five core workflows. Components for calling (1) SNVs, (2) indels, (3) SVs and (4)
 599 SCNAs in each of the three variant calling workflows are listed. Because we utilized a large
 600 number of compute environments with various configurations of cores and RAM, the average
 601 runtime for each pipelines varied with large standard deviations (Suppl Fig. 7-10). The runtime
 602 for the Broad pipeline included the 24 hours required to run GATK co-cleaning of BAMs. The
 603 measured runtime included time to download input files, but not the time to upload result files.
 604 (#) MuSE was developed at MD Anderson Cancer Center and Baylor College of Medicine.
 605

	BWA	Sanger	DKFZ/EMBL	Broad	OxoG
Analytical components in workflow	BWA-Mem Picard Biobambam samtools	CaVEMan ¹ cgpPindel ² BRASS ³ ascatNgs ⁴	dkfz_snv ¹ Platypus ² DELLY ³ ACE-seq ⁴	GATK cocleaning MuTect ¹ MuSE ^{1,#} Snowman ^{2,3} dRanger ³	OxoG VariantBam
Workflow controller	SeqWare	SeqWare	Roddy, SeqWare	Galaxy	SeqWare
Recommended compute requirements	4 cores, 15GB RAM	16 cores, 4.5GB RAM/core	16 cores, 64GB RAM	32 cores, 244GB RAM	8 cores, 64GB RAM
Average runtime across all compute environments	2.0 +/- 1.7 days	5.3 +/- 5.5 days	3.2 +/- 1.7 days	5.1 +/- 2.2 days	2.6 +/- 1.3 hours
Benchmark on AWS	5.8 days on 4-core m1.xlarge	2.2 days on 32-core r3.8xlarge	1.7 days on 32-core r3.8xlarge	3.7 days on 32-core r3.8xlarge	4 hours on 8-core m2.4xlarge
Core hours per run	557	1690	1306	2842	32
Output files per run	120GB	2 GB	5 GB	35 GB	1.5 GB

