# Large-Scale Uniform Analysis of Cancer Whole Genomes in Multiple Computing Environments — Source link ⧉

Christina K. Yung, Brian O'Connor, Sergei Yakneen, Junjun Zhang ...+120 more authors

**Institutions:** Ontario Institute for Cancer Research, Oregon Health & Science University, German Cancer Research Center, University of Tokyo ...+19 more institutions

Related papers:

- Pan-cancer analysis of whole genomes

- The cancer genome atlas pan-cancer analysis project

- International Cancer Genome Consortium Data Portal—a one-stop shop for cancer genomics data

- Signatures of mutational processes in human cancer

- Mutational heterogeneity in cancer and the search for new cancer-associated genes

# Large-Scale Uniform Analysis of Cancer Whole Genomes in Multiple Computing Environments

Christina K. Yung[1,*], Brian D. O'Connor[1,2,*], Sergei Yakneen[1,3,*], Junjun Zhang[1,*], Kyle Ellrott[4], Kortine Kleinheinz[5,6], Naoki Miyoshi[7], Keiran M. Raine[8], Romina Royo[9], Gordon B. Saksena[10], Matthias Schlesner[5], Solomon I. Shorser[1], Miguel Vazquez[11], Joachim Weischenfeldt[3,12], Denis Yuen[1], Adam P. Butler[8], Brandi N. Davis-Dusenbery[13], Roland Eils[14,6], Vincent Ferretti[1], Robert L. Grossman[15], Olivier Harismendy[16,17], Youngwook Kim[18], Hidewaki Nakagawa[19], Steven J. Newhouse[20], David Torrents[9,21], Lincoln D. Stein[1,22,‡] on behalf of the PCAWG Technical Working Group[23] and the PCAWG Network

*These authors contributed equally to this work.*

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

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

**Abstract**

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

**Introduction**

The International Cancer Genome Consortium (ICGC)/The Cancer Genome Atlas (TCGA) Pan-Cancer Analysis of Whole Genomes (PCAWG) study has characterized the pattern of mutations in over 2,800 cancer whole genomes. Extending TCGA Pan-Cancer analysis project, which focused on molecular aberrations in protein coding regions only[1], PCAWG undertook the study of whole genomes, allowing for the discovery of driver mutations in cis-regulatory sites and non-coding RNAs, examination of the patterns of large-scale structural rearrangements, identification of signatures of exposure, and elucidation of interactions between somatic mutations and germline polymorphisms.

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

2

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 donors[2]. 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 types[3]. 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 approach[4].

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 redistributed[5].  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    **<u>Phase 1: Data Marshalling and Upload</u>**

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

5

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 protocol[6]. 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

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

We began the process of sequence alignment about two months after the uploading process had begun. Both tumor and matched normal reads were subjected to uniform sequence alignment using BWA-MEM[7] on top of a common GRCh37-based reference genome that was enhanced with decoy sequences, viral sequences, and the revised Cambridge reference genome for the mitochondria.

Efforts by the project QC group demonstrated that employing multiple variant callers in ensemble fashion improved calling sensitivity[3], thus the aligned tumor/normal pairs were subjected to somatic variant calling using three "best practice" software pipelines. These pipelines were developed by the Sanger Institute[8-11]; jointly by DKFZ[12] and the European Molecular Biology Laboratory (EMBL)[13]; and the Broad Institute[14] with contribution from MD Anderson Cancer Center-Baylor College of Medicine[15]. Each pipeline represents the best practices from the authoring organizations and include the current versions of each institute's flagship tools. Each pipeline consists of multiple software tools for calling of single and multiple nucleotide variants (SNVs and MNVs), small insertions/deletions (indels), structural variants (SVs) and somatic copy number alterations (SCNAs). The minimum compute requirements, median runtime and the analytical algorithms for each pipeline are shown in Table 2.

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

7

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 algorithm[16], 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 algorithm[17]. 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 paper[4].  The algorithm for consensus SVs is described in the

194    PCAWG-6 marker paper[18].  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 algorithms[19].

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   paper[10].

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 engine[20] for bundling software and executing workflows,

10

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 Dockstore[21]

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 Language[22] (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

13

309     Data Distribution / Data Portal

310     While GNOS was used for the core pipelines, Synapse[23] 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     specimens[24]. 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).

16

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**

17

18

434

## **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. & Beroukhim, R. VariantBam: filtering and profiling

471    of next-generational 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. *F1000Res* **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

