1

2

3 **Bystro: rapid online variant annotation and natural-language filtering at whole-genome**

4 **scale**

5 Alex V. Kotlar[1], Cristina E. Trevino[1], Michael E. Zwick[1], David J. Cutler[1], and Thomas S.

6 Wingo[1,2,3,*]

7 [1]Department of Human Genetics, Emory University School of Medicine

8 [2]Division of Neurology, Atlanta VA Medical Center.

9 [3]Department of Neurology, Emory University School of Medicine.

10

11 *Corresponding author:

12 Thomas S. Wingo

13 505K Whitehead Building

14 615 Michael Street NE

15 Atlanta, GA 30322-1047

16 404-727-4905 (office)

17 404-727-3728 (fax)

18 *thomas.wingo@emory.edu*

19
20

1

**Abstract**

Accurately selecting relevant alleles in large sequencing experiments remains technically challenging. Bystro (*https://bystro.io/*) is the first online, cloud-based application that makes variant annotation and filtering accessible to all researchers for terabyte-sized whole-genome experiments containing thousands of samples. Its key innovation is a general-purpose, natural-language search engine that enables users to identify and export alleles and samples of interest in milliseconds. The search engine dramatically simplifies complex filtering tasks that previously required programming experience or specialty command-line programs. Critically, Bystro's annotation and filtering capabilities are orders of magnitude faster than previous solutions, saving weeks of processing time for large experiments.

**Keywords**

Natural-language search, genomics, bioinformatics, annotation, filtering, web, online, cloud, big data

**Background**

While genome-wide association studies (GWAS) and whole-exome sequencing (WES) remain important components of human disease research, the future lies in whole-genome sequencing (WGS), as it inarguably provides more complete data. The central challenge posed by WGS is one of scale. Genetic disease studies require thousands of samples to obtain adequate power, and the resulting WGS datasets are hundreds of gigabytes in size and contain tens of millions of variants. Manipulating data at this scale is difficult. To find the alleles that contribute to traits of interest, two steps must occur. First, the variants identified in a sequencing experiment need to be described in a process called annotation, and second, the relevant alleles need to be selected based on those descriptions in a procedure called variant filtering.

2

46       Annotating and filtering large numbers of variant alleles requires specialty software.

47    Existing annotators, such as ANNOVAR[1], SeqAnt[2], VEP[3], and GEMINI[4] have played an

48    important research role, and are sufficient for small to medium experiments (e.g.,10s to 100s of

49    WES samples). However, they require significant computer science training to use in offline,

50    distributed computing environments, and have substantial restrictions in terms of performance

51    and the maximum size of the data they will annotate online. Existing variant filtering solutions

52    are even more limited, with most analyses requiring researchers to program custom scripts,

53    which can result in errors that impact reproducibility[5]. Therefore, annotation and filtering are

54    not readily accessible to most scientists, and even bioinformaticians face challenges of

55    performance, cost and complexity.

56       Here we introduce an application called Bystro that significantly simplifies variant

57    annotation and filtering, while also improving performance by orders of magnitude and saving

58    weeks of processing time on large data sets. It is the first program capable of handling

59    sequencing experiments on the scale of thousands of whole-genome samples and tens of

60    millions of variants online in a web browser, and integrates the first, to our knowledge, publicly-

61    available, online natural-language search engine for filtering variants and samples from these

62    experiments. The search engine enables real-time (sub-second), nuanced variant filtering, both

63    across all samples and per sample, using simple phrases and interactive, web-based filters.

64    Bystro makes it possible to efficiently find alleles of interest in any sequencing experiment

65    without computer science training, improving reproducibility while reducing annotation and

66    filtering costs.

67

68    **Results**

69       To compare Bystro's capabilities with other recent programs, we submitted 1000

70    Genomes[6] Phase 1 and Phase 3 VCF files for annotation and filtering (Figure 1). Phase 1

71    contains 39.4 million variants from 1,092 WGS samples, while Phase 3 includes 84.9 million

3

72    alleles from 2,504 WGS samples. We first evaluated the online capabilities of the web-based

73    versions of Bystro, wANNOVAR[7], VEP, and GEMINI (running on the Galaxy[8] platform).

74    Bystro was the only program able to complete either 1000 Genomes Phase 1 or Phase 3 online,

75    and was also the only application to handle a $6x10^6$ variant subset of Phase 3, a size

76    representative of modest whole-genome experiments. When tested with $5x10^4 - 1x10^6$ variant

77    subsets of 1000 Genomes Phase 3, Bystro was approximately 144 – 212x faster than

78    GEMINI/Galaxy in generating a downloadable annotation and searchable result database, and

79    was significantly easier to use, as it did not require a separate annotation step (Figure 2). When

80    tested on a small trio data set, Bystro was able to identify *de novo* variants without any

81    additional software, and was 45x faster than GEMINI's de_novo tool (Additional file 1: Table

82    S1). Bystro and GEMINI/Galaxy produced similarly detailed outputs, with Bystro offering fewer,

83    but more complete and recent sources, as well as more detailed annotations for some classes

84    of data (Additional file 1: Table S2 ; Additional file 2). Notably GEMINI was found to work only

85    with the hg19 human genome assembly, whereas Bystro supports hg19, hg38, and a variety of

86    model organisms.

Figure 1 | **Using Bystro online to find alleles of interest in sequencing experiments. A**)

After logging in (*https://bystro.io/*), users upload one of more VCF or SNP-format files -

containing alleles from a sequencing experiment - from a computer or a connected Amazon

S3 bucket. Datasets of over 890GB, containing thousands of samples and tens of millions of

variants are supported. The data is rapidly annotated in the cloud, using descriptions from

public sources (e.g. RefSeq, dbSNP, Clinvar, and others). The annotated results can be

filtered using Bystro's natural-language search engine, and any search results can be saved

as new annotations. Annotated experiments and saved results can be viewed online,

downloaded as tab-delimited text, or uploaded back to linked Amazon S3 buckets. **B**) An

example of using Bystro's natural-language search engine to filter 1000 Genomes Phase 3

(*https://bystro.io/public*). To do so, users may type natural phrases, specific terms, numerical

ranges, or apply filters on any annotated field. Queries are flexible, allowing misspelled terms

such as "earl-onset" to accurately match. Complex tasks, such as identifying *de novo*

variants can be achieved by using Boolean operators (AND, OR, NOT, +, -), exact-match

filters, and user-defined terms. For instance, after labeling the "proband" and their "parents",

the user could simply search *proband –parents*, or combine with additional parameters for

more refined queries, i.e. *proband –parents missingness < .1 gnomad.exomes.af_nfe < .001.*

87

5

Figure 2 | **Online performance comparison of Bystro, VEP, wANNOVAR, and GEMINI**.