606

607 **Supplementary Information**

608

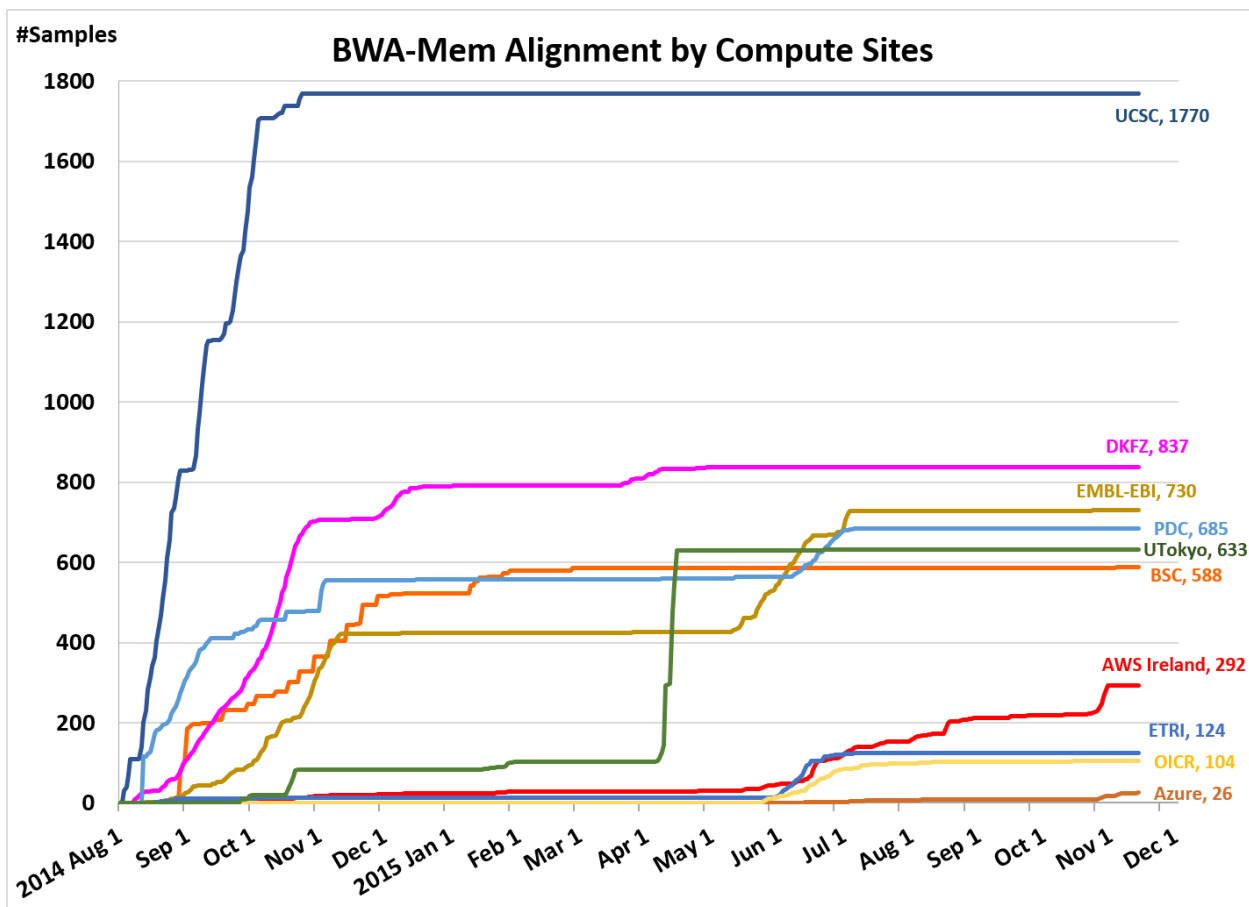


609

610

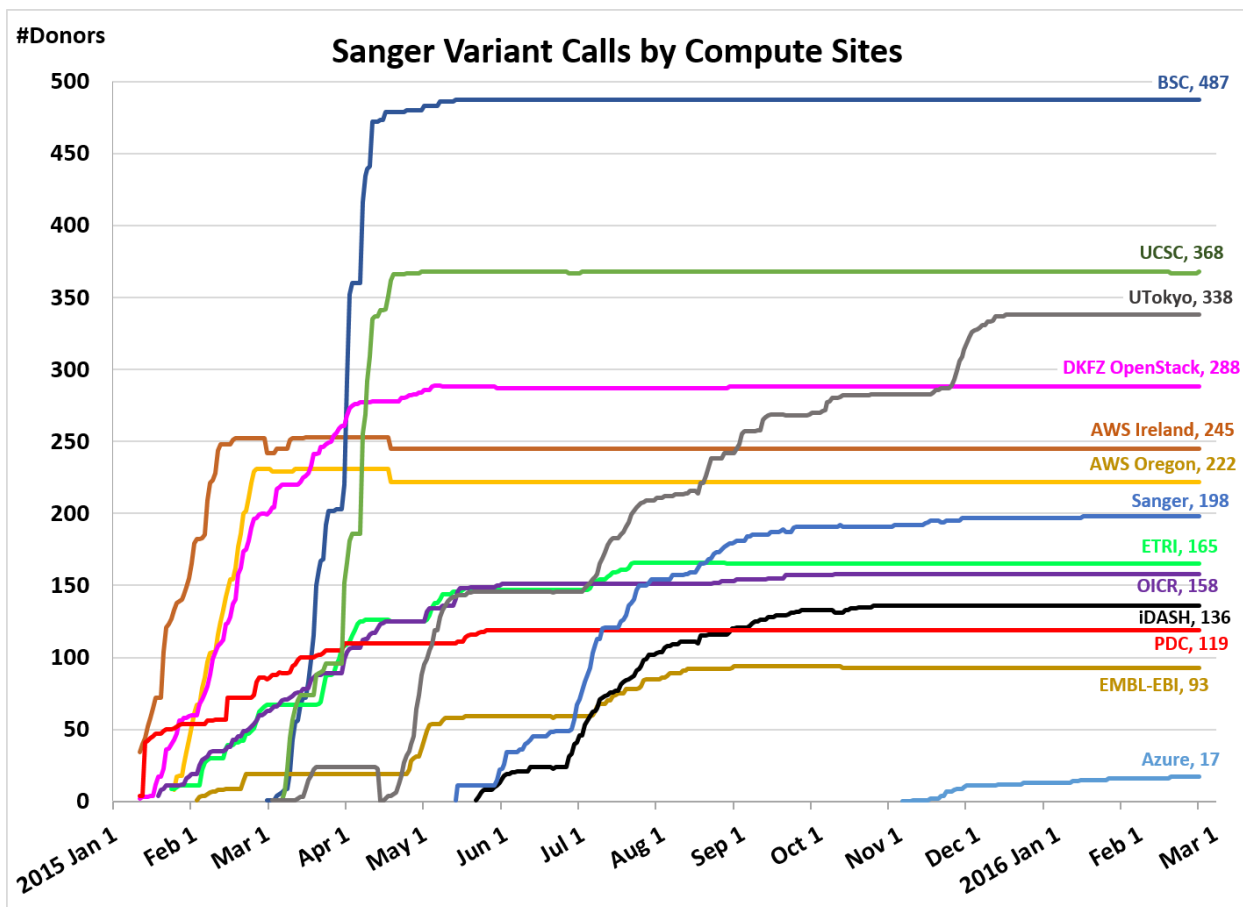
611 Supplementary Figure 1: Whole genomes from 2,834 donors across 39 cancer types were
612 collected from 48 ICGC and TCGA projects in 14 jurisdictions.

