

# 1 **GenomeChronicler: The Personal Genome Project UK Genomic** 2 **Report Generator Pipeline**

3 **José Afonso Guerra-Assunção<sup>1\*</sup>, Lucia Conde<sup>2</sup>, Ismail Moghul<sup>3</sup>, Amy P. Webster<sup>3</sup>, Simone**  
4 **Ecker<sup>3</sup>, Olga Chervova<sup>3</sup>, Christina Chatzipantsiou<sup>4</sup>, Pablo P. Prieto<sup>4</sup>, Stephan Beck<sup>3</sup>, Javier**  
5 **Herrero<sup>2</sup>**

6 <sup>1</sup>Infection and Immunity, University College London, London, United Kingdom

7 <sup>2</sup>Bill Lyons Informatics Centre, UCL Cancer Institute, University College London, London, United  
8 Kingdom

9 <sup>3</sup>Medical Genomics, UCL Cancer Institute, University College London, London, United Kingdom

10 <sup>4</sup>Lifebit, The Bower, 207 Old Street, London, United Kingdom

11 **\* Correspondence:**

12 Corresponding Author

13 a.guerra@ucl.ac.uk

14 **Keywords: Personal Genomics, PGP-UK, Genomic Report, Open Consent, Participant**  
15 **Engagement, Open Source, Cloud Computing.**

16 **Abstract**

17 In recent years, there has been a significant increase in whole genome sequencing data of individual  
18 genomes produced by research projects as well as direct to consumer service providers. While many  
19 of these sources provide their users with an interpretation of the data, there is a lack of free, open  
20 tools for generating similar reports exploring the data in an easy to understand manner.

21 GenomeChronicler was written as part of the Personal Genome Project UK (PGP-UK) to address this  
22 need. PGP-UK provides genomic, transcriptomic, epigenomic and self-reported phenotypic data  
23 under an open-access model with full ethical approval. As a result, the reports generated by  
24 GenomeChronicler are intended for research purposes only and include information relating to  
25 potentially beneficial and potentially harmful variants, but without clinical curation.

26 GenomeChronicler can be used with data from whole genome or whole exome sequencing producing  
27 a genome report containing information on variant statistics, ancestry and known associated  
28 phenotypic traits. Example reports are available from the PGP-UK data page  
29 (personalgenomes.org.uk/data).

30 The objective of this method is to leverage on existing resources to find known phenotypes  
31 associated with the genotypes detected in each sample. The provided trait data is based primarily  
32 upon information available in SNPedia, but also collates data from ClinVar, GETevidence and  
33 gnomAD to provide additional details on potential health implications, presence of genotype in other  
34 PGP participants and population frequency of each genotype.

35 The whole pipeline is self-contained, and runs without internet connection, making it a good choice  
36 for privacy conscious projects that can run GenomeChronicler within their off-line safe-haven

37 environments. GenomeChronicler can be run for one sample at a time, or in parallel making use of  
38 the nextflow workflow manager.

39 The source code is available from GitHub (<https://github.com/PGP-UK/GenomeChronicler>),  
40 container recipes are available for Docker and Singularity, as well as a pre-built container from  
41 SingularityHub (<https://singularity-hub.org/collections/3664>) enabling easy deployment in a variety  
42 of settings. Users without access to computational resources to run GenomeChronicler can access the  
43 software from the Lifebit CloudOS platform (<https://cloudos.lifebit.ai>) enabling the production of  
44 reports and variant calls from raw sequencing data in a scalable fashion.

## 45 1 Introduction

46 The publication of the first draft human genome sequence ('Initial Sequencing and Analysis of the  
47 Human Genome' 2001) brought along the promise of a revolution in how we see ourselves as  
48 individuals and how future medical care should take into account our genetic background. Almost ten  
49 years later, the perspective of widespread personal genomics was still to be achieved (Venter 2010).

50 There is now a wide range of clinical and non-clinical genetic tests that are routinely employed to  
51 detect individuals' carrier status for certain disease genes or particular mutations of clinical  
52 relevance. Many more associations between genotype and phenotype have been highlighted by  
53 research, sometimes with uncertain clinical relevance or simply describing personal traits like eye  
54 color (Pontikos et al. 2017; Kuleshov et al. 2019).

55 Over the past few years we have seen a dramatic reduction of the cost to sequence the full human  
56 genome. This reduction in cost enables many more projects to start using whole genome sequencing  
57 (WGS) approaches, as well as the marked rise in the number of personal genomes being sequenced.

58 Personal genomics is very much part of the public consciousness as can be seen by the rampant rise  
59 of direct to consumer (DTC) genomic analysis offerings on the market. In this context it is  
60 unsurprising that the analysis of one's own genome provides a valuable educational opportunity  
61 (Salari et al. 2013; Linderman et al. 2018) as well as increases participant engagement as part of  
62 biomedical trials (Sanderson et al. 2016).

63 The personal genome project is one of the initiatives enabled by the increased popularity of whole  
64 genome sequencing and its lowering costs. The global PGP network currently consists of 5 projects  
65 spread around the world, managed independently but joined by a common goal of providing open  
66 access data containing genomic, environmental and trait information  
67 (<https://www.personalgenomes.org/>).