### Additional Members of the PCAWG Technical Working Group

486    Javier  Bartolomé Rodriguez[1], Keith A. Boroevich[2], Rich  Boyce[3], Angela N. Brooks[4], Alex

487    Buchanan[5], Ivo Buchhalter[6,7], Niall J. Byrne[8], Andy  Cafferkey[9], Peter J. Campbell[10], Zhaohong

488    Chen[11], Sunghoon  Cho[12], Wan  Choi[13], Peter  Clapham[14], Francisco M. De La Vega[15,16], Jonas

489    Demeulemeester[17,18], Michelle T. Dow[19], Lewis J. Dursi[8,20], Juergen  Eils[21], Claudiu  Farcas[22],

490    Francesco Favero[23], Nodirjon Fayzullaev[8], Paul  Flicek[3], Nuno A. Fonseca[3], Josep L.l. Gelpi[24,25],

491    Gad Getz[26,27], Bob Gibson[8], Michael C. Heinold[7,6], Julian M. Hess[26], Oliver Hofmann[28], Jongwhi

492    H. Hong[29], Thomas J. Hudson[30,31], Daniel  Huebschmann[6,7], Barbara  Hutter[32,33], Carolyn M.

493    Hutter[34], Seiya Imoto[35], Sinisa Ivkovic[36], Seung-Hyup Jeon[13], Wei Jiao[8], Jongsun Jung[37], Rolf

494    Kabbe[6], Andre Kahles[38,39], Jules Kerssemakers[40], Hyunghwan Kim[13], Hyung-Lae Kim[41,42],

495    Jihoon Kim[11], Jan O. Korbel[43,3], Michael Koscher[40], Antonios Koures[11], Milena Kovacevic[36],

496    Chris Lawerenz[6], Ignaty Leshchiner[26], Dimitri G. Livitz[26], George L. Mihaiescu[8], Sanja

497    Mijalkovic[36], Ana Mijalkovic Lazic[36], Satoru Miyano[44], Hardeep K. Nahal[8], Mia Nastic[36],

498    Jonathan Nicholson[14], David Ocana[3], Kazuhiro Ohi[44], Lucila Ohno-Machado[22], Larsson

499    Omberg[45], B.F. Francis Ouellette[8,46], Nagarajan Paramasivam[6,47], Marc D. Perry[8], Todd D. Pihl[48],

500    Manuel Prinz[6], Montserrat Puiggròs[24], Petar Radovic[36], Esther Rheinbay[26,49], Mara W.

501    Rosenberg[26,49], Charles Short[3], Heidi J. Sofia[50], Jonathan Spring[51], Adam J. Struck[5], Grace

502    Tiao[26], Nebojsa Tijanic[36], Peter Van Loo[17,18], David Vicente[1], Jeremiah A. Wala[26,52], Zhining

503    Wang[53], Johannes Werner[6], Ashley Williams[11], Youngchoon Woo[13], Adam J. Wright[8], Qian
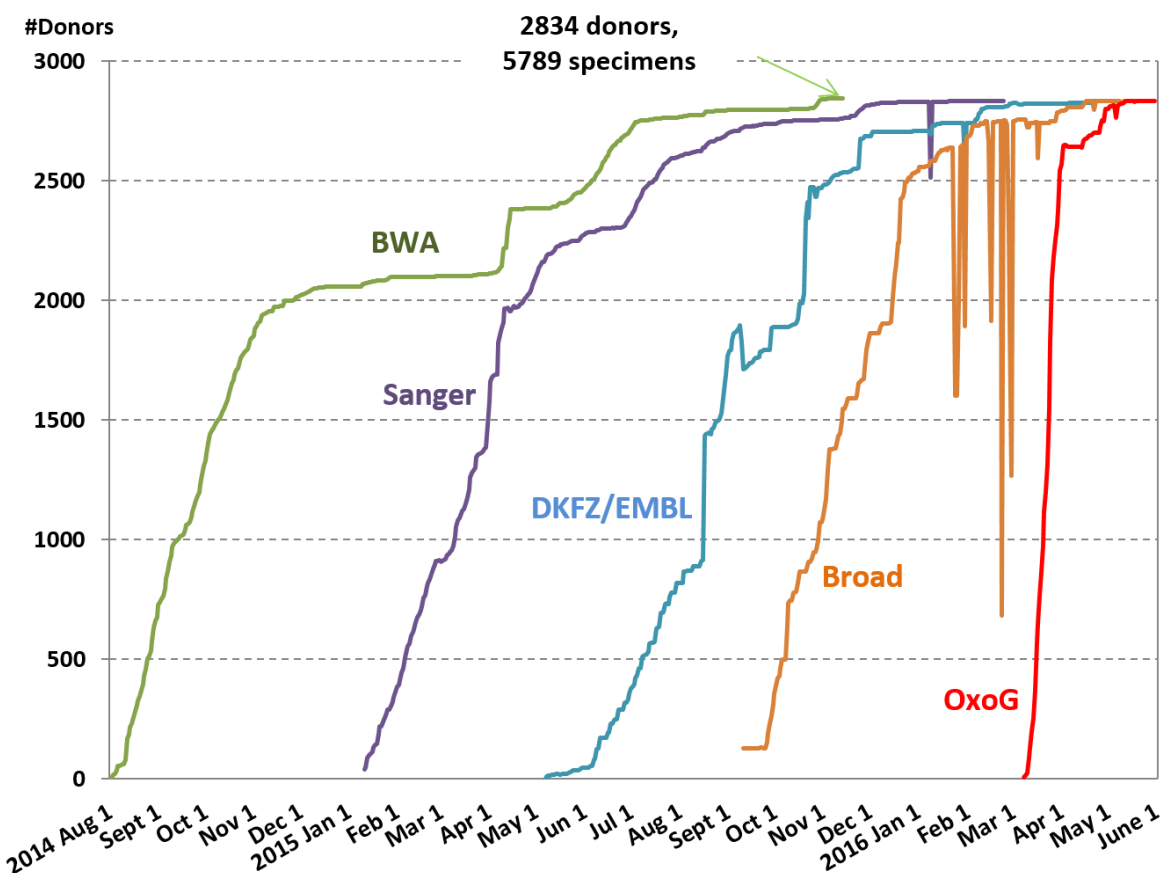
504    Xiang[8]

505

506    [1]Department of Operations, Barcelona Supercomputing Center, Barcelona, Catalunya, 8034, Spain. [2]Laboratory for

507    Medical Science Mathematics, RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa, 230-0045,

508    Japan. [3]European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton, Cambridge, CB10

509    1SD, United Kingdom. [4]Biomolecular Engineering, University of California Santa Cruz, Santa Cruz, California,

510    95065, USA. [5]Department of Computational Biology, Oregon Health and Science University, Portland, Oregon,

511    97239, USA. [6]Division of Theoretical Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg,

512    Baden-Württemberg, 69120, Germany. [7]Department for Bioinformatics and Functional Genomics, Institute for

513    Pharmacy and Molecular Biotechnology and BioQuant, Heidelberg University, Heidelberg, Baden-Württemberg,

514    69120, Germany. [8]Informatics and Biocomputing Program, Ontario Institute for Cancer Research, Toronto, Ontario,

515    M5G 0A3, Canada. [9]Technical Services Cluster, European Molecular Biology Laboratory, European Bioinformatics

516    Institute, Hinxton, Cambridge, CB10 1SD, United Kingdom. [10]Cancer Genome Project, Wellcome Trust Sanger

517    Institute, Hinxton, Cambridgeshire, CB10 1SA, United Kingdom [11]Department of Medicine, University of

518    California San Diego, San Diego, California, 92093, USA. [12]PDXen Biosystems Inc., Seoul, 4900, South Korea.

519    [13]Electronics and Telecommunications Research Institute, Daejon, 34129, South Korea. [14]Informatics Support

520    Group, Wellcome Trust Sanger Institute, Hinxton, Cambridgeshire, CB10 1SA, United Kingdom. [15]Department of

521    Biomedical Data Science, Stanford University School of Medicine, Stanford, California, 94305, USA. [16]Annai

522    Systems, Inc., Carlsbad, California, 92011, USA. [17]The Francis Crick Institute, London, NW1 1AT, United

523    Kingdom. [18]Department of Human Genetics, University of Leuven, B-3000 Leuven, Belgium [19]Biomedical

524    Informatics, University of California San Diego, San Diego, California, 92093, USA. [20]The Centre for

525    Computational Medicine, The Hospital for Sick Children, Toronto, Ontario, M5G 0A4, Canada. [21]Theoretical

526    Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg, Baden-Württemberg, 69120, Germany.

527    [22]Health System Department of Biomedical Informatics, University of California San Diego, La Jolla, California,

528    92093, USA. [23]BRIC/Finsen Laboratory, Rigshospitalet, Copenhagen, 2200, Denmark. [24]Department of Life

529    Sciences, Barcelona Supercomputing Center, Barcelona, Catalunya, 8034, Spain. [25]Department of Biochemistry and

530    Molecular Biomedicine, University of Barcelona, Barcelona, Catalunya, 8028, Spain. [26]Cancer Program, Broad

531    Institute of MIT and Harvard, Cambridge, Massachusetts, 02142, USA. [27]Cancer Center and Department of

532    Pathology, Massachusetts General Hospital, Boston, Massachusetts, 02114, USA. [28]Center for Cancer Research,

533    University of Melbourne, Melbourne, VIC 3001, Australia. [29]Genome Data Integration Center, Syntekabio Inc.,

534    Daejon, 34025, South Korea. [30]Genomics Program, Ontario Institute for Cancer Research, Toronto, Ontario, M5G

535    0A3, Canada. [31]Oncology Discovery and Early Development, AbbVie, Redwood City, California, 94063, USA.

536    [32]Division of Applied Bioinformatics, German Cancer Research Center (DKFZ), Heidelberg, Baden-Württemberg,

537    69120, Germany. [33]Division of Applied Bioinformatics, National Center for Tumor Diseases, Heidelberg, Baden-