Bystro, wANNOVAR, VEP, and GEMINI (running on Galaxy) we run under similar conditions. Total processing time was recorded for 1000 Genomes Phase 3 WGS VCF files, containing either the full data set (2,504 samples, $8.49 \times 10^7$ variant sites), or subsets (2,504 samples and $5 \times 10^4$, $3 \times 10^5$, $1 \times 10^6$, and $6 \times 10^6$ variants). Only Bystro successfully processed more than $1 \times 10^6$ variants online: wANNOVAR (not shown) could not complete the smallest $5 \times 10^4$ variant subset; VEP could not complete more than $5 \times 10^4$ variants; and GEMINI/Galaxy could not complete more than $1 \times 10^6$ variants. Online, VEP outputted a restricted subset of annotation data compared to its offline version. GEMINI and Bystro (but not VEP) outputted whole-genome CADD scores, while only Bystro also returned whole-genome PhyloP and PhastCons conservation scores. Bystro was faster than GEMINI/Galaxy by 144x-212x across all time points.

88

89      We next tested offline performance on identical servers to gauge performance in the

90      absence of web-related file-size and networking limitations. Bystro was 113x faster than

91      ANNOVAR and up to 790x faster than VEP, annotating all $8.5 \times 10^7$ variants and 2,504 samples

92      from Phase 3 in less than 3 hours (Table 1). Furthermore, ANNOVAR was unable to finish

93      either Phase 1 or Phase 3 annotations due to memory requirements (exceeding 60GB of RAM),

94      and VEP annotated Phase 3 at a rate of 10 variants per second, indicating that it would need at

95      least 98 days to complete. Critically, Bystro's run time grew linearly with the number of

96      submitted genotypes, suggesting that it could handle even hundreds of thousands of samples

97      within days.

98      While offering significantly faster performance, Bystro also provided 3.5x the number of

99      annotation output fields as ANNOVAR and 5.6x that of VEP (Additional file 3**).** Notably, unlike

100     ANNOVAR or VEP, Bystro annotated each sample relative to its genotype, reporting

101     homozygosity, heterozygosity, missingness, sample minor allele frequency, and labeling each

6

102    sample as homozygous, heterozygous, or missing. In contrast, ANNOVAR provided only

103    sample minor allele frequency, while VEP reported no sample-level data. We note that VEP is

104    capable of providing per-sample annotations (heterozygosity/homozygosity status), but we were

105    unable to use this feature for performance reasons. A detailed comparison of the exact settings

106    used is given (Additional file 2 ; Additional file 3).

107        To investigate annotation accuracy, we next compared Bystro with ANNOVAR and VEP

108    on a previously-analyzed synthetic dataset[9]. Overall, excellent concordance between all

109    methods was noted (Additional files 4, 5, and 6). For instance, in comparison with ANNOVAR,

110    allele position (>98%), allele identity (100%), and variant effects (>99%) were highly consistent

111    across all classes of variation, for sites that Bystro did not exclude for quality reasons

112    (Additional file 4).

113        In cases where the annotators disagreed, Bystro gave the more correct interpretations.

114    For instance, Bystro and VEP excluded reference sites (ALT: "."), while ANNOVAR annotated

115    such loci as "synonymous SNV"; it is of course incorrect to call reference sites variant

116    (Additional file 4 ; Additional file 5). In cases of insertions and deletions, which are often

117    ambgiuously represented in VCF files due to the format's padding requirements, Bystro always

118    provided the parsimonious left-shifted representation, while ANNOVAR and VEP occasionally

119    right-shifted variants (Additional file 4 ; Additional file 5). This is evident at

120    chr15:42680000CA>CAA, where both ANNOVAR and VEP called the insertion as occuring after

121    the first "A", with 2 bases of padding, rather than the simpler option after the first base, "C", with

122    1 base of padding (Additional file 1: Table S3). Similar results were found at multiallelic loci with

123    complex indels (Additional file 1: Table S4).

124        Similarly, in cases where Bystro and ANNOVAR or VEP disagreed on variant

125    consequences, Bystro always appeared correct relative to the underlying transcript set. For

126    example, in the case of the simple insertion chr19:41123094G>GG, Bystro correctly identified

127    all three overlapping transcripts (NM_003573;NM_001042544;NM_001042545), and noted the

7

128    variant as coding (exonic) relative to all three. In contrast, ANNOVAR called the allele as

129    disrupting a splice site, despite the fact that the nerest intron, and therefore splice site, was

130    37bp downstream (Additional file 1: Figure S1).

131         Additionally, Bystro's strict VCF quality control measures substantially improved

132    annotation accuracy.This is evident in the case of gnomAD, a VCF-format dataset that

133    represents the largest experiment on human genetic variation. While Bystro and ANNOVAR

134    provided identical gnomAD data for 93.7% of tested alleles, the remaining 6.3% were low-

135    quality gnomAD results that were included in ANNOVAR and excluded from Bystro (Additional

136    file 4). For instance, in the case of chr16:2103394C>T, ANNOVAR reported rs760688660,

137    which failed gnomAD's random forest qc step. We note that a 6.3% false-positive rate is similar

138    to the frequency of common variation, and significantly larger than the frequency of rare

139    variants, making ANNOVAR's gnomAD annotations a potentially unreliable source of data for

140    both common and rare variant filtering.

141

Table 1 | **Bystro, VEP, ANNOVAR offline command-line performance**.

| Software | Dataset | Samples | Variants | Variants/s | Bystro vs |
|---|---|---|---|---|---|
| **Bystro** | **1000G Phase 3 chr1** | **2504** | **$1\times10^6$** | **8156 ± 195** | - |
| | **1000G Phase 3 chr1** | **2504** | **$2\times10^6$** | **8484 ± 67.9** | - |
| | **1000G Phase 3 chr1** | **2504** | **$4\times10^6$** | **8516 ± 57.2** | - |
| | **1000G Phase 3 chr1** | **2504** | **$6.5\times10^6$** | **7779 ± 21.8** | - |
| | **1000G Phase 1** | **1092** | **$3.9\times10^7$** | **5417 ± 76.8** | |
| | **1000G Phase 3** | **2504** | **$8.5\times10^7$** | **7904 ± 15.9** | - |
| **VEP** | 1000G Phase 1 | 1092 | $3.9\times10^7$ | 18.67 ± 0.58 | 290x |
| | 1000G Phase 3 | 2504 | $8.5\times10^7$ | 10.00 ± 0.00 | 790x |
| **ANNOVAR** | 1000G Phase 3 chr1 | 2504 | $1\times10^6$ | 74.67 ± 0.21 | 109x |
| | 1000G Phase 3 chr1 | 2504 | $2\times10^6$ | 75.32 ± 0.06 | 113x |
| | 1000G Phase 3 chr1 | 2504 | $4\times10^6$ | 75.15 ± 0.39 | 113x |
| | 1000G Phase 3 chr1 | 2504 | $6.5\times10^6$ | NA | NA |
| | 1000G Phase 1 | 1092 | $3.9\times10^7$ | NA | NA |
| | 1000G Phase 3 | 2504 | $8.5\times10^7$ | NA | NA |