68 Data analysis within PGP-UK poses interesting ethical challenges, as all the data and genome reports  
69 are intended to become freely and openly available on the World Wide Web. However, until  
70 completion and approval of the reports, the data must be treated as confidential private information.  
71 Prior to enrollment, all participants are well informed and tested for their understanding of the  
72 potential risks of participating in a project of this nature. Upon receipt of their report, participants  
73 have three options. First, they can trigger the release of their report and data themselves by selecting  
74 the 'release immediately' option in their personal accounts. To date, 67% of participants have  
75 selected this release option. Second, they can withdraw from the study in which case no release  
76 occurs and all data will be deleted. This option has never been selected by any participant. Third, the  
77 participants default to a cool-off period of four weeks to explore their data and reports and to seek all

78 the required clarifications. If neither option one or two are selected by the end of the cool-off period,  
79 the data and reports will be released automatically.

80 There are several resources aimed at users of DTC genetic testing companies on the internet  
81 including Promethease and Genomelink ('Promethease' 2019; 'Genomelink | Upload Raw DNA Data  
82 for Free Analysis On 25 Traits' 2019). There are some other tools with a focus on clinical aspects or  
83 particular diseases (Nakken et al. 2018), as well as academic databases containing genotypes of other  
84 individuals (Greshake et al. 2014), pharmacogenomic information (Klein and Ritchie 2018) or  
85 genotype to phenotype links (Ramos et al. 2014; Pontikos et al. 2017; Kuleshov et al. 2019) that can  
86 be useful for the interpretation of personal genomes. Many of these are linked into resources like  
87 SNPedia (Cariaso and Lennon 2012), allowing a wide range of exploration options for the known  
88 associations of each genotype from multiple perspectives.

89  
90 Surprisingly, we found no pre-existing solution that would allow the annotation and evaluation of  
91 variants on the whole genome level, assessment of ancestry and more fine-grained analysis of  
92 variants that have been previously associated with specific phenotypes. In particular one that could be  
93 run locally ensuring full control of the data before the results are scrutinized and approved.

94 GenomeChronicler represents, to the best of our knowledge, the first pipeline that can be run off-line  
95 or in the cloud, to generate personal genomics reports that are not limited to disease only, from whole  
96 genome or whole exome sequencing data.

97 GenomeChronicler contains a database of positions of interest for ancestry or phenotype. The  
98 genotype at each of these positions is inferred from the user provided data that has been mapped to  
99 the human genome. These genotypes are then compared to local versions of a series of publicly  
100 available resources to infer ancestry and likely phenotypes for each individual participant. These  
101 results are then presented as a PDF document containing hyperlinks where more information about  
102 each variant and phenotype can be found. A visual representation of the pipeline and its underlying  
103 resources is shown in Figure 1.

104 This pipeline will continue to be developed and used to generate genome reports by PGP-UK (Beck  
105 et al. 2018). We envision this project will also be useful to other research endeavors that want to  
106 provide personal genomes information to their participants to increase engagement; e.g. to altruistic  
107 individuals who have obtained their whole genome sequencing data from a DTC or health care  
108 provider and are looking for an ethics-approved framework to share their data;

## 109 **2 Materials and Methods**

### 110 **2.1 Data Preprocessing Requirements**

111 The GenomeChronicler pipeline was designed to run downstream of a standardised germline variant  
112 calling pipeline. GenomeChronicler requires a pre-processed BAM file and optionally, the summary  
113 HTML report produced by the Ensembl Variant Effect Predictor (McLaren et al. 2016).

114 GenomeChronicler can be run with any variant caller provided that the reference dataset is matched  
115 to the reference genome used (the included GenomeChronicler databases currently use GRCh38). It  
116 is also imperative that the BAM or CRAM files used have had their duplicates removed and quality  
117 recalibrated prior to being used for GenomeChronicler.

118 To simplify this entire process and to make the tool more accessible to users who may not know how  
119 to run a germline variant calling pipeline, GenomeChronicler can also be run in a fully automated  
120 mode where the germline variant calling pipeline is also run and the whole process is managed by the  
121 nextflow workflow management system. In this scenario, GenomeChronicler uses the Sarek pipeline  
122 (Garcia et al. 2018) to process raw FASTQ files in a manner that follows the GATK variant calling  
123 best practices guidelines (Van der Auwera et al. 2013). Manual inspection of the initial quality  
124 control steps of Sarek is recommended prior to perusing the final results.

125 The combined version of Sarek + GenomeChronicler written using the nextflow workflow manager  
126 (Di Tommaso et al. 2017) is available both on Github ([https://github.com/PGP-  
127 UK/GenomeChronicler-Sarek-nf](https://github.com/PGP-UK/GenomeChronicler-Sarek-nf)) and on Lifebit CloudOS.