538    Württemberg, 69120, Germany. [34]Division of Genomic Medicine, National Human Genome Research Institute,

22

539    Bethesda, Maryland, 20852, USA. [35]Health Intelligence Center, Institute of Medical Science, University of Tokyo,
540    Tokyo, 108-8639, Japan. [36]Seven Bridges, Cambridge, Massachusetts, 02142, USA. [37]Genome Data Integration
541    Center, Syntekabio Inc., Daejon, 34025, South Korea [38]Department of Computer Science, ETH Zurich, Zurich,
542    Zurich, 8092, Switzerland. [39]Computational Biology Center, Memorial Sloan Kettering Cancer Center, New York,
543    New York, 10065, USA. [40]German Cancer Research Center (DKFZ), Heidelberg, Baden-Württemberg, 69120,
544    Germany. [41]Department of Biochemistry, Ewha Womans University, Seoul, O7985, South Korea. [42]PGM21, Seoul,
545    O7985, South Korea. [43]Genome Biology Unit, European Molecular Biology Laboratory, Heidelberg, Baden-
546    Württemberg, 69120, Germany. [44]Human Genome Center, Institute of Medical Science, University of Tokyo,
547    Tokyo, 108-8639, Japan. [45]Systems Biology, Sage Bionetworks, Seattle, Washington, 98112, USA. [46]Department of
548    Cell and Systems Biology, University of Toronto, Toronto, Ontario, M5S 3G5, Canada. [47]Medical Faculty
549    Heidelberg, Heidelberg University, Heidelberg, Baden-Württemberg, 69120, Germany. [48]CSRA Incorporated,
550    Fairfax, Virginia, 22042, USA. [49]Cancer Center, Massachusetts General Hospital, Boston, Massachusetts, 02114,
551    USA. [50]National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland, 20892-
552    9305, USA. [51]Center for Data Intensive Science, University of Chicago, Chicago, Illinois, 60637, USA.
553    [52]Department of Cancer Biology, Dana-Farber Cancer Institute, Boston, Massachusetts, 02115, USA. [53]TCGA
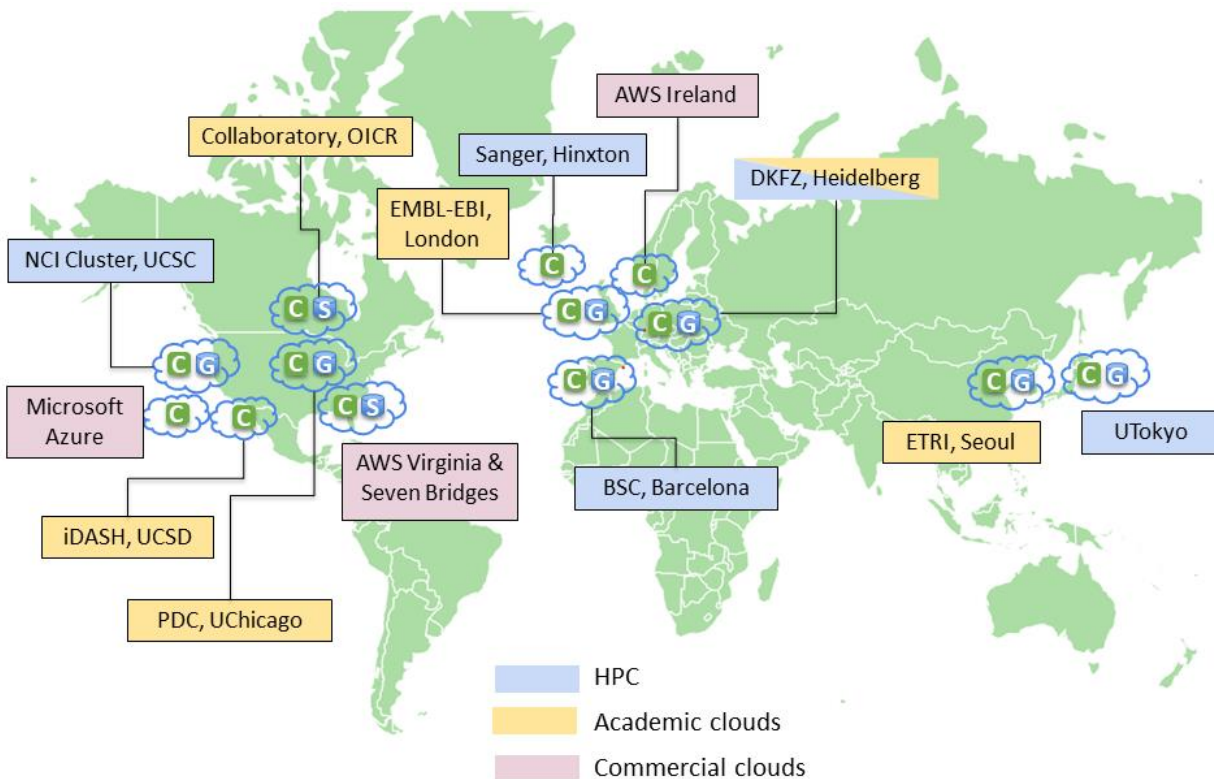554    Program Office, National Cancer Institute, Bethesda, Maryland, 20892, USA.
555

556 **<u>Figures</u>**



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
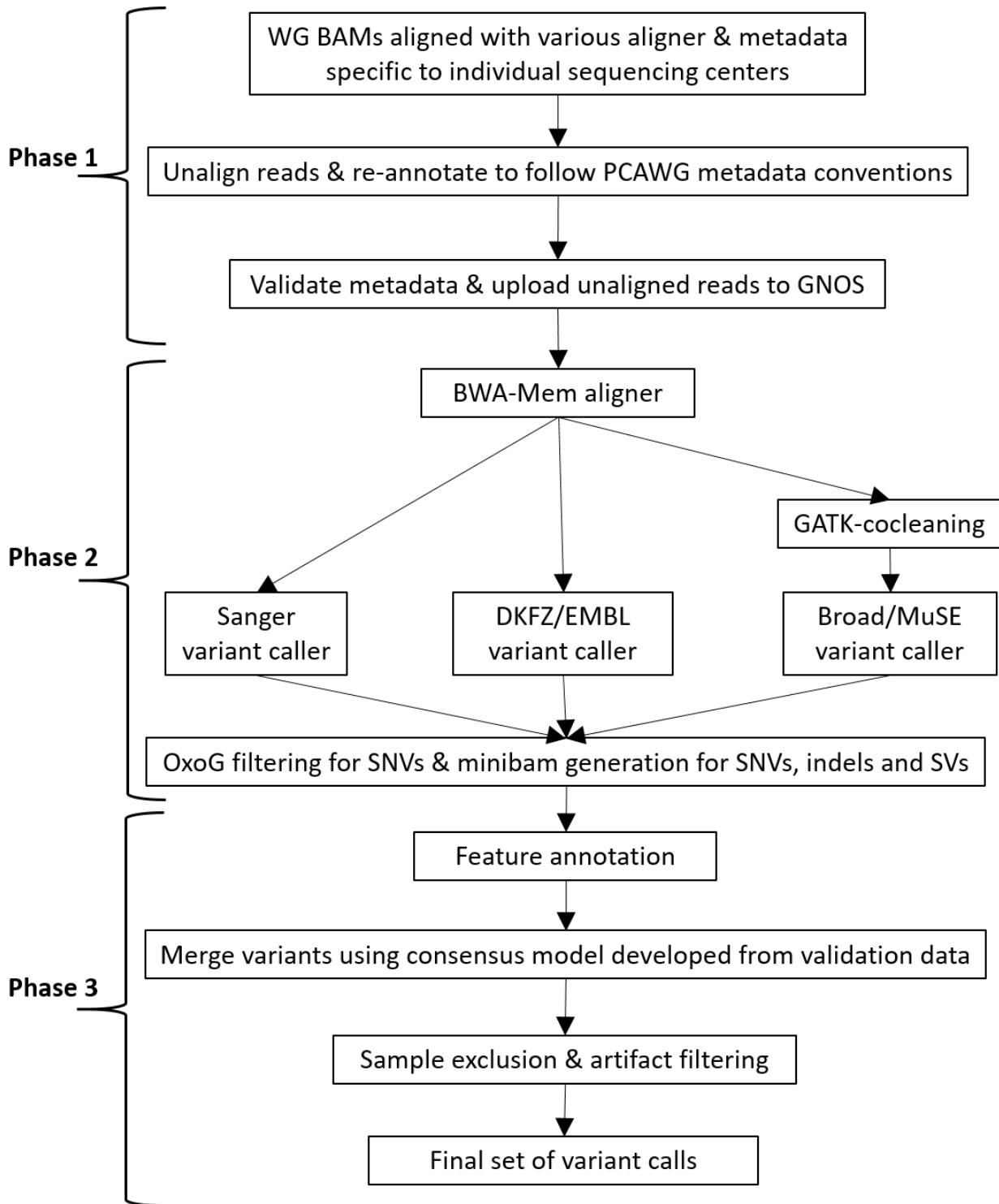568 been completed in under 6 months.