Bystro, VEP, and ANNOVAR were similarly configured with 8 threads on Amazon i3.2xlarge servers. "Dataset" refers to the VCF file used. "Variants/s" is the average of three trials. VEP performance was recorded after $2\times10^5$ sites in consideration of time. In runs of $1\times10^6$ or more annotated sites, VEP performance did not deviate from the $2\times10^5$ value. ANNOVAR could not complete the full Phase 1, Phase 3, or Phase 3 chromosome 1 datasets due to memory limitations. Thus, ANNOVAR was compared to Bystro on subsets of 1000 Genomes Phase 3 chromosome 1. Bystro run times included time taken to compress outputs. 1000 Genomes Phase 1 performance reflects IO limitations.

142    Next, we explored the Bystro search engine's ability to filter the 84.9 million annotated

143    Phase 3 variants. Bystro's search engine was unique in its natural-language capabilities, and no

144    other tested online program could handle the full Phase 3 dataset, or subsets as large as $6\times10^6$

145    variants (Figure 2). First, we used Bystro's search engine to find all alleles in exonic regions by

146    entering the term "exonic" (933,343 alleles, 0.030 ± .001 seconds, Table 2). The search engine

147    calculated a transition to transversion ratio of 2.96 for the query, consistent with previously

148    observed values in coding regions. To refine results to rare, predicted deleterious alleles, we

149    queried "cadd > 20 maf < .001 pathogenic expert review missense" (65 alleles, 0.029 ± 0.025s,

9

150    Table 2). This search query could be written using partial words ("pathogen"), possessive nouns

151    ("expert's"), different tenses ("reviews"), and synonyms ("nonsynonymous") without changing

152    the results.

Table 2 | **Online comparison of Bystro and recent programs in filtering 8.49x10$^7$ variants from 1000 Genomes**

| Group | Search query | Time (s) | Variants | Tr:Tv |
|---|---|---|---|---|
| 1 | exonic | 0.030 ± 0.030 | 993,343 | 2.96 |
| 2 (a) | cadd > 20 maf < .001 pathogenic expert review missense | 0.029 ± 0.009 | 65 | 1.71 |
| 2 (b) | cadd > 20 maf < .001 pathogenic **expert's** review **non-synonymous** | 0.036 ± 0.019 | 65 | 1.71 |
| 2 (c) | cadd > 20 maf < .001 **pathogen** expert-**reviewed nonsynonymous** | 0.044 ± 0.025 | 65 | 1.71 |
| 3 (a) | early onset breast cancer | 0.046 ± 0.029 | 4,335 | 2.51 |
| 3 (b) | **early-onset** breast cancer | 0.037 ± 0.020 | 4,335 | 2.51 |
| 3 (c) | **Early onset** breast **cancers** | 0.033 ± 0.015 | 4,335 | 2.51 |
| 4 (a) | Pathogenic nonsense Ehlers-Danlos | 0.038 ± 0.027 | 1 | NA |
| 4 (b) | **pathogenic** nonsense **E.D.S** | 0.078 ± 0.087 | 1 | NA |
| 4 (c) | **pathogenic stopgain eds** | 0.040 ± 0.022 | 1 | NA |

The full 1000 Genomes Phase 3 VCF file (853GB, 8.49x10$^7$ variants, 2,504 samples) was filtered in the publicly-available Bystro web application using the Bystro natural-language search engine. VEP, GEMINI, and wANNOVAR (not shown) were also tested, but were unable to annotate this data set or filter it. Bystro's search engine uses a natural language parser that allows for unstructured queries: queries in groups 2, 3, and 4 show phrasing variations that did not affect results returned, as would be expected for a search engine that could handle normal language variation. "Tr:Tv" is the transition to transversion ratio automatically calculated for each query by the search engine. The transition to transversion ratio of 2.96 for the "exonic" query is close to the ~2.8-3.0 ratio expected in coding regions, suggesting that the search engine accurately identified exonic (coding) variants.

10

153       To test the search engine's ability to accurately match variants from full-text disease

154    queries, we first searched "early-onset breast cancer", returning the expected alleles in *BRCA1*

155    and *BRCA2* (4,335 variants, .037 ± .020s, Table 2). Notably, the queried phrase "early-onset

156    breast cancer" did not exist within the annotation, and instead matched closely-related RefSeq

157    transcript names, such as "Homo sapiens breast cancer 2, early onset (BRCA2), mRNA." We

158    next explored Bystro's ability to handle synonyms and acronyms. To test the hypothesis that

159    Bystro could interpret common ontologies, we queried "pathogenic nonsense E.D.S", where

160    "nonsense" is a common synonym for "stopGain" (a term annotated by the Bystro annotation

161    engine), and "E.D.S" is an acronym for "Ehlers-Danlos Syndrome". Bystro successfully parsed

162    this query, returning a single *PLOD1* variant found in 1000 Genomes Phase 3 that introduces an

163    early stop codon in all three of its overlapping transcripts, and which has been reported in

164    Clinvar as "pathogenic" for "Ehlers-Danlos syndrome, type 4" (1 variant, .038s ± .027s, Table 2).

165       Since no other tested program could load or filter the 1000 Genomes Phase 3 VCF file

166    online, we next compared Bystro to GEMINI (running on the Galaxy platform) on subsets of

167    1000 Genomes Phase 3. In contrast with GEMINI's structured SQL queries, Bystro enabled

168    shorter and more flexible searches. For instance, to return all missense, rare variants with

169    CADD Phred scores larger than 15, GEMINI required a 162 character SQL query, while Bystro

170    needed only 36 characters. Bystro also demonstrated synonym support, returning identical

171    results for "missense" and "nonsynonymous" queries. Critically, Bystro's search engine enabled

172    real-time (sub-second) filtering, performing approximately four orders of magnitude faster than

173    GEMINI on Galaxy while searching and returning similar volumes of data (Table 3).

174       To test the accuracy of Bystro's search engine relative to the underlying annotation, we

175    first compared Bystro's natural-language queries with Bystro's "Filters", which provide a

176    complimentary, exact-match filtering option. All results were identical between the two methods

177    (Additional file 1: Table S5). To control for the possibility that Bystro's "Filters" were biased, we

178    created separate Perl filtering scripts that searched for exact matches within the underlying tab-

11

179    delimited text annotation. Again, results were completely concordant (Additional file 1: Table

180    S5). Finally, to control for the possibility that both Bystro's "Filters" and the Perl scripts were

181    biased due to the programmer, we compared Bystro's natural-language queries with Excel

182    filters on a smaller dataset that could be manually examined. The queries were found

183    completely specific in this comparison as well (Additional file 1: Table S6; Additional file 7).