613



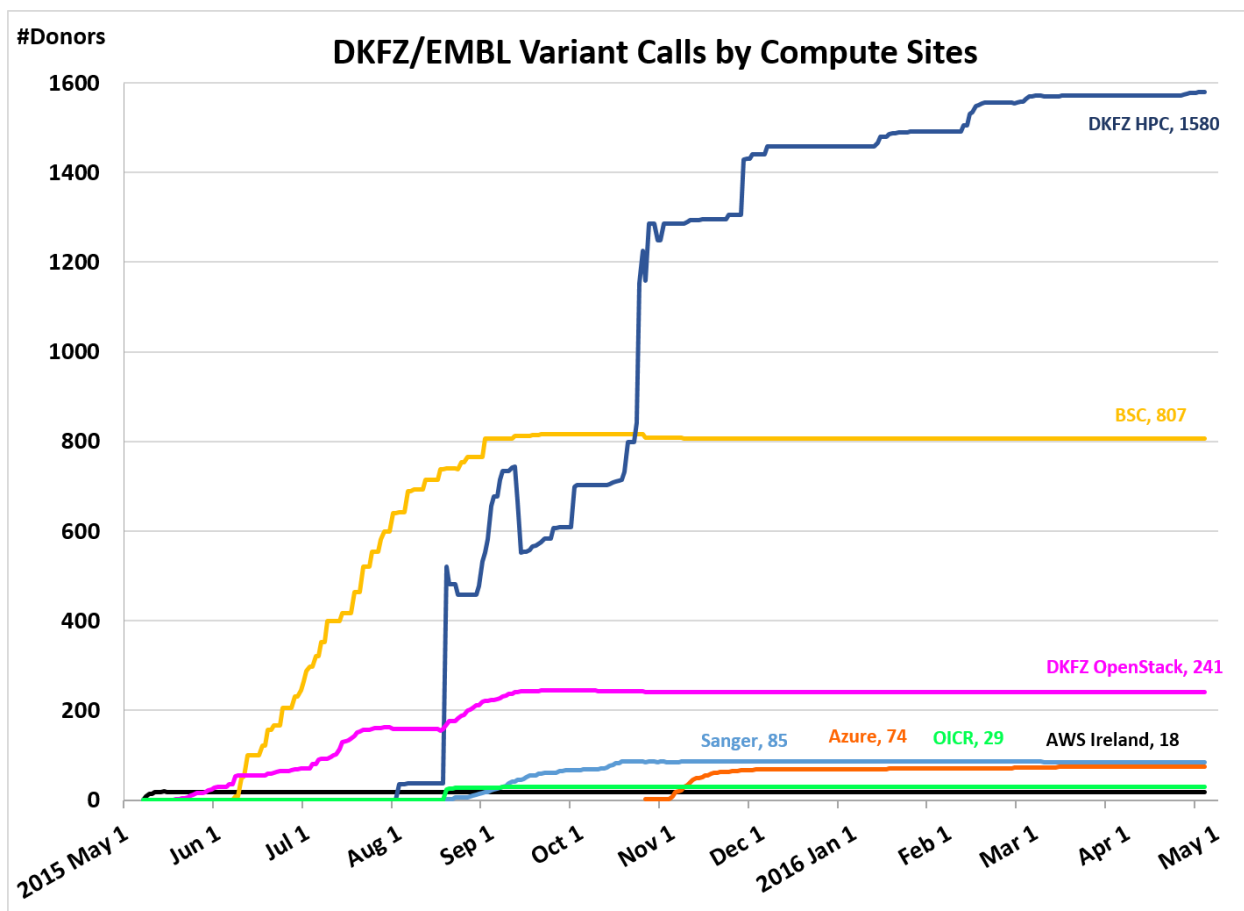
614
615
616

Supplementary Figure 2: Progress of BWA-Mem alignment over time at 7 compute sites.