## 128 2.2 Ancestry Inference

129 We infer an individuals' ancestry through a Principal Components Analysis (PCA) which is a widely  
130 used approach for identifying ancestry difference among individuals (Novembre et al. 2008).  
131 For each sample of interest, we merged the genotypes with a reference dataset consisting of  
132 genotypes from the 1000 genomes project samples (The 1000 Genomes Project Consortium 2015),  
133 containing individuals from 26 different worldwide populations and applying PCA on the merged  
134 genotype matrix.

135 Prior to merging data was filtered to keep only unrelated samples. In order to avoid strand issues  
136 when merging the datasets, all ambiguous (A/T and C/G) SNPs were removed, as well as non-  
137 biallelic SNPs, SNPs with >5% of missing data, rare variants (MAF < 0.05) and SNPs out of Hardy-  
138 Weinberg equilibrium (pval < 0.0001). From the remaining SNPs, a subset of unlinked SNPs are  
139 selected by pruning those with  $r^2 > 0.1$  using 100-SNP windows shifted at 5-SNP intervals.  
140 These genotypes are used to run PCA based on the variance-standardized relationship matrix, selecting  
141 20 as the number of PCs to be extracted. We then project the data over the first 3 principal  
142 components to identify clusters of populations and highlight the sample of unknown ancestry on the  
143 resulting plot.

144 Here, we used PLINK (Purcell et al. 2007) to process the genotype data and the R Statistical  
145 Computing platform for plotting the final PCA figures to illustrate the ancestry of each sample. An  
146 example of the distribution of the reference samples on the PCA is show in Figure 2.  
147

## 148 2.3 Linked Databases

### 149 2.3.1 SNPedia

150 SNPedia is a large public repository of manually added as well as automatically mined genotype to  
151 phenotype links sourced from existing literature. SNPedia (Cariaso and Lennon 2012) is the core  
152 resource behind the phenotype tables in GenomeChronicler; it not only provides annotations for  
153 single-gene phenotypes, but also for a few phenotypes involving multiple loci referred to as genosets  
154 in the produced reports.

### 155 2.3.2 ClinVar

156 ClinVar (Landrum and Kattman 2018) is a database hosted by the NCBI that focuses exclusive of  
157 variants related to health that has been running since 2013. In comparison to SNPedia, ClinVar is a  
158 much smaller database that is closely linked to the clinical relevance of each variant. ClinVar is  
159 curated more strictly with clinical review – something that is not available for the other data sources  
160 used by GenomeChronicler.

161 **2.3.3 GETevidence**

162 GETevidence was developed as part of the Personal Genome Project Harvard (Mao et al. 2016) to  
163 showcase the variants present within its participants and to allow manual annotation and  
164 interpretations of the results. For some of the genotypes present, it also contains manual annotations  
165 that have been added by the users or curation team. GETevidence allow individuals to compare their  
166 genotypes against those from other personal genomes available within the Harvard project.

167 **2.3.4 gnomAD**

168 Spanning several human populations, the Genome Aggregation Database (gnomAD) (Karczewski et  
169 al. 2019) aggregates data from multiple sources to produce an atlas of variation across the human  
170 genome. Extensively annotated and now covering most of the latest assembly of the human genome,  
171 these links enable easy access to information such as allele frequencies for the genotype across  
172 different populations around the world, as well as some annotation context for each variant, regarding  
173 potential effect on genes if relevant and how selection forces are constraining the genomic region.

174 **2.4 Database Updates**

175 The underlying databases required to run GenomeChronicler are provided within the package. A set  
176 of scripts to regenerate these SQLite databases is also provided within the source code. When the  
177 databases are generated, a set of positions of the interest is compiled so that when genotyping is  
178 performed only relevant positions are computed to save computational time.

179  
180 SNPedia provides an API to query its records in a systematic way. The other linked databases  
181 provide regular dumps of the whole dataset, enabling easy assessment for which dbSNP rs identifiers  
182 are represented within the full database. The use of rs identifiers and genotypes to link between the  
183 different databases enables an unambiguous way to compare information between different  
184 resources.

185 **2.5 Genotype assessment and reporting**

186 In many scenarios, during normal genomics data processing, only VCF files are produced, which do  
187 not include any information regarding the positions of the genome that match the reference sequence.  
188 These become indistinguishable from positions in the genome where there is no read coverage.

189 To ensure comparable results between runs, the genotype information (gVCF) is computed,  
190 following GATK best practices, during each run of GenomeChronicler. To reduce the computational  
191 burden of computing genotypes, only a subset of genomic positions that we know are meaningful,  
192 from the ancestry and phenotype databases, are computed thus saving computational time and storage  
193 space.

194 **2.6 The Genome Report Template**

195 GenomeChronicler is a multistep modular approach in which the final report is only compiled as the  
196 very last step, integrating data from all previous steps. To give users the possibility of fully  
197 customising the report layout and the amount and content of extra information provided in each  
198 report, GenomeChronicler uses a template file written in the LaTeX typesetting language. This can  
199 be modified to better suit the user, for example: To include project branding and introductory texts to  
200 put the report into perspective; To integrate more analyses from other pipelines ran independently  
201 from GenomeChronicler provided the results are in a format that can be typeset using LaTeX and are

202 present in predefined locations that can be sourced from the template file; To deactivate certain  
203 sections of the report that are not relevant; Or simply to modify the structure of the report produced.