24

569
570 Figure 2: Geographical distribution of compute centers (C), GNOS servers (G), and
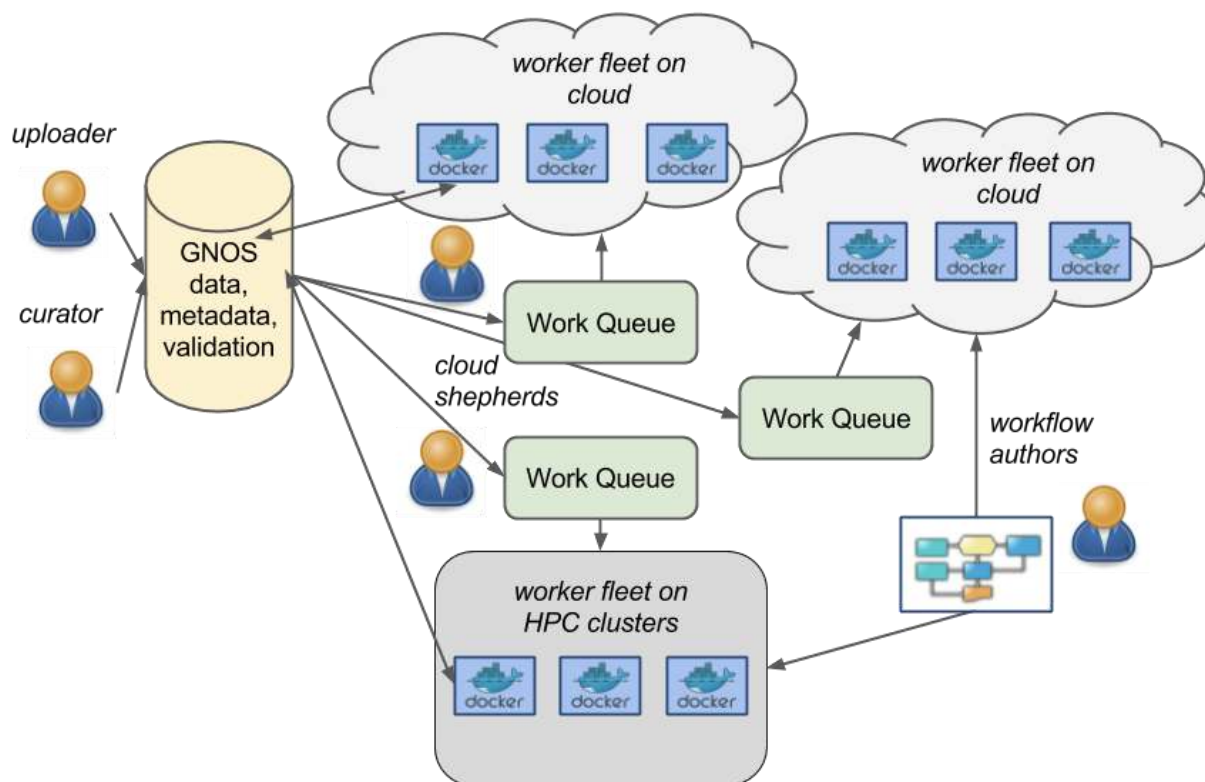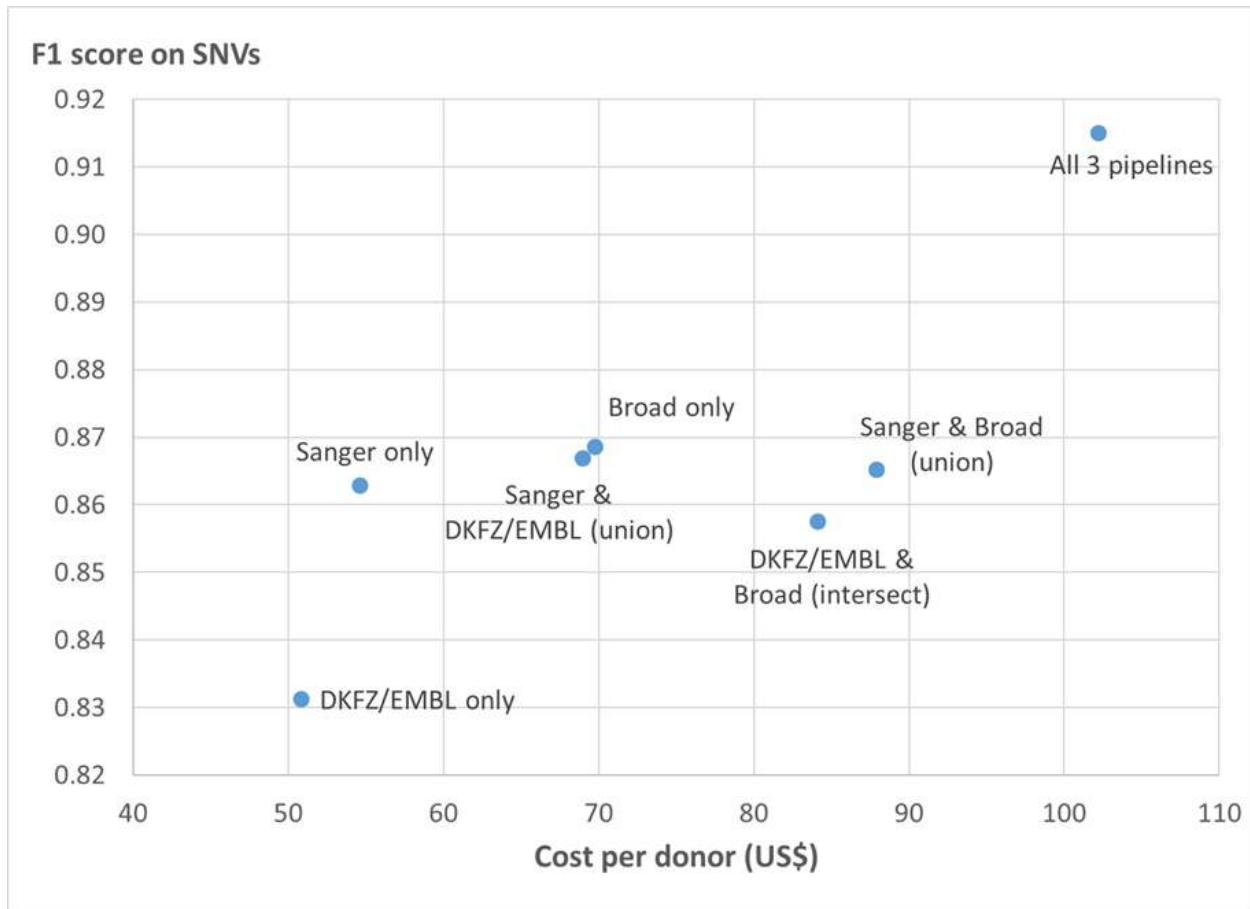571 S3-compatible data storage (S).
572

573
574
575  Figure 3: The uniform analysis of whole genomes involves three broad phases. Phase 1: Data
576  marshalling and upload.  Phase 2: Sequence alignment and variant calling. Phase 3: Variant
577  merging and filtering.  The algorithms for merging SNVs and indels are described in the
578  PCAWG-1 paper, SVs in the PCAWG-6 paper, and CNVs in the PCAWG-11 paper.
579

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

583
584
585    Figure 5: Costs for analyzing a tumor/normal pair through BWA-Mem, different combinations of
586    variant calling pipelines, and OxoG filtering.  The cost is calculated based on AWS instances at
587    average spot pricing we experienced during the project, and includes egress costs to transfer
588    the result files.  PCAWG ran all 3 variant calling pipelines and achieved an F1 score of 0.9151
589    for SNVs.  If running only one or two pipelines, there will be savings in cost but sacrifice in
590    accuracy.  Detailed cost analysis is shown in Suppl Table 3.
591

592 **Tables**
593
594 Table 1. Compute resources. * Shared between environments. ** Transient storage used for
595 local data processing.
596