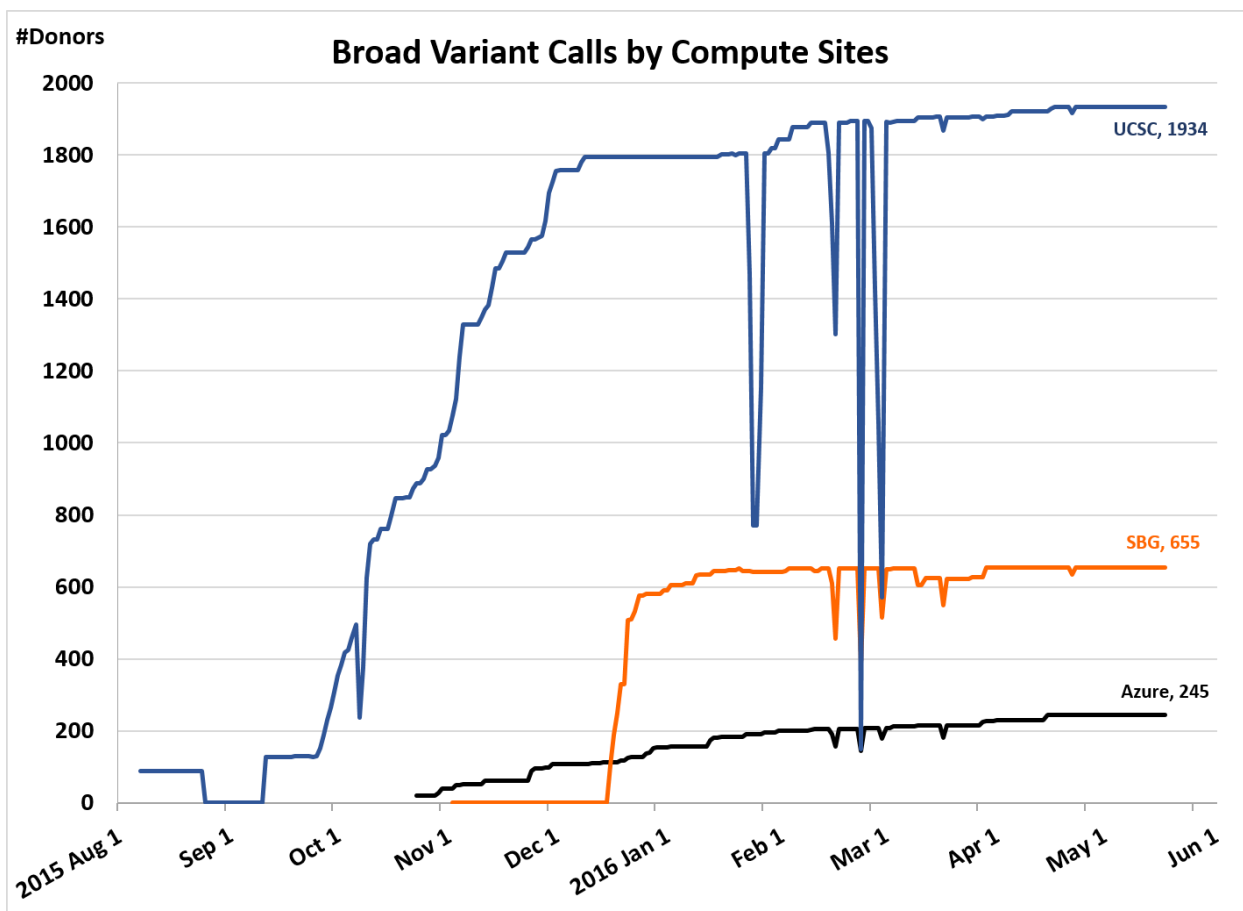
617
618
619
620

Supplementary Figure 3: Progress of Sanger variant calling workflow over time at 13 compute sites.



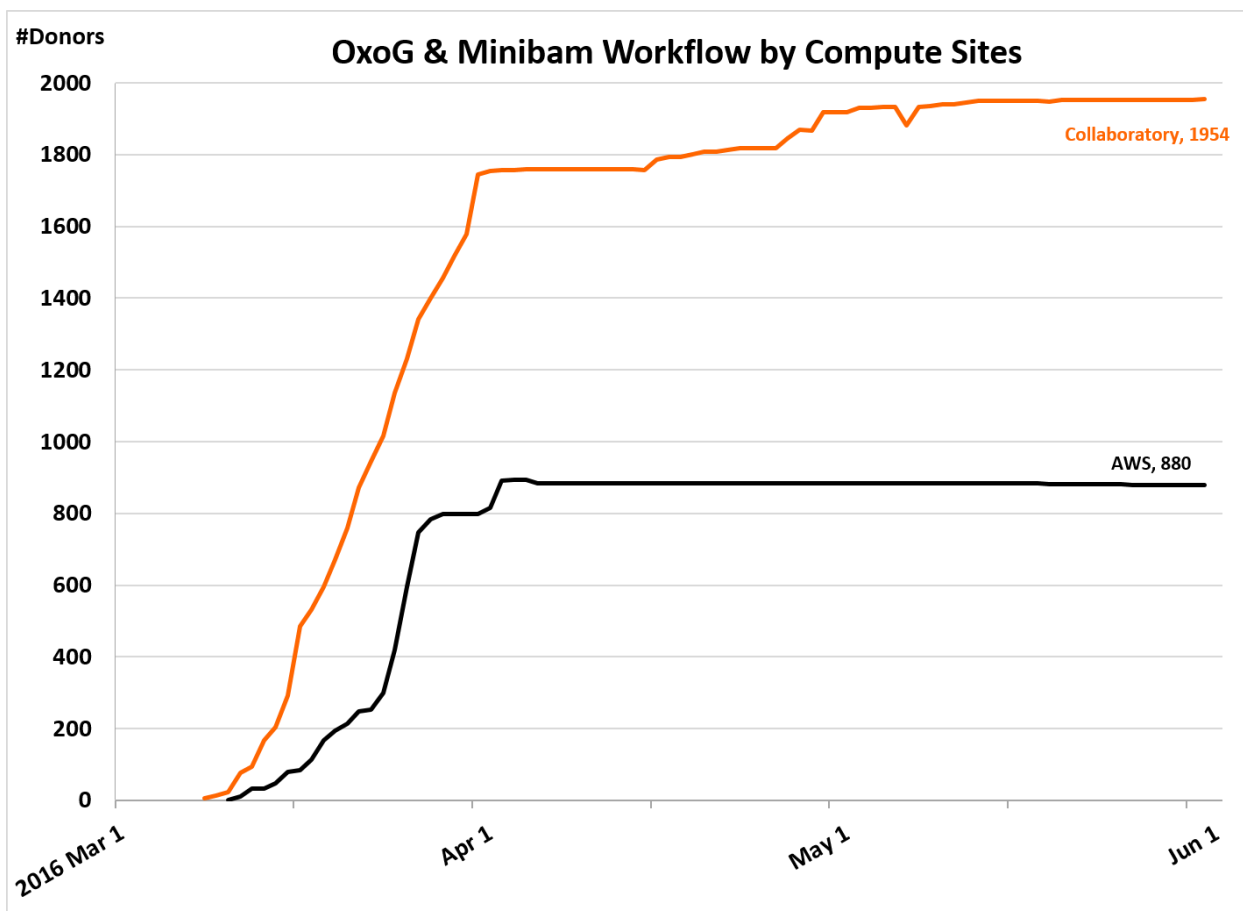
621
622
623
624

Supplementary Figure 4: Progress of DKFZ/EMBL variant calling workflow over time at 7 compute sites.