204 **2.7 Output Files**

205 The main output of the GenomeChronicler pipeline is a full report in PDF format, containing  
206 information from all sections of the pipeline that have run as set by the LaTeX template provided  
207 when running the script. Additionally an Excel file containing the genotype phenotype link  
208 information, and all corresponding hyperlinks is also produced, allowing the user to reorder and/or  
209 filter out results as they see fit in a familiar environment. While most intermediate files are  
210 automatically removed at the end of the GenomeChronicler run, the original PDF version of the  
211 ancestry PCA plot is retained, as well as a file containing the sample name within the results  
212 directory to ease automation and a log file containing output produced whilst running  
213 GenomeChronicler.

214 **2.8 Accessing GenomeChronicler**

215 Just like the PGP-UK data, all the code for GenomeChronicler is freely available. To make it easier  
216 to implement, several options are available to remove the need for installing dependencies and  
217 underlying packages, or even the need to own computer hardware capable of handling the processing  
218 of a human genome. The range of options available is detailed below and illustrated in Figure 1.

219 **2.8.1 Running GenomeChronicler Locally**

220 **2.8.1.1 From the available source code**

221 The source code for GenomeChronicler is available on GitHub at <https://github.com/PGP-UK/GenomeChronicler>. The pre-compiled accessory databases are available as links within a setup  
222 script that will help download all the required information.

224 GenomeChronicler has a series of dependencies including LaTeX, R and Perl. The provided  
225 Singularity recipe file can act as a useful list of required packages, in particular for those installing it  
226 on a Debian/Ubuntu based system.

227 **2.8.1.2 Using a pre-compiled container**

228 For those that have access to a machine where the Singularity (Kurtzer, Sochat, and Bauer 2017)  
229 container solution is installed, a container with all dependencies pre-installed and ready to use can be  
230 obtained from SingularityHub (Sochat, Prybol, and Kurtzer 2017). This can be performed by running  
231 the command: `singularity pull shub://PGP-UK/GenomeChronicler`.

232 Once downloaded, the main script (GenomeChronicler\_mainDruid.pl) can be run with the desired  
233 data and options to produce genome reports.

234 **2.8.2 Running GenomeChronicler on Cloud**

235 To enable reproducible, massively parallel, cloud native analyses, GenomeChronicler has also been  
236 implemented as a Nextflow pipeline. The implementation abstracts the installation overhead from the  
237 end user, as all the dependencies are already available via pre-built containers, integrated seamlessly  
238 in the Nextflow pipeline. The source code for this implementation is available on GitHub at  
239 <https://github.com/PGP-UK/GenomeChronicler-nf>, as a standalone nextflow process. To provide an  
240 end-to-end FASTQ to PGP-UK reports pipeline, we also implemented an integration of

241 GenomeChronicler, with a curated and widely used by the bioinformatics community pipeline,  
242 namely Sarek (Garcia et al. 2018; Ewels et al. 2019). This PGP-UK implementation of Sarek is  
243 available on GitHub at <https://github.com/PGP-UK/GenomeChronicler-Sarek-nf>. The  
244 aforementioned pipeline, is available in the collection of curated pipelines on the Lifebit CloudOS  
245 platform (<https://cloudos.lifebit.ai/app/home>). Lifebit CloudOS enables users without any prior cloud  
246 computing knowledge to deploy analysis in the Cloud. In order to run the pipeline the user only  
247 needs to specify input files, desired parameters and select resources from an intuitive graphical user  
248 interface. After the completion of the analysis on Lifebit CloudOS, the user has a permanent  
249 shareable live link that includes performance and file metadata, the associated github repository  
250 revision and also links to the generated results. The relevant analysis page can be used to repeat the  
251 exact same analysis. The analysis page for the PGP-UK user with id uk35C650 can be accessed in  
252 the following permalink <https://cloudos.lifebit.ai/public/jobs/5e3582dae3474100f4665c7a>. Each  
253 analysis can have different privacy settings allowing the user to choose if the results are publicly  
254 visible, making it easier for sharing, or private use, thus maintaining data confidentiality.

### 255 3 Results

256 The main resulting document is a multipage PDF file containing sections relating to variants of  
257 unknown significance, ancestry estimation (as exemplified in Figure 2) and variants with an  
258 associated phenotype, separated by either potentially beneficial or potentially harmful as well as  
259 phenotypes affected by multiple variants, referred to as genosets.

260 To date, more than one hundred such reports have been produced and made available available as  
261 part of the PGP-UK (Beck et al. 2018) and are publicly available in the projects open access data  
262 page (<https://www.personalgenomes.org.uk/data>). This collection contributes to the educational  
263 potential of the project as a whole. On one hand it allows participants of PGP-UK and other users of  
264 the GenomeChronicler tool to compare their genome report results to those of other individuals. On  
265 the other hand, it allows individuals that are interested in the subject but did not have their genome  
266 sequenced to explore what a personal genome looks like.