**Table 3 | Online comparison of Bystro and GEMINI/Galaxy in filtering 1x10⁶ variants**

| # | *Program* | Query | Time (s) | Variants | Ts/Tv |
|---|-----------|-------|----------|----------|-------|
| *1* | **Bystro** | **cadd > 15 alt:(a \|\| c \|\| t \|\| g)** | **.004 ± 0** | **28,099** | **2.512** |
| 1 | GEMINI | SELECT * FROM variants JOIN variant_impacts ON variants.variant_id = variant_impacts.variant_id WHERE cadd_scaled > 15 | 442 ± 87 | 22,063 | NA |
| *2* | **Bystro** | **gnomad.exomes.af < .001 cadd > 15 missense** | **.007 ± .003** | **6,840** | **3.083** |
| 2 | GEMINI | SELECT * FROM variants JOIN variant_impacts ON variants.variant_id = variant_impacts.variant_id WHERE cadd_scaled > 15 AND aaf_exac_all < .001 AND variant_impacts.impact = 'missense_variant' | 77.6 ± 18.6 | 5,160 | NA |
| *3* | **Bystro** | **gnomad.exomes.af < .001 cadd > 15 nonsynonymous** | **.006 ± .001** | **6,840** | **3.083** |
| 3 | GEMINI | SELECT * FROM variants JOIN variant_impacts ON variants.variant_id = variant_impacts.variant_id WHERE cadd_scaled > 15 AND aaf_exac_all < .001 AND variant_impacts.impact = '**nonsynonymous_variant**' | NA | 0 | NA |

Bystro was compared to the latest hosted version of GEMINI (v0.8.1, on the Galaxy platform) in filtering the 1x10⁶ variant subset of 1000 Genomes Phase 3, which was the largest tested file that GEMINI/Galaxy could process. GEMINI requires structured SQL queries, while Bystro allows for shorter, unstructured search. In query #1, Bystro searched for CADD scores only within single-nucleotide polymorphisms (using alt:(a || c || t || g), or equivalently the regex query alt:/[actg]/), to normalize results with GEMINI, which provides no CADD data for insertions and deletions. In queries #2 and #3, Bystro's search engine returned identical results for the synonymous terms "missense" and "nonsynonymous", despite annotating such sites only as "nonsynonymous". In contrast, GEMINI required the specific term 'missense_variant'. GEMINI/Galaxy and Bystro returned different results because the latest version of GEMINI on Galaxy (0.8.1) uses outdated annotation sources. Comparisons between Bystro and GEMINI/Galaxy are further limited as GEMINI doesn't provide a natural-language parser, annotation field filters, an interactive result browser, per-query statistics, or the ability to filter saved search results. Notably, Bystro also performed substantially faster, returning all results in less than 1 second.

184

13

185

186    **Discussion**

187    The Bystro annotation and filtering capabilities are primarily exposed through a public

188    web application (*https://bystro.io/*), and are also available for custom, offline installation. To

189    ensure data safety, Bystro follows industry recommendations for password management, in-

190    transit data security, and at-rest data security. Input and output files are encrypted at rest on

191    Amazon EFS file systems, using AES 256-bit encryption, and every request for annotation or

192    search data is authenticated by the web server using short-lived identity tokens. To further

193    protect user data, annotation and search services are not directly open to the Internet, but

194    require routing and authentication through the web server. Furthermore, all web traffic is

195    encrypted using TLS (HTTPS), and password hashing follows the National Institute of

196    Standards and Technology (NIST) recommended PBKDF2-HMAC-SHA512 strategy.

197    Creating an annotation online is as simple as selecting the genome and assembly used

198    to make the variant call format (VCF)[10] or SNP[11] format files, and uploading these files from

199    a computer or Amazon S3 bucket, which can be easily linked to the web application. Annotation

200    occurs in the cloud, where distributed instances of the Bystro annotation engine process the

201    data and send the results back to the web application for storage and display (Figure 1).

202    The Bystro annotation engine is open source, and supports diverse model organisms

203    including *Homo sapiens* (hg19, hg38), *M. musculus* (mm9, mm10), *R. macaque* (rheMac8), *R.*

204    *norvegicus* (rn6), *D. melanogaster* (dm6), *C. elegans* (ce11), *S. cerevisiae* (sacCer3). To

205    annotate, it rapidly matches alleles from users' submitted files to descriptions from RefSeq[12],

206    dbSNP[13], PhyloP[14], PhastCons[14], Combined Annotation-Dependent Depletion (CADD),

207    Clinvar[15], and gnomAD[16]. For custom installations, Bystro supports Ensembl, RefSeq, or

208    UCSC Known Genes transcript sets, and can be flexibly configured include annotations from

209    any files in genePredExt, wigFix, BED, or VCF formats.

14

210     The annotation engine is aware of alternate splicing, and annotates all variants relative

211     to each alternate transcript. When provided sample information, Bystro also annotates all

212     variants relative to all sample genotypes. In such cases, at every site it labels each sample as

213     homozygous, heterozygous, or missing, and also calculates the heterozygosity, homozogosity,

214     missingness, and sample minor allele frequency. Furthermore, in contrast with current programs

215     that require substantial VCF file pre-processing, Bystro automatically removing low-quality sites,

216     normalizes variant representations, splits multi-allelic variants, and checks the reference allele

217     against the genome assembly. Critically, Bystro's algorithm guarantees parsimonious (left-

218     shifted) variant representations, even for multi-allelic sites containing complex insertions and

219     deletions.