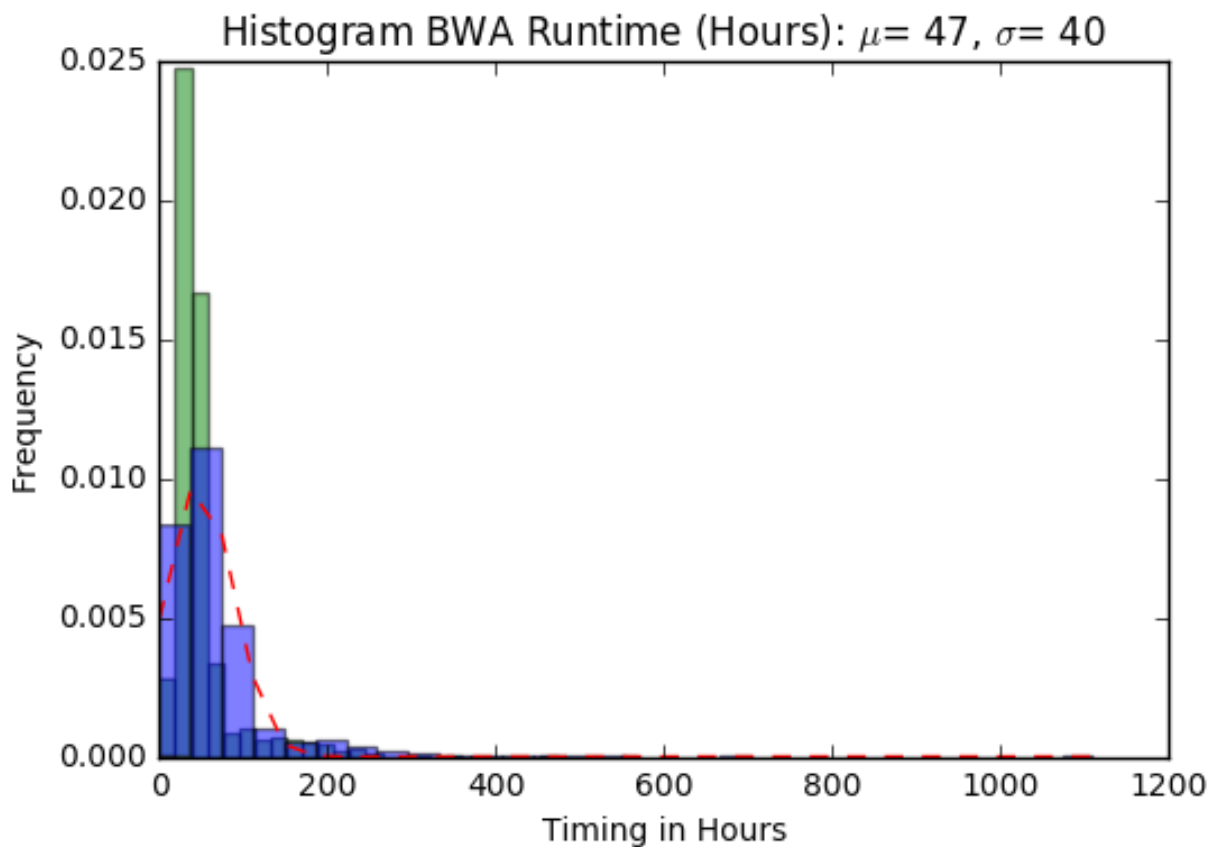
625
626
627
628

Supplementary Figure 5: Progress of Broad variant calling workflow over time at 3 compute sites.