267 Methods such as GenomeChronicler also allow other research projects in possession of sequencing  
268 data collected from a single individual to easily produce genome reports, customisable with static  
269 text providing information about the project that can be other to the template file, or even the addition  
270 of links to other databases that are relevant.

### 271 4 Conclusions

272 Here we present GenomeChronicler, a computational pipeline to produce genome reports including  
273 variant calling summary data, ancestry inference, and phenotype annotation from genotype data for  
274 personal genomics data obtained through whole genome or whole exome sequencing.

275 The pipeline is modular, fully open source, and available as containers and on the Lifebit CloudOS  
276 computing platform, enabling easy integration with other projects, regardless of computational  
277 resources available and bioinformatics expertise.

278 This work was developed as part of PGP-UK, and incorporates feedback from early participants to  
279 improve the usefulness of the reports produced, and of participant engagement. It is designed to be  
280 easily expandable, adaptable to other contexts and most of all, suit projects with a wide range of  
281 ethical requirements, from those that need the data to be processed inside a safe-haven environment  
282 to those that process all the data in the public domain.

283 Future directions for this work will include the integration of other omics data types that are  
284 produced within PGP-UK, as well as potentially expanding the databases that are linked by default  
285 when running the pipeline.

286 We hope that GenomeChronicler will be useful to other projects and interested individuals. As it is  
287 open source, the pipeline can easily adapt custom templates to satisfy any curiosity-driven analyses  
288 and increase the level of genomic understanding in general. It can also be of interest to educational  
289 groups such as Open Humans (Greshake Tzovaras et al. 2019). Open Humans  
290 (<https://www.openhumans.org/>) is a vibrant community of researchers, patients, data and citizen  
291 scientists who want to learn more about themselves.

## 292 **5 Conflict of Interest**

293 Pablo P. Prieto is CTO of Lifebit and Christina Chatzipantsiou is an employee of Lifebit. All other  
294 authors declare that the research was conducted in the absence of any commercial or financial  
295 relationships that could be construed as a potential conflict of interest.

## 296 **6 Author Contributions**

297 J.A.G.-A. led the development and implementation of the method and wrote the manuscript with  
298 input from all authors. J.A.G.-A., L.C. contributed computer code. C.C. contributed the nextflow and  
299 Lifebit CloudOS integrations with support from P.P.P.. J.A.G.-A., L.C., I.M., A.P.W., S.E., JH.,  
300 O.C. and S.B. contributed to the conceptual development of the method and usability. All authors  
301 read and approved the manuscript.

## 302 **7 Funding**

303 PGP-UK gratefully acknowledges support from the Frances and Augustus Newman Foundation,  
304 Dangoor Education and the National Institute for Health Research (NIHR) UCLH Biomedical  
305 Research Centre (BRC369/CN/SB/101310). The views expressed are those of the authors and not  
306 necessarily those of the NIHR or the Department of Health and Social Care.

## 307 **8 Acknowledgments**

308 The authors acknowledge the use of the UCL Legion High Performance Computing Facility  
309 (Legion@UCL) and associated support services. The authors thank all PGP-UK participants for their  
310 contributions to the project and feedback on the items that feature in the reports produced by  
311 GenomeChronicler.

## 312 **9 Data Availability Statement**

313 The datasets analyzed and used for the development of the approach here described are deposited at  
314 the European Nucleotide Archive (ENA) hosted by the EMBL-EBI under the umbrella accession  
315 PRJEB24961. [<https://www.ebi.ac.uk/ena/data/view/PRJEB24961>]. The PGP-UK pilot data was  
316 described in a data descriptor published in Scientific Data (Chervova et al. 2019). The source code  
317 for the software is deposited and maintained in GitHub and available at [<https://github.com/PGP-UK/GenomeChronicler>]. The nextflow integrated version is available at [<https://github.com/PGP-UK/GenomeChronicler-nf>] and finally, the version also containing the Sarek variant calling pipeline  
318 is available at [<https://github.com/PGP-UK/GenomeChronicler-Sarek-nf>]. Reports generated using  
319  
320

321 this approach for PGP-UK samples are archived in the PGP-UK data page  
322 <https://www.personalgenomes.org.uk/data>.

323 **10 References**

324 Beck, Stephan, Alison M. Berner, Graham Bignell, Maggie Bond, Martin J. Callanan, Olga  
325 Chervova, Lucia Conde, et al. 2018. 'Personal Genome Project UK (PGP-UK): A Research  
326 and Citizen Science Hybrid Project in Support of Personalized Medicine'. *BMC Medical  
327 Genomics* 11 (1): 108. <https://doi.org/10.1186/s12920-018-0423-1>.

328 Cariaso, Michael, and Greg Lennon. 2012. 'SNPedia: A Wiki Supporting Personal Genome  
329 Annotation, Interpretation and Analysis'. *Nucleic Acids Research* 40 (Database issue):  
330 D1308-1312. <https://doi.org/10.1093/nar/gkr798>.

331 Chervova, Olga, Lucia Conde, José Afonso Guerra-Assunção, Ismail Moghul, Amy P. Webster,  
332 Alison Berner, Elizabeth Larose Cadieux, et al. 2019. 'The Personal Genome Project-UK, an  
333 Open Access Resource of Human Multi-Omics Data'. *Scientific Data* 6 (1): 1–10.  
334 <https://doi.org/10.1038/s41597-019-0205-4>.