220     The Bystro annotation engine is designed to scale to any size experiment, offering the

221     speed of distributed computing solutions such as Hail[17], but with less complexity. Current well-

222     performing annotators - such as ANNOVAR and SeqAnt - load significant amounts of data into

223     memory to improve performance. However, when these programs use multiple threads to take

224     advantage of multicore CPUs they may exceed available memory (in some cases over 60GB),

225     resulting in a sharp drop in performance or system crash. To solve this, Bystro annotates

226     directly from an efficient memory-mapped database (LMDB), using only a few megabytes per

227     thread, and because memory-mapped databases naturally lend themselves to the caching

228     frequently accessed data, Bystro achieves most of the benefits of in-memory solutions, but

229     without the per-thread penalties. This approach allows Bystro to take excellent advantage of

230     multicore CPUs, while also enabling it to perform well on inexpensive, low-memory machines.

231     Critically, when multiple files are submitted to it simultaneously, the Bystro annotation engine

232     can automatically distribute the work throughout the cloud (or a user-configured computer

233     cluster), gaining additional performance by processing the files on multiple computers (Figure

234     1). Furthermore, in reflection of the large sizes of both input sequencing experiments and the

235     corresponding annotation outputs - on the order of terabytes for modern whole-genome

15

236    experiments - Bystro accepts compressed input files, and directly writes compressed outputs.

237    This ability to directly write compressed annotations with no uncompressed intermediate is

238    critical given the rapid growth in sequencing experiment size.

239        When the web application receives a completed annotation, it saves the data and

240    creates a permanent results page. Detailed information about the annotation, such as the

241    database version used for the annotation is stored in a log file that the user may download.

242    Users may then explore several quality control metrics, including the transition to transversion

243    ratio on a per-sample or per-experiment basis. They may also download the results as tab-

244    delimited text to their computer, or upload them to any connected Amazon S3 bucket. In parallel

245    with the completion of an annotation, the Bystro search engine automatically begins indexing

246    the results. Once finished, a search bar is revealed in the results page, allowing users to filter

247    their variants using the search engine (Figure 1).

248        Unlike existing filtering solutions, Bystro's Elasticsearch-based natural-language search

249    engine accepts unstructured, "full-text" queries, and relies on a sophisticated language parser to

250    match annotated variants. This allows it to offer the flexibility of modern search engines like

251    Google and Bing, while remaining specific enough for the precise identification of alleles

252    relevant to the research question. The Bystro search engine matches terms regardless of

253    capitalization, punctuation, or word tense, and accurately finds partial terms within long

254    annotation values. Like the annotation engine, the search engine is also exceptionally fast,

255    automatically distributing indexed annotations throughout the cloud, enabling users to sift

256    through millions of variants from large whole-genome sequencing experiments in milliseconds.

257        In order to provide flexible, but specific matches without relying on structured SQL

258    queries, the search engine identifies the data type of every value in the annotation. Text

259    undergoes stemming and lemmatization, which reduces the influence of grammatical variation,

260    and is then tokenized into left-edge n-grams, which allows for flexible matching. Numerical data

261    is stored in the smallest integer or float format that can accommodate it, allowing for rapid and

16

262   accurate range queries. For complex queries, the search engine supports Boolean operators

263   (AND, OR), regular expressions, and Levenshtein-edit distance fuzzy matches. It also has a

264   built-in dictionary of synonyms, for instance equating "stopgain" and "nonsense".

265        In some cases, text will match accurately, but not specifically; this most often happens

266   with short, generic terms. For instance, querying "intergenic" alone may match the word

267   "intergenic" in "long intergenic non-protein coding RNA" in refSeq's description field, as well as

268   "intergenic" in the refSeq's siteType field. To help improve accuracy in such cases, Bystro

269   provides three, closely related features: 1) "Aggregations" allows users to see the top 200

270   values for any text field, or equivalently the min, max, mean, standard deviation (and other

271   similar statistics) for any numerical field. This allows users to quickly and precisely understand

272   the composition of search results, as well as to generate summary statistics. 2) "Filters" allows

273   users to refine queries, by forcing the inclusion or exclusion of any values found in any field. For

274   instance, rather than query "intergenic", it may be easier and more precise to simply click on the

275   "refSeq.siteType" filter, and select the "intergenic" value. Any number of "Filters" may be

276   combined with any natural-language query, containing up to 1 million words. 3) Bystro allows

277   field names within a natural-language query for added specificity. For example, rather than

278   searching for "intergenic", the user could type "refSeq.siteType:intergenic", to indicate that they

279   wished to match "intergenic" specifically in the refSeq.siteType annotation field.

280        Bystro's search engine also includes several features to increase flexibility beyond the

281   contents of the annotation: 1) "Custom Synonyms" allows users to define their own terms and

282   annotations. Among other uses, this make it is possible to label trios, which can be used to

283   easily identify *de novo* variants and test allele transmission models. 2) "Search Tools" are small

284   programs, accessible by a single mouse click, that dynamically modify any query to generate

285   complex result summaries. Some of their functions include identifying compound heterozygotes.

286   3) "Statistical Filters" dynamically perform statistical tests on the variants returned from any

287 query. For instance, the "HWE" filter allows users to exclude variants out of Hardy-Weinberg

288 Equilibrium. This is an often-needed quality control step.

289 Most importantly, there is no limit to the number of query terms and "Filters" that can be

290 combined, and users can save and download the results of any search query, which enables

291 recursive filtering on a single dataset. The saved results are indexed for search, and hyperlinked

292 to the annotations that they were generated from, forming permanent records that can be used

293 to reproduce complex analyses. This multi-step filtering provides functionality similar to custom

294 command-line filtering script pipelines, but is significantly faster, less error prone, and

295 accessible to researchers without programming experience.

296

297 While Bystro's annotation and filtering performance is currently unparalleled by any other

298 approach, other software (such as Hail[17]) could achieve similar performance by implementing

299 distributed computing algorithms like MapReduce[18], and spreading annotation workloads

300 across many servers. Bystro demonstrates that these workarounds are unnecessary to achieve

301 reasonable run-times for large datasets online or offline. Additionally, while Bystro's natural-

302 language search engine significantly reduces the difficulty of variant filtering, it does not handle

303 language idiosyncrasies as robustly as more mature solutions like Google's, and may return

304 unexpected results when search queries are very short and non-specific, since such queries

305 may have multiple correct matches. This is easily avoided by using longer phrases, by using

306 "Custom Synonyms" to define more specific terms, by examining the composition of results

307 using "Aggregations", or by applying "Filters" to precisely filter results. Such considerations and

308 options are well-documented in Bystro's online user guide (*https://bystrio.io/help*).

309

310 **Conclusions**

311 To date, identifying alleles of interest in sequencing experiments has been time-

312 consuming and technically challenging, especially for whole-genome sequencing experiments.