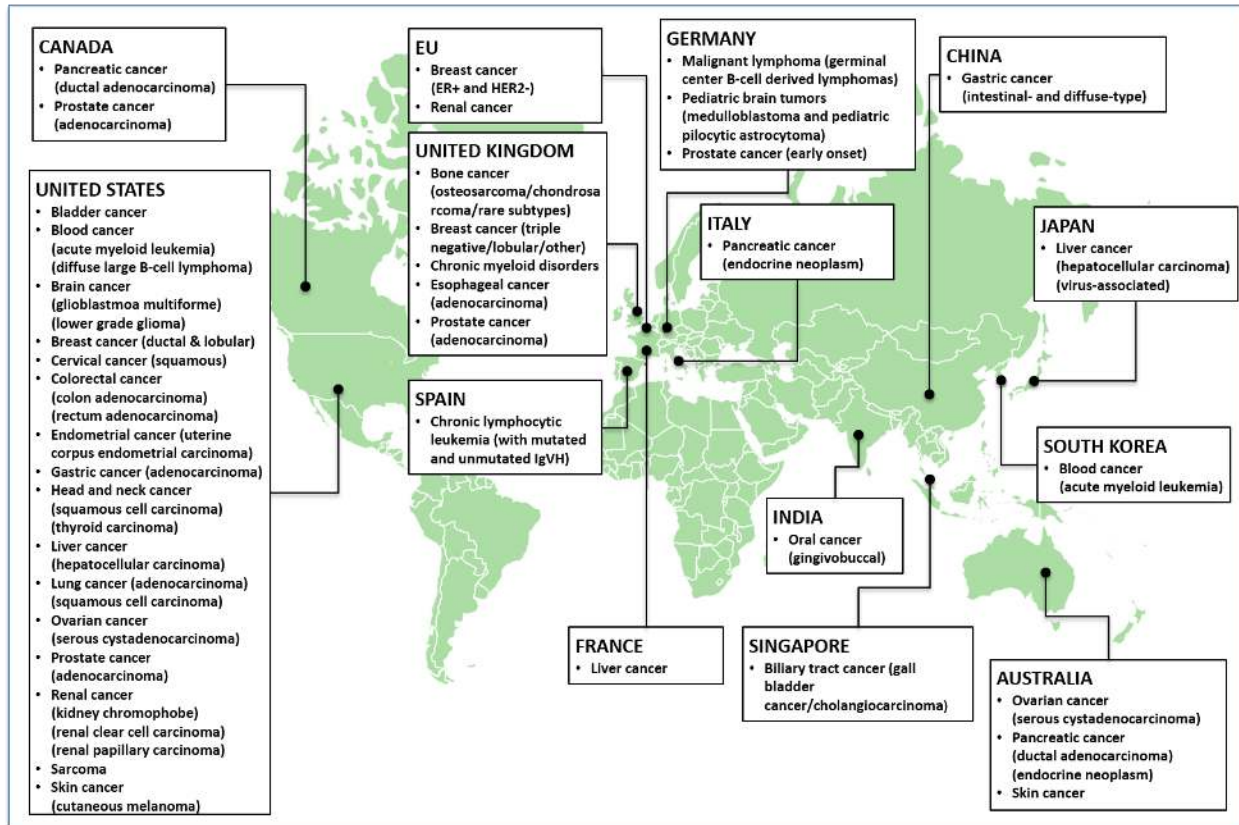|  | Type | Allocated CPU/Cores | Allocated memory | Data Co-location Repository | Local Storage Amount |
|---|---|---|---|---|---|
| AWS | Cloud | variable | variable | Y | 420TB |
| Azure | Cloud | variable | variable | N | - |
| BSC | HPC | 1000 | 7.75TB | Y | 300TB |
| Collaboratory | Cloud | 350 | 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

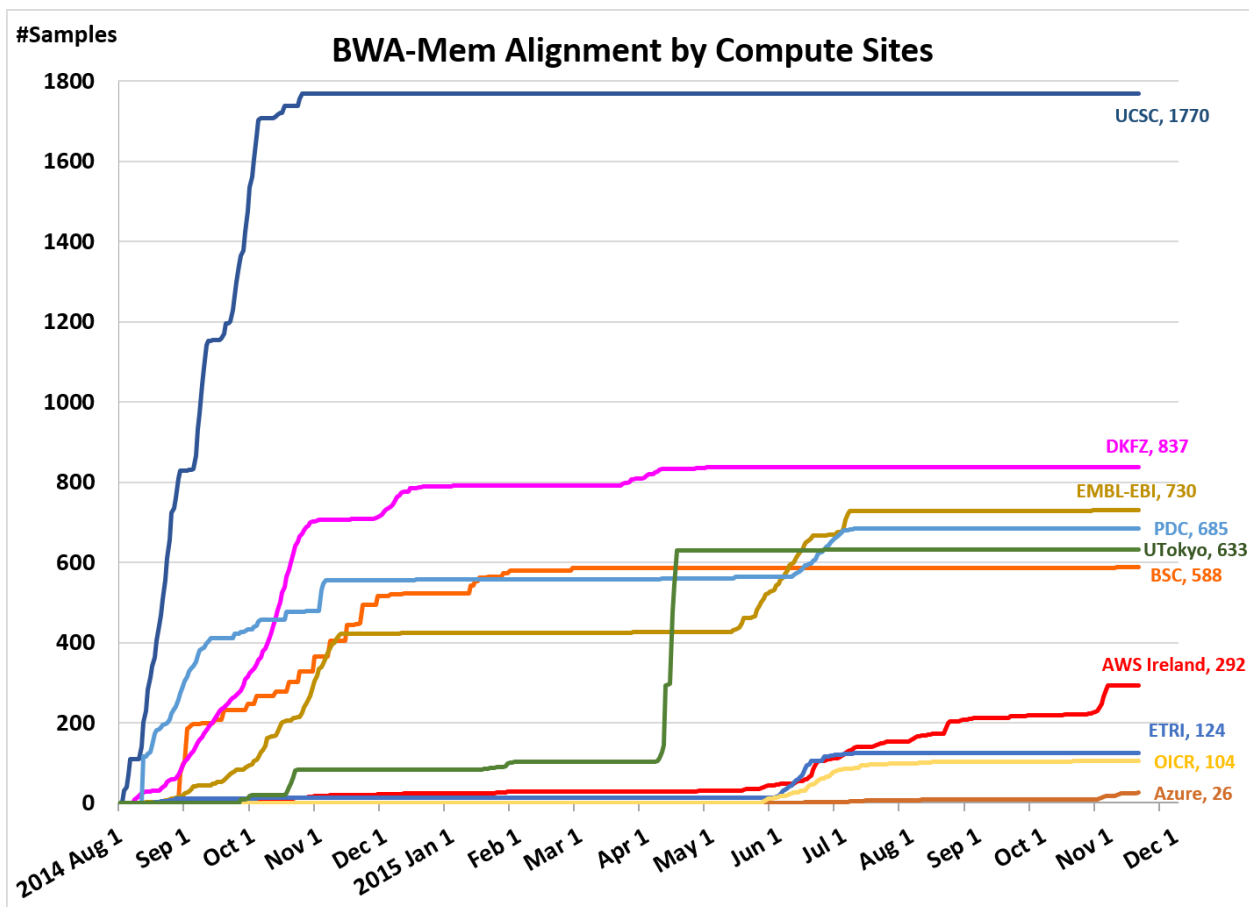| | BWA | Sanger | DKFZ/EMBL | Broad | OxoG |
|---|---|---|---|---|---|
| Analytical components in workflow | BWA-Mem Picard Biobambam samtools | CaVEMan[1] cgpPindel[2] BRASS[3] ascatNgs[4] | dkfz_snv[1] Platypus[2] DELLY[3] ACE-seq[4] | GATK cocleaning MuTect[1] MuSE[1,#] Snowman[2,3] dRanger[3] | 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