335 Di Tommaso, Paolo, Maria Chatzou, Evan W. Floden, Pablo Prieto Barja, Emilio Palumbo, and  
336 Cedric Notredame. 2017. 'Nextflow Enables Reproducible Computational Workflows'.  
337 *Nature Biotechnology* 35 (4): 316–19. <https://doi.org/10.1038/nbt.3820>.

338 Ewels, Philip A., Alexander Peltzer, Sven Fillinger, Johannes Alneberg, Harshil Patel, Andreas  
339 Wilm, Maxime Ulysse Garcia, Paolo Di Tommaso, and Sven Nahnse. 2019. 'Nf-Core:  
340 Community Curated Bioinformatics Pipelines'. *BioRxiv*, May, 610741.  
341 <https://doi.org/10.1101/610741>.

342 Garcia, Maxime, Szilveszter Juhos, Malin Larsson, Pall I. Olason, Marcel Martin, Jesper Eisfeldt,  
343 Sebastian DiLorenzo, et al. 2018. 'Sarek: A Portable Workflow for Whole-Genome  
344 Sequencing Analysis of Germline and Somatic Variants'. *BioRxiv*, May, 316976.  
345 <https://doi.org/10.1101/316976>.

346 'Genomelink | Upload Raw DNA Data for Free Analysis On 25 Traits'. 2019. 2019.  
347 <https://genomelink.io/>.

348 Greshake, Bastian, Philipp E. Bayer, Helge Rausch, and Julia Reda. 2014. 'OpenSNP–A  
349 Crowdsourced Web Resource for Personal Genomics'. *PLOS ONE* 9 (3): e89204.  
350 <https://doi.org/10.1371/journal.pone.0089204>.

351 Greshake Tzovaras, Bastian, Misha Angrist, Kevin Arvai, Mairi Dulaney, Vero Estrada-Galíñanes,  
352 Beau Gunderson, Tim Head, et al. 2019. 'Open Humans: A Platform for Participant-Centered  
353 Research and Personal Data Exploration'. *GigaScience* 8 (6).  
354 <https://doi.org/10.1093/gigascience/giz076>.

355 'Initial Sequencing and Analysis of the Human Genome'. 2001. *Nature* 409 (6822): 860–921.  
356 <https://doi.org/10.1038/35057062>.

357 Karczewski, Konrad J., Laurent C. Francioli, Grace Tiao, Beryl B. Cummings, Jessica Alföldi,  
358 Qingbo Wang, Ryan L. Collins, et al. 2019. 'Variation across 141,456 Human Exomes and  
359 Genomes Reveals the Spectrum of Loss-of-Function Intolerance across Human Protein-  
360 Coding Genes'. *BioRxiv*, August, 531210. <https://doi.org/10.1101/531210>.

361 Klein, Teri E., and Marylyn D. Ritchie. 2018. ‘PharmCAT: A Pharmacogenomics Clinical  
362 Annotation Tool’. *Clinical Pharmacology and Therapeutics* 104 (1): 19–22.  
363 <https://doi.org/10.1002/cpt.928>.

364 Kuleshov, Volodymyr, Jialin Ding, Christopher Vo, Braden Hancock, Alexander Ratner, Yang Li,  
365 Christopher Ré, Serafim Batzoglou, and Michael Snyder. 2019. ‘A Machine-Compiled  
366 Database of Genome-Wide Association Studies’. *Nature Communications* 10 (1): 1–8.  
367 <https://doi.org/10.1038/s41467-019-11026-x>.

368 Kurtzer, Gregory M., Vanessa Sochat, and Michael W. Bauer. 2017. ‘Singularity: Scientific  
369 Containers for Mobility of Compute’. *PLOS ONE* 12 (5): e0177459.  
370 <https://doi.org/10.1371/journal.pone.0177459>.

371 Landrum, Melissa J., and Brandi L. Kattman. 2018. ‘ClinVar at Five Years: Delivering on the  
372 Promise’. *Human Mutation* 39 (11): 1623–30. <https://doi.org/10.1002/humu.23641>.

373 Linderman, Michael D., Saskia C. Sanderson, Ali Bashir, George A. Diaz, Andrew Kasarskis, Randi  
374 Zinberg, Milind Mahajan, Sabrina A. Suckiel, Micol Zweig, and Eric E. Schadt. 2018.  
375 ‘Impacts of Incorporating Personal Genome Sequencing into Graduate Genomics Education:  
376 A Longitudinal Study over Three Course Years’. *BMC Medical Genomics* 11 (1): 5.  
377 <https://doi.org/10.1186/s12920-018-0319-0>.

378 Mao, Qing, Serban Ciotlos, Rebecca Yu Zhang, Madeleine P. Ball, Robert Chin, Paolo Carnevali,  
379 Nina Barua, et al. 2016. ‘The Whole Genome Sequences and Experimentally Phased  
380 Haplotypes of over 100 Personal Genomes’. *GigaScience* 5 (1): 42.  
381 <https://doi.org/10.1186/s13742-016-0148-z>.