18

313     Bystro increases performance by orders of magnitude and improves ease of use through three

314     key innovations: 1) a low-memory, high-performance, multithreaded variant annotator that

315     automatically distributes work in cloud or clustered environments; 2) an online architecture that

316     handles significantly larger sequencing experiments than previous solutions; and 3) the first

317     publicly-available, general-purpose, natural-language search engine for variant filtering in

318     individual research experiments. Bystro annotates large experiments in minutes, and its search

319     engine is capable of matching variants within whole-genome datasets in milliseconds, enabling

320     real-time data analysis. Bystro's features enable practically any researcher – regardless of their

321     computational experience - to analyze large sequencing experiments (e.g. thousands of whole-

322     genome samples) within less than a day, and small ones (e.g. hundreds of whole-exome

323     samples) in seconds. As genome sequencing continues the march toward ever-larger datasets

324     and becomes more frequently used in diverse research settings, Bystro's combination of

325     performance and ease of use will prove invaluable for reproducible, rapid research.

326

327     **Methods**

328

329     **Accessing Bystro**

330        For most users, we recommend the Bystro web application (*https://bystro.io*), as it gives

331     full functionality, supports arbitrarily large datasets, and provides a convenient interface to the

332     natural-language search engine. Users with computational experience can download the Bystro

333     open-source package (*https://github.com/akotlar/bystro*). Using the provided installation script or

334     Amazon AMI image, Bystro can be easily deployed on an individual computer, computational

335     cluster, or any Amazon Web Services (AWS) EC2 instance. Bystro has very low memory and

336     CPU requirements, but benefits from fast SSD drives. As such we recommend at AWS

337     instances with provisioned I/O EBS drives, RAID 0 non-provisioned EBS, or i2/i3-class EC2

338     instances.

19

339

340     Detailed documentation on Bystro's use, as well as example search queries can be

341     found at https://bystro.io/help.

342

**Bystro comparisons with ANNOVAR, wANNOVAR, VEP, and GEMINI/Galaxy**

344

**Bystro Database**

346     Bystro databases were created using the open-source package

347     (https://github.com/akotlar/bystro). The hg19 and hg38 databases contains RefSeq, dbSNP,

348     PhyloP, PhastCons, Combined Annotation-Dependent Depletion (CADD), and Clinvar fields, as

349     well as custom annotations (Additional file 8). A complete listing of the original source data is

350     enumerated in the Git repository (https://github.com/akotlar/bystro/tree/master/config). Other

351     organism databases contain a subset of these sources, based on availability. Pre-built, up-to-

352     date versions of these databases are publicly available (https://github.com/akotlar/bystro).

353

**WGS Datasets**

355     Phase 1 and Phase 3 autosome and chromosome X VCF files were downloaded from

356     http://www.internationalgenome.org/data/. Phase 1 files were concatenated using bcftools[19]

357     "concat" function. Phase 3 files were concatenated using a custom Perl script

358     (https://github.com/wingolab-org/GenPro/blob/master/bin/mergeSnpFiles). The Phase 1 VCF file

359     was 895GB (139GB compressed), and the Phase 3 data was 853GB (15.6GB compressed).

360     The larger size of Phase 1 can be attributed to the inclusion of extra genotype information (the

361     genotype likelihood). The full Phase 3 chromosome 1 VCF file ($6.4 \times 10^6$ variants, 1.2GB

362     compressed), and $5 \times 10^4$-$4 \times 10^6$ variant allele subsets (8-655MB compressed) were also tested.

363     All Phase 1 and Phase 3 data correspond to the GRCh37/hg19 human genome assembly. All

364     data used are available (Additional file 9).

365

**Online annotation comparisons**

367    For online comparisons, the latest online versions offered at time of writing were used.

368    Bystro beta10 (September 2017), wANNOVAR (April 2017), VEP (April 2017), and GEMINI

369    (Galaxy version 0.8.1, released February 2016, latest as of October 2017) were tested online

370    with the full 1000 Genomes Phase 1 and Phase 3 VCF files, unless they were unable to upload

371    the files due to file size restrictions (Additional file 2). Bystro was found to be the only program

372    capable of uploading and processing the full Phase 1 and Phase 3 data sets, or subsets of

373    Phase 3 larger than $1 \times 10^6$ variants.

374

375    To conduct Bystro online annotations, a new user was registered within the public Bystro

376    web application (*https://bystro.io/*). Phase 1 and Phase 3 files were submitted in triplicate, one

377    replicate at a time, using the default database configuration (Additional file 2). Indexing was

378    automatically performed by Bystro upon completion of each annotation. The Phase 3 annotation

379    is publicly available to be tested (*https://bistro.io/public*).

380

381    The public Bystro server was configured on an Amazon i3.2xlarge EC2 instance. The

382    server supported 8 simultaneous users. Throughout the duration of each experiment, multiple

383    users had concurrent access to this server, increasing experiment variance, and limiting

384    observed performance.

385

386    Online Variant Effect Predictor (VEP) submissions were done using the VEP web

387    application (*http://www.ensembl.org/info/docs/tools/vep/index.html*). VEP has a 50MB

388    (compressed) file size limit. Due to gateway timeout issues and this file size limit, data sets

389    larger than $5 \times 10^4$ variants failed to complete (Additional file 2).

390

21

391    Online ANNOVAR submissions were handled using the wANNOVAR web application.

392    wANNOVAR could not accept the smallest tested file, the $5 \times 10^4$ variant subset of Phase 3

393    chromosome 1 (8MB compressed) due to file size restrictions (Additional file 2).

394    Galaxy submission was made using the public Galaxy servers. Galaxy provides

395    ANNOVAR, but its version of this software failed to complete any annotations, with the error

396    "unknown option: vcfinput". Annotations on Galaxy were therefore performed using GEMINI,

397    which provides annotations similar to Bystro's. Galaxy has a total storage allocation of 250GB

398    (after requisite decompression), and both Phase 1 and Phase 3 exceed this size. Galaxy was

399    therefore tested with the full $6.4 \times 10^6$ variant Phase 3 chromosome 1 VCF file. Galaxy's FTP

400    server was able to upload the file; however, Galaxy was unable to load the data into GEMINI,

401    terminating after running for 36 hours, with the message "This job was terminated because it ran

402    longer than the maximum allowed job run time" (Additional file 2). Subsets of Phase 3

403    chromosome 1 containing $5 \times 10^4$, $3 \times 10^5$, and $1 \times 10^6$ variants were therefore tested. Three

404    repetitions of the $5 \times 10^4$ variant submission were made. In consideration of the duration of

405    execution, two repetitions were made of the $3 \times 10^5$ and $1 \times 10^6$ variants submissions. Since

406    Galaxy does not record completion time, QuickTime was used to record each submission.

407

408    Bystro, VEP, and GEMINI online annotation times included the time to generate both a user-

409    readable tab-delimited text annotation and a searchable database. GEMINI required an extra

410    step to do so, using the query SELECT * FROM variants JOIN variant_impacts ON

411    variants.name = variant_impacts.name.