30

609
610

Supplementary Figure 1: Whole genomes from 2,834 donors across 39 cancer types were collected from 48 ICGC and TCGA projects in 14 jurisdictions.
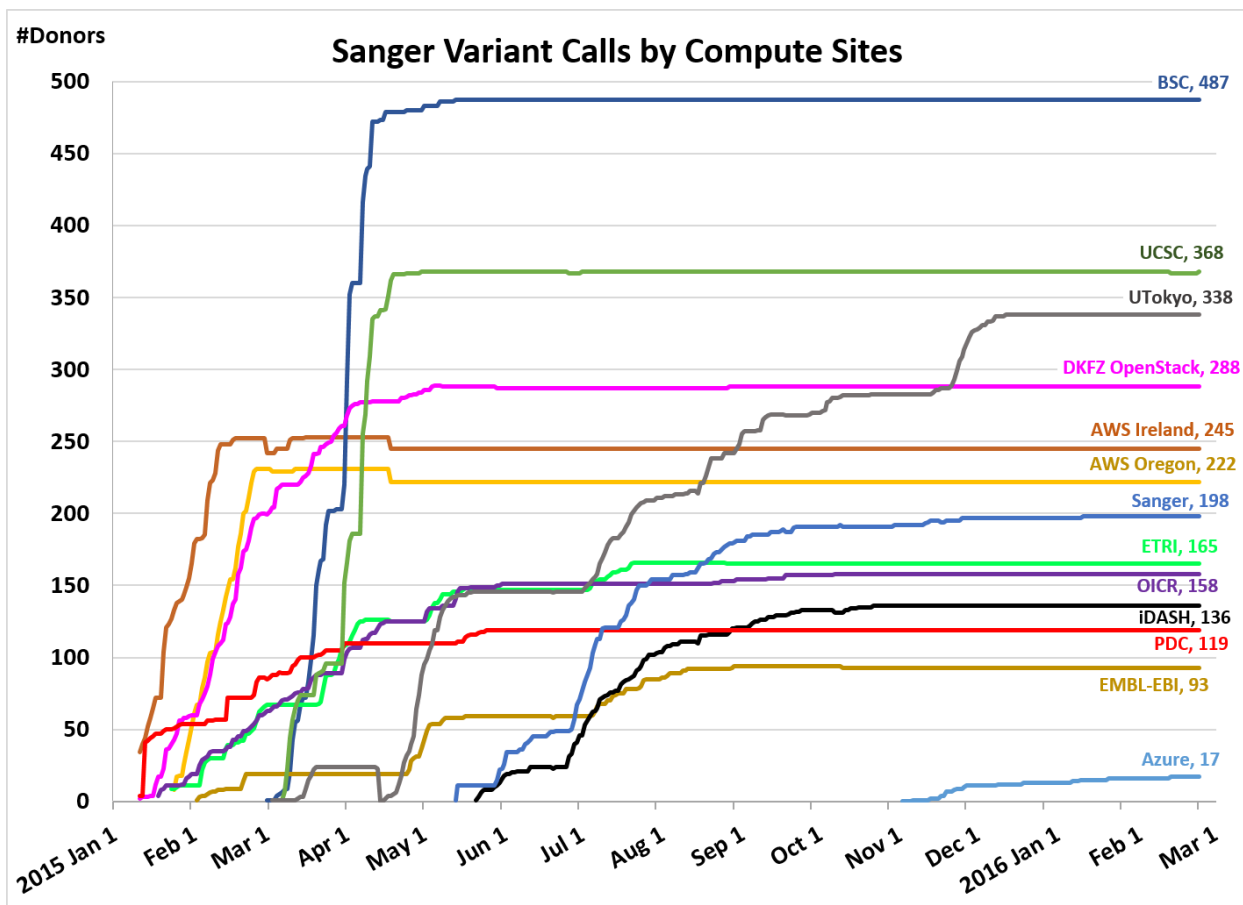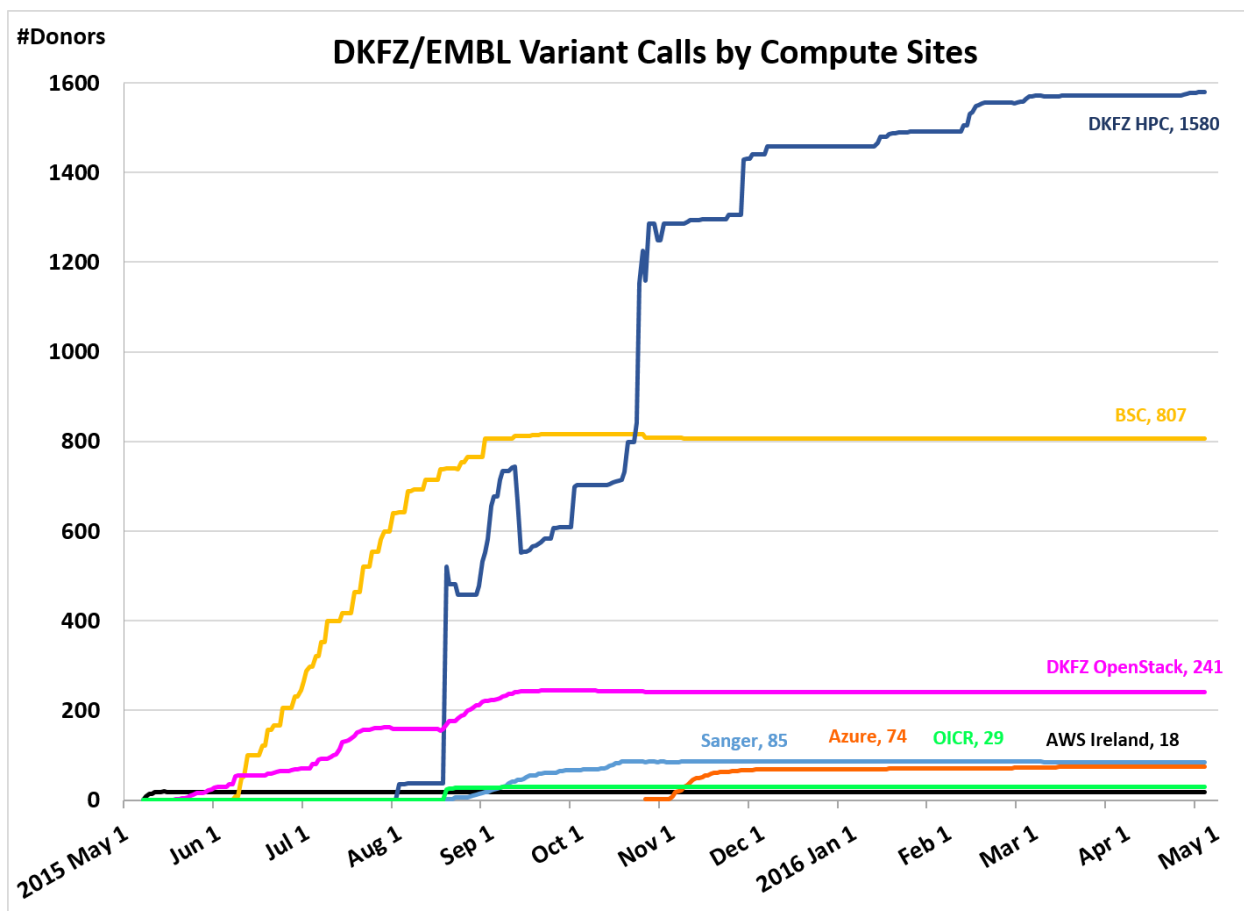
613