382 McLaren, William, Laurent Gil, Sarah E. Hunt, Harpreet Singh Riat, Graham R. S. Ritchie, Anja  
383 Thormann, Paul Flicek, and Fiona Cunningham. 2016. ‘The Ensembl Variant Effect  
384 Predictor’. *Genome Biology* 17 (1): 122. <https://doi.org/10.1186/s13059-016-0974-4>.

385 Nakken, Sigve, Ghislain Fournous, Daniel Vodák, Lars Birger Aasheim, Ola Myklebost, and Eivind  
386 Hovig. 2018. ‘Personal Cancer Genome Reporter: Variant Interpretation Report for Precision  
387 Oncology’. *Bioinformatics (Oxford, England)* 34 (10): 1778–80.  
388 <https://doi.org/10.1093/bioinformatics/btx817>.

389 Novembre, John, Toby Johnson, Katarzyna Bryc, Zoltán Kutalik, Adam R. Boyko, Adam Auton,  
390 Amit Indap, et al. 2008. ‘Genes Mirror Geography within Europe’. *Nature* 456 (7218): 98–  
391 101. <https://doi.org/10.1038/nature07331>.

392 Pontikos, Nikolas, Jing Yu, Ismail Moghul, Lucy Withington, Fiona Blanco-Kelly, Tom Vulliamy,  
393 Tsz Lun Ernest Wong, et al. 2017. ‘Phenopolis: An Open Platform for Harmonization and  
394 Analysis of Genetic and Phenotypic Data’. *Bioinformatics* 33 (15): 2421–23.  
395 <https://doi.org/10.1093/bioinformatics/btx147>.

396 ‘Promethease’. 2019. 2019. <https://www.promethease.com/>.

397 Purcell, Shaun, Benjamin Neale, Kathe Todd-Brown, Lori Thomas, Manuel A. R. Ferreira, David  
398 Bender, Julian Maller, et al. 2007. ‘PLINK: A Tool Set for Whole-Genome Association and  
399 Population-Based Linkage Analyses’. *American Journal of Human Genetics* 81 (3): 559–75.  
400 <https://doi.org/10.1086/519795>.

401 Ramos, Erin M, Douglas Hoffman, Heather A Junkins, Donna Maglott, Lon Phan, Stephen T Sherry,  
402 Mike Feolo, and Lucia A Hindorff. 2014. ‘Phenotype–Genotype Integrator (PheGenI):  
403 Synthesizing Genome-Wide Association Study (GWAS) Data with Existing Genomic

404        Resources'. *European Journal of Human Genetics* 22 (1): 144–47.  
405        <https://doi.org/10.1038/ejhg.2013.96>.

406        Salari, Keyan, Konrad J. Karczewski, Louanne Hudgins, and Kelly E. Ormond. 2013. 'Evidence That  
407        Personal Genome Testing Enhances Student Learning in a Course on Genomics and  
408        Personalized Medicine'. *PLoS One* 8 (7): e68853.  
409        <https://doi.org/10.1371/journal.pone.0068853>.

410        Sanderson, Saskia C., Michael D. Linderman, Sabrina A. Suckiel, George A. Diaz, Randi E. Zinberg,  
411        Kadija Ferryman, Melissa Wasserstein, Andrew Kasarskis, and Eric E. Schadt. 2016.  
412        'Motivations, Concerns and Preferences of Personal Genome Sequencing Research  
413        Participants: Baseline Findings from the HealthSeq Project'. *European Journal of Human  
414        Genetics* 24 (1): 14–20. <https://doi.org/10.1038/ejhg.2015.118>.

415        Sochat, Vanessa V., Cameron J. Prybol, and Gregory M. Kurtzer. 2017. 'Enhancing Reproducibility  
416        in Scientific Computing: Metrics and Registry for Singularity Containers'. *PLOS ONE* 12  
417        (11): e0188511. <https://doi.org/10.1371/journal.pone.0188511>.

418        The 1000 Genomes Project Consortium. 2015. 'A Global Reference for Human Genetic Variation'.  
419        *Nature* 526 (7571): 68–74. <https://doi.org/10.1038/nature15393>.

420        Van der Auwera, Geraldine A., Mauricio O. Carneiro, Chris Hartl, Ryan Poplin, Guillermo del  
421        Angel, Ami Levy-Moonshine, Tadeusz Jordan, et al. 2013. 'From FastQ Data to High  
422        Confidence Variant Calls: The Genome Analysis Toolkit Best Practices Pipeline'. *Current  
423        Protocols in Bioinformatics / Editorial Board, Andreas D. Baxevanis ... [et Al.]* 11 (1110):  
424        11.10.1-11.10.33. <https://doi.org/10.1002/0471250953.bi1110s43>.

425        Venter, J. Craig. 2010. 'Multiple Personal Genomes Await'. *Nature* 464 (7289): 676–77.  
426        <https://doi.org/10.1038/464676a>.