412    **Variant filtering comparisons**

413    After Bystro completed each annotation, it automatically indexed the results for search.

414    The time taken to index this data was recorded. Once this was completed, the Bystro web

415    application's search bar was used to filter the annotated sequencing experiments. The query

416    time, as well as the number of results and the transition to transversion ratio for each query,

22

417    were automatically generated by the search engine and recorded. Query time did not take into

418    account network latency between the search server and the web server. All queries were run six

419    times and averaged. The public search engine, which processed all queries, was hosted on a

420    single Amazon i3.2xlarge EC2 instance.

421

422        Since VEP, wANNOVAR, and Galaxy/GEMINI could not complete Phase 1 or Phase 3

423    annotations, variant filtering on these data sets could not be attempted. For small experiments

424    VEP and GEMINI can filter based on exact matches, while wANNOVAR provides only pre-

425    configured phenotype and disease model filters. VEP could annotate and filter at most only

426    $5x10^4$ variants and was therefore excluded from query comparisons.

427        Galaxy/GEMINI was tested with subsets of 1000 Genomes Phase 3 of $1x10^6$ variants

428    (the largest tested data set that Galaxy could handle), with the described settings (Additional file

429    2). In all GEMINI queries a JOIN operation on the variant_impacts table was used to return all

430    variant consequences, and all affected transcripts, as Bystro does by default. Similarly, Bystro's

431    CADD query was restricted to single nucleotide polymorphisms (using alt:(A || C || T || G)), as its

432    behavior diverges from GEMINI's at insertions and deletions: Bystro returns all possible CADD

433    Phred scores at such sites, whereas GEMINI returns a missing value. Bystro returns all values

434    to give users added flexibility: its search engine can accurately search within arrays (lists) of

435    data. Furthermore, as GEMINI on Galaxy only provided the Ensembl transcript set, for all query

436    comparisons with GEMINI, Bystro was configured to use Ensembl 90, which was the latest

437    version available at time of revision. It is important to note that the latest version of GEMINI on

438    Galaxy (0.8.1) dates to February 2016, and its databases are several years older: CADD (v1.0,

439    2014), Ensembl (v75, February 2014), ExAc (v0.3, October 2014), whereas Bystro uses up-to-

440    date resources. As a result of searching more up to date Ensembl (v90), population allele

441    frequency (gnomAD 2.0.1, the successor to ExAc 1.0), and CADD (v1.3) data, Bystro's queries

442    returned more data.

23

443    Since Galaxy does not report run times, QuickTime software was used to record each

444    run, and the query time was calculated as the difference between the time the search

445    submission entered the Galaxy queue, to the time that it was marked completed.

446    Galaxy/GEMINI queries were each run more than 6 times. Because run times varied by more

447    than 17x, the fastest consecutive 6 runs were averaged to minimize the influence of Galaxy

448    server load.

449

450    All comparisons with the Bystro search engine are limited, because no other existing

451    method provides natural-language parsing, and either rely on built-in scripts or require the user

452    to learn a specific language (SQL).

453

454    **Filtering accuracy comparison**

455    The latest version of Bystro (beta 10, September 2017) was used. For the 1000

456    Genomes query accuracy checks, the same underlying Ensembl-based Bystro annotation and

457    search index was used as in the Bystro/GEMINI filtering comparison. Direct comparison to

458    GEMINI were not made, in reflection of the age of the latest GEMINI Galaxy version (v0.8.1,

459    with database sources dating to 2014). All Bystro queries from that comparison were saved,

460    downloaded, and compared with Bystro "Filters", which are exact-match alternatives to Bystro's

461    natural-language queries, as well as custom Perl filtering scripts that also require exact

462    matches. A second query accuracy step was conducted, on the Yen et al 2017[9] VCF file. This

463    file was annotated using the standard RefSeq Bystro database. The same queries used in the

464    Bystro/GEMINI comparison were re-created on this smaller annotation, saved, downloaded, and

465    compared with Bystro "Filters" and Excel filters. Excel filters were created in Excel 2016 (Mac),

466    and required exact matches. All Excel-filtered and all Bystro query results were manually

467    inspected for concordance (Additional file 7). All scripts generated and used in the comparison

468    may be found at *https://github.com/akotlar/bystro-paper.*

469

470     **Offline annotation comparisons**

471         To generate offline performance data, the latest versions of each program available at

472     time of writing were used. Bystro beta10 (September 2017), VEP 86 (March 2017), and

473     ANNOVAR (March 2017) were each run on separate, dedicated Amazon i3.2xlarge EC2

474     instances (Additional file 3). All programs' databases were updated to the latest versions

475     available as of March 2017 (VEP, ANNOVAR), or September 2017 (Bystro). All programs were

476     configured to use the RefSeq transcript set.

477

478         Each instance contained 4 CPU cores (8 threads), 60GB RAM, and a 1920GB NVMe

479     SSD. Each instance was identically configured. All programs were configured to as closely

480     match Bystro's output as possible, although Bystro output more total annotation fields

481     (Additional file 3). Each data set tested was run 3 times. The annotation time for each run was

482     recorded, and averaged to generate the mean variant per second (variant/s) performance.

483     Submissions were recorded using the terminal recorder asciinema, and both memory and cpu

484     usage were recorded using the **free** and **top** commands set to a 30 second timeout.

485

486         VEP was configured to use 8 threads and to run in "offline" mode to maximize

487     performance, as recommended[3]. In each of three recorded trials, VEP was set to annotate

488     from RefSeq and CADD, and to check the reference assembly (Additional file 3). Based on

489     VEP's observed performance, adding PhastCons annotations was not attempted. VEP's

490     performance was measured by reading the program's log, which records variant/second

491     performance every $5x10^3$ annotated sites. In consideration of time, VEP was stopped after at

492     least $2x10^5$ variants were completed, and the $2x10^5$ variants performance was recorded.

493

25

494       ANNOVAR was configured to annotate RefSeq, CADD, PhastCons 100way, PhyloP

495     100way, Clinvar, avSNP, and ExAc version 0.3 (Additional file 3). ANNOVAR's avSNP database

496     was used in place of dbSNP, as recommended. We configured ANNOVAR to report allele

497     frequencies from ExAc, because it does not do so from either avSNP or dbSNP databases.