614
615    Supplementary Figure 2: Progress of BWA-Mem alignment over time at 7 compute sites.
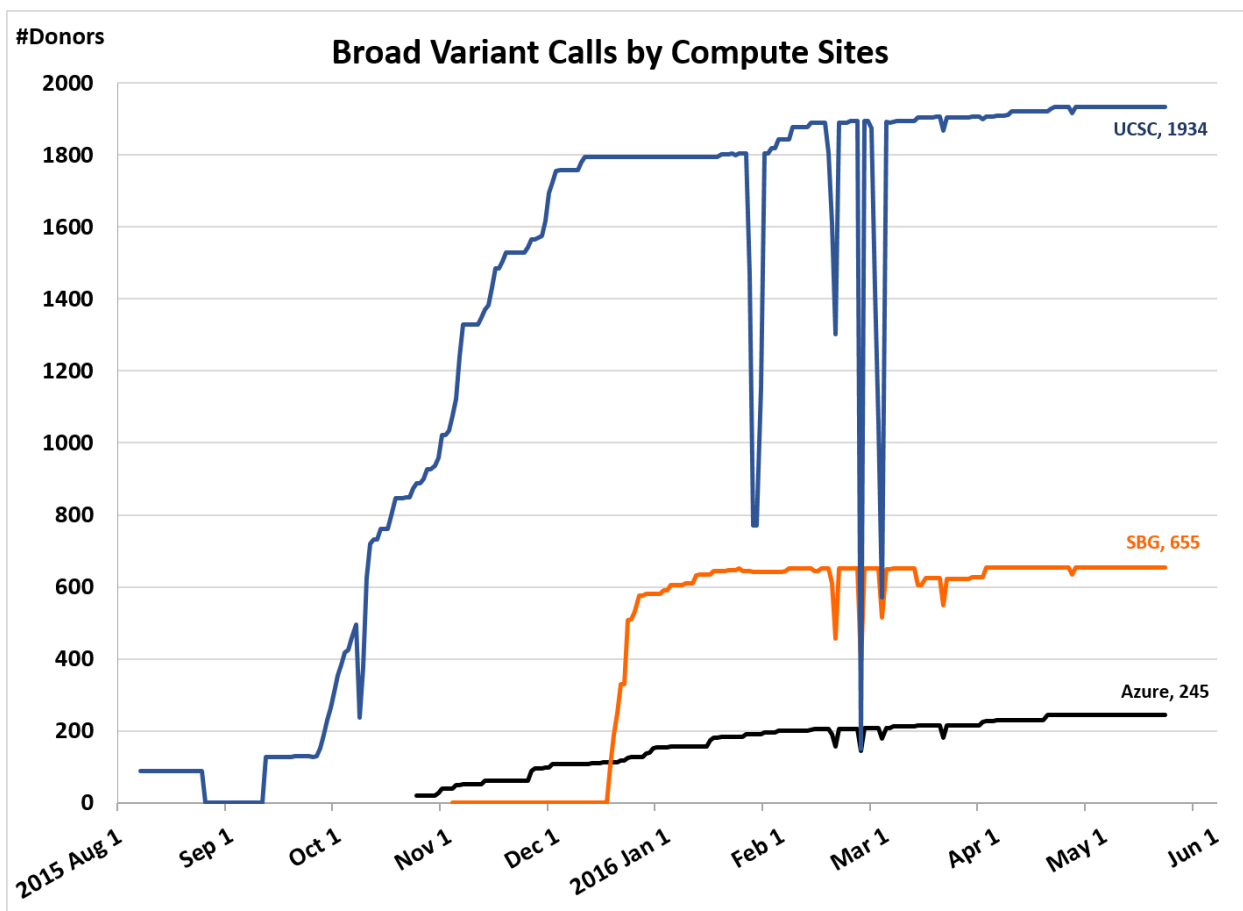616

**Sanger Variant Calls by Compute Sites**

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

33

**#Donors**

**DKFZ/EMBL Variant Calls by Compute Sites**

DKFZ HPC, 1580

BSC, 807

DKFZ OpenStack, 241

Sanger, 85    Azure, 74    OICR, 29    AWS Ireland, 18

621
622    Supplementary Figure 4: Progress of DKFZ/EMBL variant calling workflow over time at 7
623    compute sites.
624

**Broad Variant Calls by Compute Sites**

625

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

627 sites.