427

## 428        11      Figure Legends

429        Figure 1: Flow Diagram of GenomeChronicler processing pipeline, illustrating the multiple entry  
430        points for the pipeline, resources integrated by default and generated outcomes. Either entry point of  
431        the pipeline can be run locally in a single machine, as a nextflow workflow or in the Cloud. All  
432        source code and integrations are freely available in their respective GitHub repositories. The stand-  
433        alone GenomeChronicler is available at (<https://github.com/PGP-UK/GenomeChronicler>), the  
434        integration of GenomeChronicler with nextflow is available at ([https://github.com/PGP-UK/GenomeChronicler-nf](https://github.com/PGP-<br/>435        UK/GenomeChronicler-nf)) and the combined GenomeChronicler with Sarek variant calling is  
436        available at (<https://github.com/PGP-UK/GenomeChronicler-Sarek-nf>). The recipe files for the  
437        Docker and Singularity containers are available within the respective GitHub repositories. The  
438        resource logos are reproduced from the respective resource websites and remain copyright of their  
439        original owner.

440        Figure2: Example Ancestry PCA plot containing the current reference data from the 1000 genomes  
441        project used by GenomeChronicler, with shaded areas broadly illustrating the origin of the  
442        populations represented.

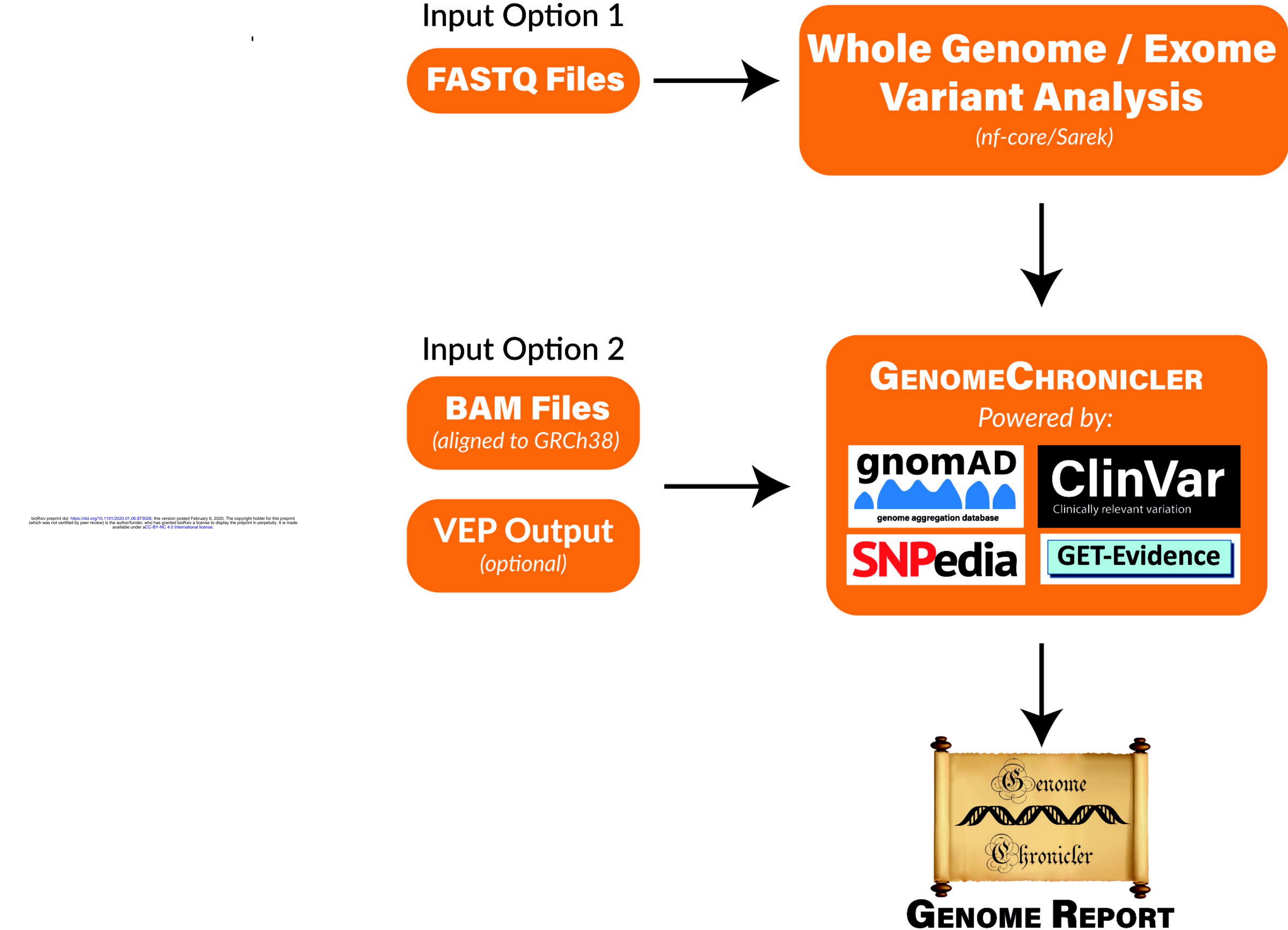


Figure 2