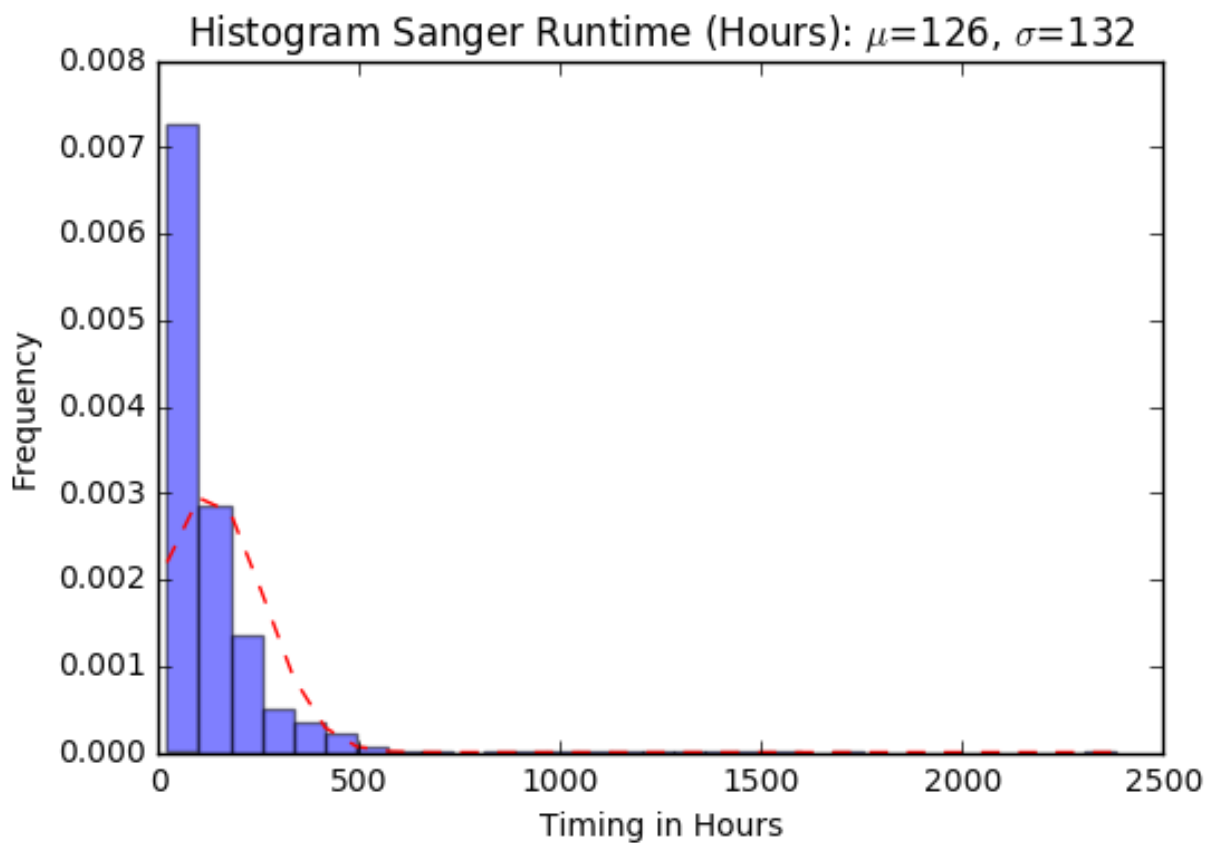
629
630
631

Supplementary Figure 6: Progress of OxoG and minibam workflow over time at 2 compute sites.