628

**#Donors**

**OxoG & Minibam Workflow by Compute Sites**
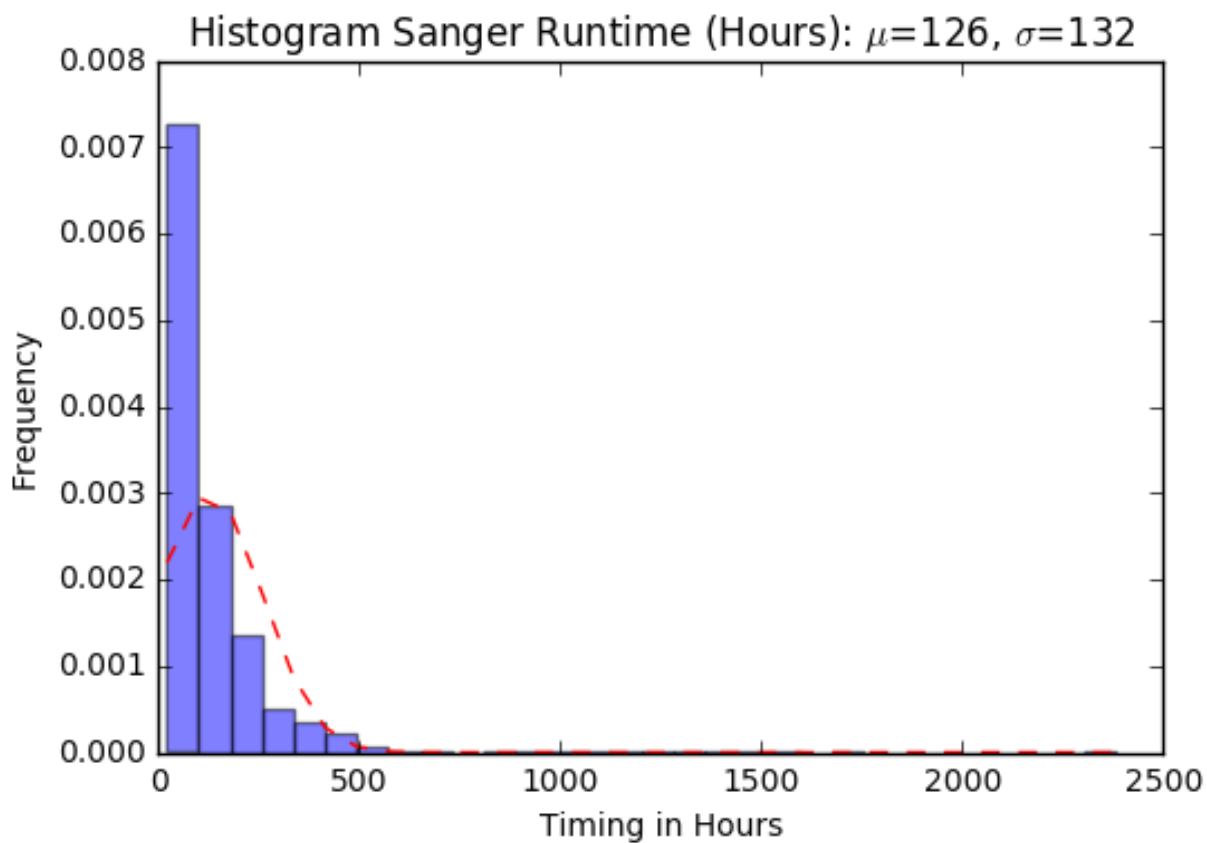
Collaboratory, 1954

AWS, 880

629
630   Supplementary Figure 6: Progress of OxoG and minibam workflow over time at 2 compute sites.
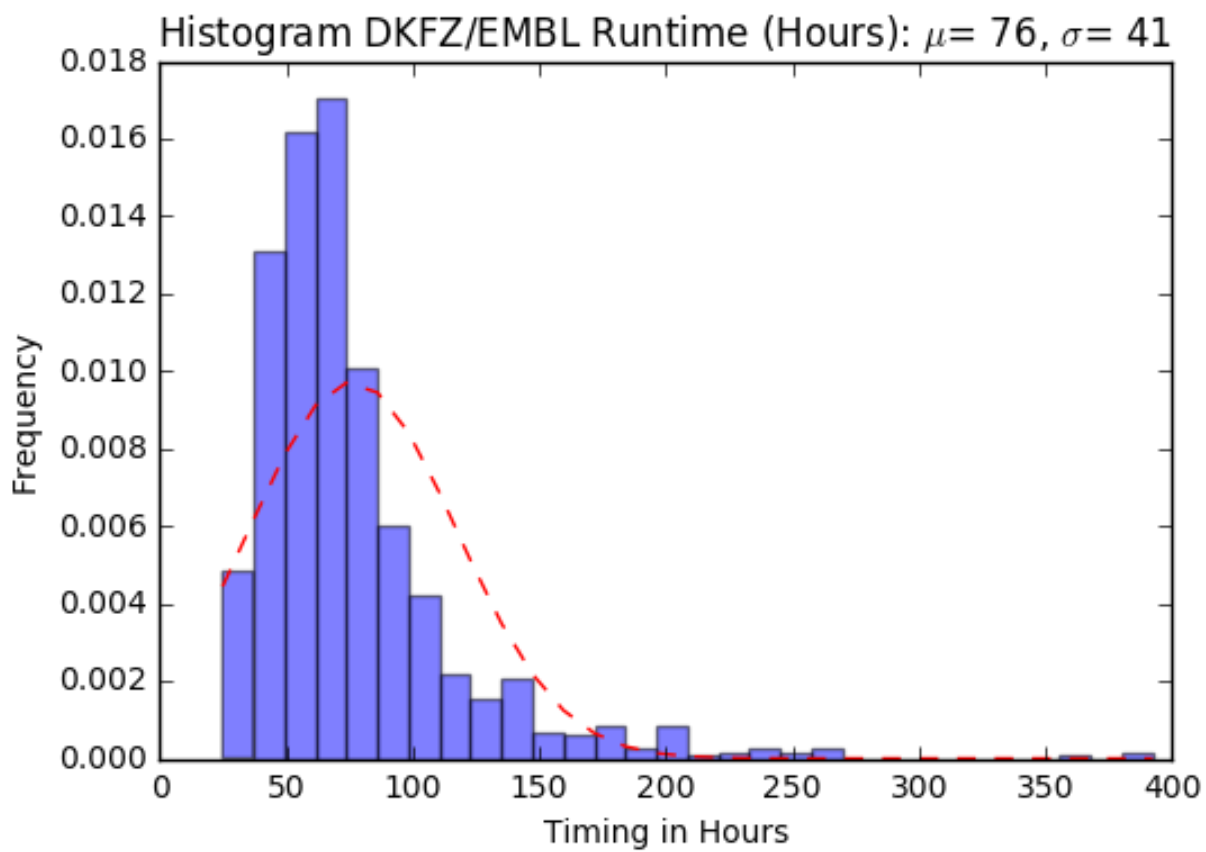631

36

632
633    Supplementary Figure 7: Average runtimes for BWA-Mem alignment workflow
634

635
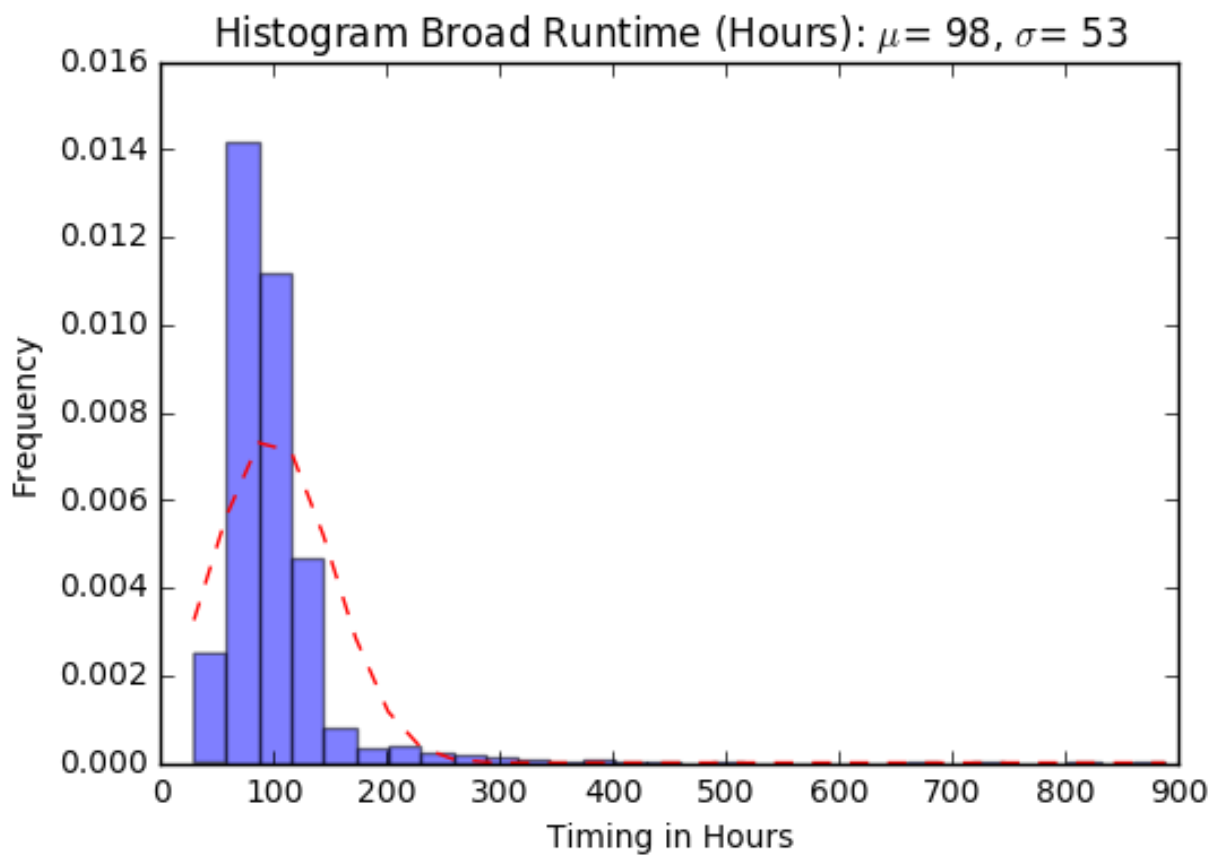636    Supplementary Figure 8: Average runtime for the Sanger somatic variant calling workflow.
637

638
639    Supplementary Figure 9: Average runtime for the DKFZ/EMBL somatic variant calling workflow.
640

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

41

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<br>0.9348 +/- 0.03785<br>0.9151 +/- 0.02820 |
| Sanger only | 54.63 | 30.19 | 24.44 | 8.2 | 0.8032 +/- 0.06515<br>0.9550 +/- 0.03855<br>0.8629 +/- 0.04795 |
| DKFZ/EMBL only | 50.84 | 26.13 | 24.71 | 7.7 | 0.7565 +/- 0.0544<br>0.9352 +/- 0.0365<br>0.8313 +/- 0.05125 |
| Broad only | 69.77 | 42.36 | 27.41 | 9.7 | 0.9095 +/- 0.01955<br>0.8386 +/- 0.06335<br>0.8687 +/- 0.04085 |
| Sanger & DKFZ/EMBL | 68.94 | 44.05 | 24.89 | 8.2 | Union<br>0.8454 +/- 0.0572<br>0.9032 +/- 0.04405<br>0.8669 +/- 0.0509<br>Intersect<br>0.7228 +/- 0.05385<br>0.9954 +/- 0.00980<br>0.8216 +/- 0.04390 |
| Sanger & Broad | 87.88 | 60.29 | 27.59 | 9.7 | Union<br>0.9374 +/- 0.01935<br>0.8183 +/- 0.06395<br>0.8653 +/- 0.04220<br>Intersect<br>0.7856 +/- 0.0566<br>0.9913 +/- 0.0111<br>0.8632 +/- 0.03755 |
| DKFZ/EMBL & Broad | 84.09 | 56.23 | 27.86 | 9.7 | Union<br>0.9339 +/- 0.01955<br>0.801 +/- 0.06505<br>0.8576 +/- 0.0429<br>Intersect<br>0.7384 +/- 0.05865<br>0.9939 +/- 0.0186<br>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-dkfz-workflow | https://dockstore.org/containers/quay.io/pancancer/pcawg-dkfz-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

45