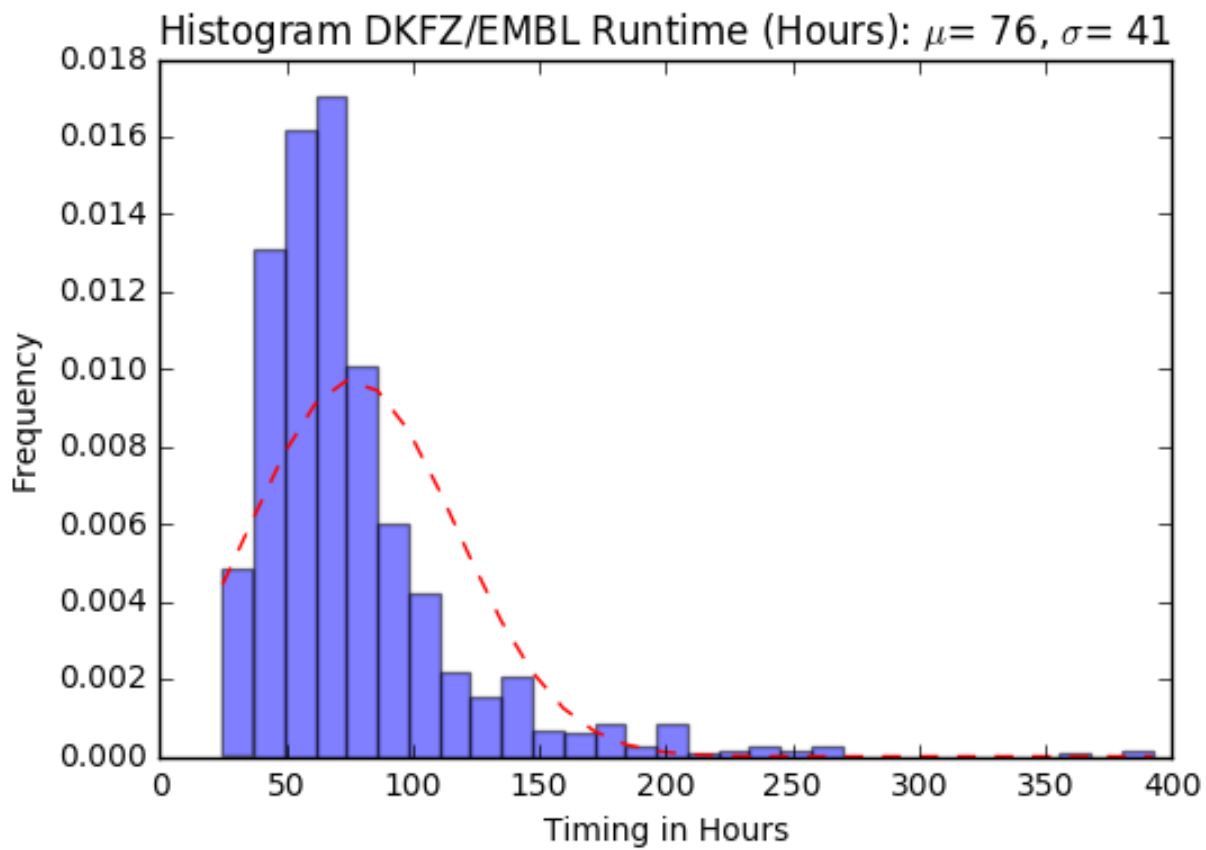
632
633
634

Supplementary Figure 7: Average runtimes for BWA-Mem alignment workflow



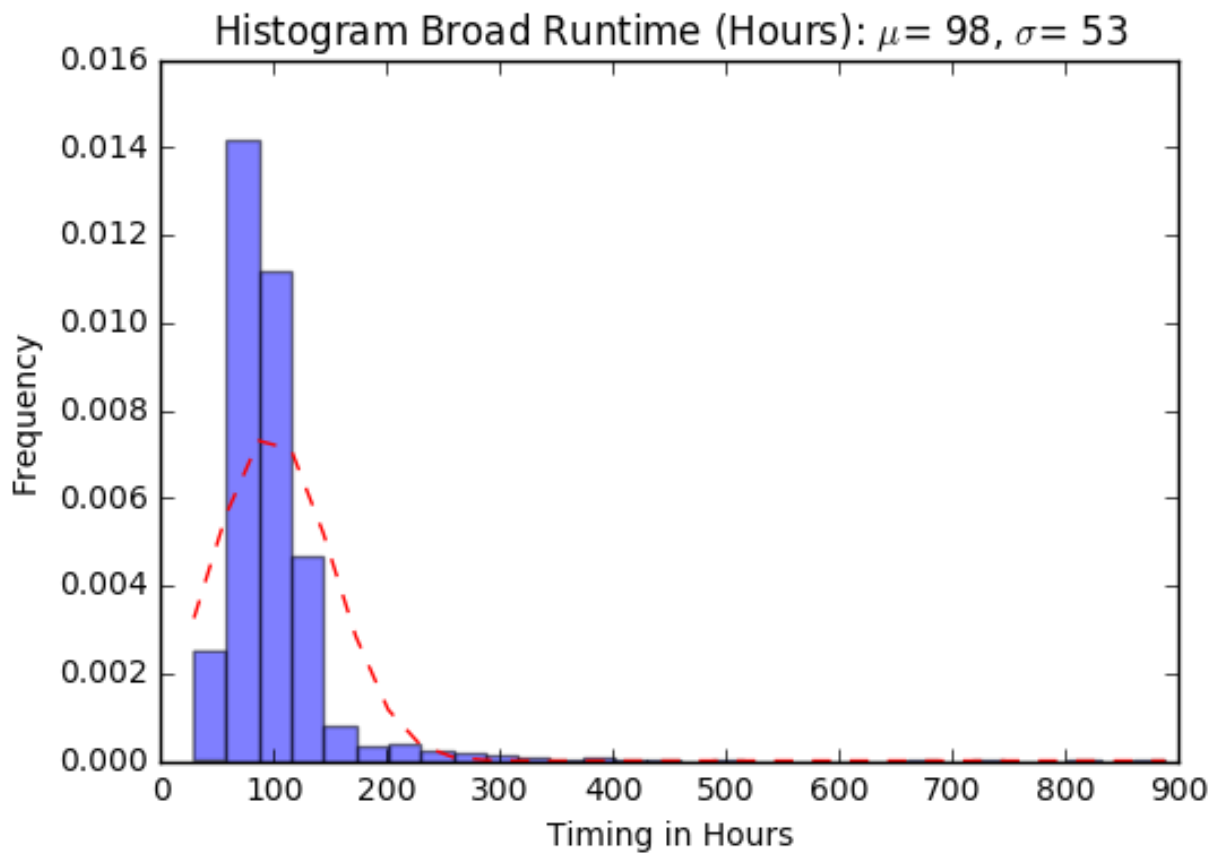
635
636
637

Supplementary Figure 8: Average runtime for the Sanger somatic variant calling workflow.



638
639
640

Supplementary Figure 9: Average runtime for the DKFZ/EMBL somatic variant calling workflow.



641
642 Supplementary Figure 10: Average runtime for the Broad somatic variant calling workflow.
643 Preceding the variant calling workflow, the GATK co-cleaning step takes an additional 24 hours.
644

645 Supplementary Table 1. Percentage samples/donors run at each site for each pipeline
646

	BWA	Sanger	DKFZ/EMBL	Broad/MuSE	OxoG
AWS Ireland	5.0	16.4	0.6		31.1
Azure	0.4	0.6	2.6	8.6	
BSC	10.2	17.2	28.5		
Collaboratory					68.9
DKFZ (HPC)			55.8		
DKFZ (OpenStack)	14.5	10.2	8.5		
EMBL-EBI	12.6	3.3			
ETRI	2.1	5.8			
iDASH		4.8			
OICR	1.8	5.6	1.0		
PDC	11.8	4.2			
Sanger		7.0	3.0		
Seven Bridges				23.1	
UCSC	30.6	13.0		68.2	
UTokyo	10.9	11.9			

647

648 Supplementary Table 2. Data distribution as of May 2017. While ETRI GNOS and CGHub
 649 served as data centres during the project, they have since been retired. Variant calls include
 650 those from individual variant calling pipelines and the final consensus callsets. Long-term
 651 repositories are denoted by asterisk (*) and will increase their data holdings over time while
 652 GNOS servers are gradually being retired. Latest information can be found at
 653 <https://dcc.icgc.org/repositories>
 654

Data Repository	ICGC Data			TCGA Data		
	% WG Alignments (534 TB)	% RNA-Seq Alignments (13 TB)	% Variant calls (520 GB)	% WG Alignments (240 TB)	% RNA-Seq Alignments (14 TB)	% Variant calls (228 GB)
BSC GNOS	100.0	30.0	0.3			
DKFZ GNOS	25.0		62.9			
EMBL-EBI GNOS	100.0	59.3	98.6			
UTokyo GNOS	54.6	17.1	1.6			
UChicago-ICGC GNOS	16.8	40.3	28.7			
UChicago-TCGA GNOS				100.0	100.0	100.0
EGA*	97.8					
Collaboratory*	100.0	100.0	100.0			
AWS*	76.7	80.1	75.1			
Bionimbus PDC*				100.0	100.0	0.2

655

656 The following set of tables show how costs are calculated for Figure 5 which compares the
657 costs and accuracies of running the different combination of variant calling pipelines.

658
659 **Supplementary Table 3a.** The average run time for each workflow was rounded up to the
660 nearest hour to reflect how AWS charges for EC2 instances that run for part of an hour. The
661 size of the output files are noted as they contribute to either egress or storage costs.

Workflow	Average wall clock run time (hours)	Size of output files (GB)	AWS EC2 Instances Used
BWA-Mem	140	134	m1.xlarge
Sanger	53	2	r3.8xlarge
DKFZ/EMBL	41	5	r3.8xlarge
Broad	89	35	r3.8xlarge
OxoG	4	1.5	m2.4xlarge

662

663

664 **Supplementary Table 3b.** The project utilized EC2 spot instances in US East (N. Virginia), US
665 West (Oregon), EU (Ireland) regions. Because spot pricing fluctuates, users should consult
666 real-time information. The average spot pricing listed here was based on our own usage
667 throughout the project.

AWS EC2 Instances	vCPU	Mem (GiB)	Storage (GB)	Average spot pricing
m1.xlarge	4	15	4 x 420	\$0.0426
r3.8xlarge	32	244	2 x 320	\$0.3382
m2.4xlarge	8	68.4	2 x 840	\$0.0834

668

669

670 **Supplementary Table 3c.** Cost calculations are based on the above spot pricing and an egress
671 cost of \$0.09 per GB. The analysis time is made up of 3 steps: (1) running the BWA-Mem
672 workflow on two separate instances to align simultaneously one tumor and one normal
673 specimen; (2) running the variant calling workflows simultaneously with the longest running
674 workflow dictating the run time of this step; (3) running the OxoG workflow after all variant
675 calling workflows are completed. If analyzing 100 donors with all 3 variant calling pipelines, the
676 analysis will involve running a fleet of 200, 300 and 100 EC2 instances, respectively in the 3
677 steps. We have no other significant storage cost as the reference files amount to ~35GB
678 costing under \$1/month in S3. An alternative to transferring the data out is to store the 312 GB
679 of data for each donor in S3 for under \$8/month.

680

Variant Calling Pipelines	Total Cost	Compute Cost	Egress Cost	Analysis Time (days)	Median Sensitivity, Precision, F1
All 3 pipelines	102.19	7.15	28.04	9.7	0.9047 +/- 0.03145 0.9348 +/- 0.03785 0.9151 +/- 0.02820
Sanger only	54.63	30.19	24.44	8.2	0.8032 +/- 0.06515 0.9550 +/- 0.03855 0.8629 +/- 0.04795
DKFZ/EMBL only	50.84	26.13	24.71	7.7	0.7565 +/- 0.0544 0.9352 +/- 0.0365 0.8313 +/- 0.05125
Broad only	69.77	42.36	27.41	9.7	0.9095 +/- 0.01955 0.8386 +/- 0.06335 0.8687 +/- 0.04085
Sanger & DKFZ/EMBL	68.94	44.05	24.89	8.2	<u>Union</u> 0.8454 +/- 0.0572 0.9032 +/- 0.04405 0.8669 +/- 0.0509 <u>Intersect</u> 0.7228 +/- 0.05385 0.9954 +/- 0.00980 0.8216 +/- 0.04390
Sanger & Broad	87.88	60.29	27.59	9.7	<u>Union</u> 0.9374 +/- 0.01935 0.8183 +/- 0.06395 0.8653 +/- 0.04220 <u>Intersect</u> 0.7856 +/- 0.0566 0.9913 +/- 0.0111 0.8632 +/- 0.03755
DKFZ/EMBL & Broad	84.09	56.23	27.86	9.7	<u>Union</u> 0.9339 +/- 0.01955 0.801 +/- 0.06505 0.8576 +/- 0.0429 <u>Intersect</u> 0.7384 +/- 0.05865 0.9939 +/- 0.0186 0.8315 +/- 0.0456

681

682 Supplementary Table 4. DOIs for PCAWG core analysis workflows
 683

Workflow/Tool	Dockstore	Latest DOI	Version	Github
pcawg-bwa-mem-workflow	https://dockstore.org/containers/quay.io/pancancer/pcawg-bwa-mem-workflow	https://doi.org/10.5281/zenodo.192377	2.6.8_1.2	https://github.com/ICGC-TCGA-PanCancer/Seqware-BWA-Workflow
pcawg-dkfst-workflow	https://dockstore.org/containers/quay.io/pancancer/pcawg-dkfst-workflow	https://doi.org/10.5281/zenodo.192376	2.0.1_cwl1.0	https://github.com/ICGC-TCGA-PanCancer/DEWrapperWorkflow
pcawg-sanger-cgp-workflow	https://dockstore.org/containers/quay.io/pancancer/pcawg-sanger-cgp-workflow	https://doi.org/10.5281/zenodo.192162	2.0.3	https://github.com/ICGC-TCGA-PanCancer/CGP-Somatic-Docker
pcawg_delly_workflow	https://dockstore.org/containers/quay.io/pancancer/pcawg_delly_workflow	https://doi.org/10.5281/zenodo.192166	2.0.1-cwl1.0	https://github.com/ICGC-TCGA-PanCancer/DEWrapperWorkflow
broad				
oxog				

684