498     When annotating Phase 1, Phase 3, or Phase 3 chromosome 1, ANNOVAR crashed by

499     exceeding the available 60GB of memory. It was therefore tested with the subsets of Phase 3

500     chromosome 1 that contained $1 \times 10^6 - 4 \times 10^6$ variants.

501

502       Bystro was configured to annotate descriptions from RefSeq, dbSNP 147, CADD,

503     PhastCons 100way, PhyloP 100way, Clinvar, and to check the reference for each submitted

504     genomic position (Additional file 3).

505

506     **Annotation accuracy comparison**

507       The latest version of Bystro (beta 10, September 2017), ANNOVAR (July 2017), and

508     VEP (version 90) at the time of revision submission were used. All programs' databases were

509     updated to the latest version available. RefSeq-based databases were downloaded using each

510     program's database builder. All programs were compared on the Yen et al 2017 VCF file [9] for

511     position, variant call, and variant effects, based on each programs' respective RefSeq database.

512     The Yen et al VCF file *fileformat* header line was modified to "VCFv4.1" to allow programs to

513     recognize it as a valid VCF file. This modified file is available: https://github.com/akotlar/bystro-

514     paper. For the SnpEff comparison, annotations were adapted from Additional File 1 of Yen et al

515     2017[9]. ANNOVAR was additionally configured with gnomAD genomes, gnomAD exomes, and

516     CADD 1.3, and compared to Bystro on the corresponding values.

517

518     **Additional Files**

519    Additional file 1: This file contains 1) a feature comparison of tested programs, 2) investigation

520    of annotation concordance between tested programs, 3) investigation of Bystro query accuracy

521    (.docx, 1.4MB)

522    Additional file 2: Description of online comparison settings (.xlsx, 859KB)

523    Additional file 3: Description of online comparison settings (.xlsx, 40KB)

524    Additional file 4: Bystro vs ANNOVAR annotation comparison details (.xslx, 87KB)

525    Additional file 5: Bystro vs VEP annotation comparison details (.xslx, 701KB)

526    Additional file 6: Bystro vs SnpEff annotation comparison details (.xslx, 63KB)

527    Additional file 7: Bystro queries vs Excel filters concordance details (.xslx, 166KB)

528    Additional file 8: Species supported at time of writing, and their configurations (.xslx, 36KB)

529    Additional file 9: URLs of 1000 Genomes Phase 1, 1000 Genomes Phase 3, and Yen et al 2017

530    VCF files used (.xslx, 47KB)

531

532    **Declarations**

533    ***Availability of data and materials***

534    The Bystro web application is freely accessible at *https://bystro.io/*, and features detailed

535    interface documentation (*https://bystro.io/help*). The Bystro annotator, search indexer,

536    distributed queue servers, and database builder source code is freely available on GitHub

537    (*https://github.com/akotlar/bystro)* and Zenodo (doi: *10.5281/zenodo.1012417*), under the

538    Apache 2 open-source license [20]. The software is written in Perl and Go programming

539    languages and runs on Linux and Mac operating systems. Detailed documentation for Bystro

540    software is provided at *https://github.com/akotlar/bystro/blob/master/README.md*. The datasets

541    generated during and/or analyzed during the current study are available in the GitHub

542    repository, *https://github.com/akotlar/bystro-paper* [6],[9].

543

544    ***Author contributions***

545    A.V.K designed, wrote, and tested Bystro and performed experiments. C.E.T wrote

546    Bystro documentation and performed quality control. M.E.Z and D.J.C. contributed to the design

547    of Bystro and experiments. T.S.W. designed and wrote Bystro and designed and performed

548    experiments. A.V.K. and T.S.W. wrote the manuscript with contributions from all authors.

549

550    ***Acknowledgements***

551    We thank Kelly Shaw and Katherine Squires for beta testing and design suggestions. We thank

552    Viren Patel and the Emory Integrated Genomics Core (EIGC) for technical support.

553

554    ***Funding***

555    This work was supported by the AWS Cloud Credits for Research program**,** the Molecules to

556    Mankind program (a project of the Burroughs Wellcome Fund and the Laney Graduate School

557    at Emory University), Veterans Health Administration (BX001820), and the National Institutes of

558    Health (AG025688, MH101720, NS091859).

559

560    ***Competing interests***

561    The authors have no competing interests to declare.

562

563    ***Ethics approval and consent to participate***

564    Not applicable
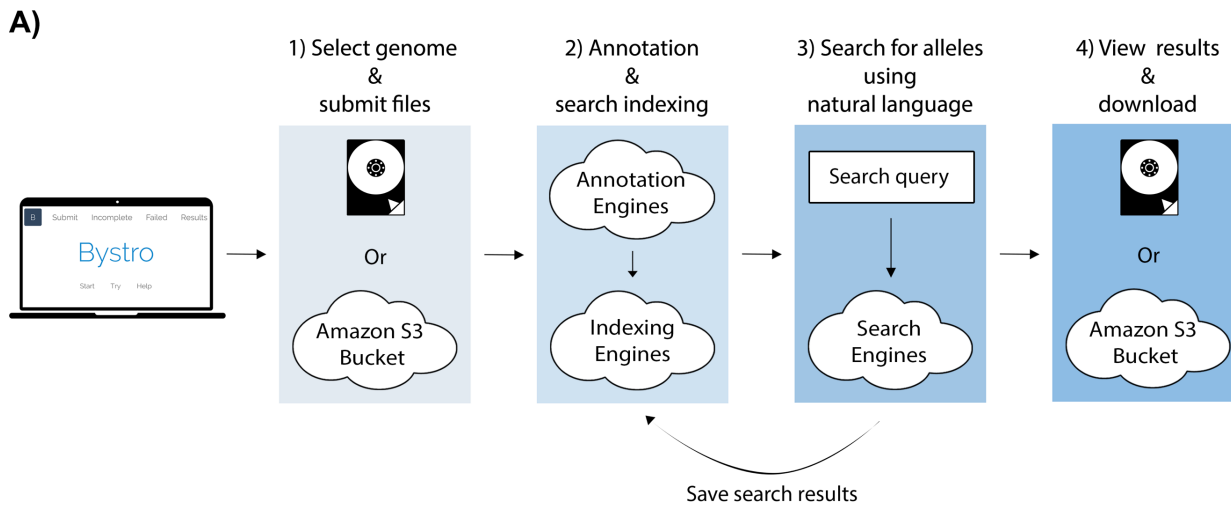
565

566    **References**

567    1.    Wang K, Li M, Hakonarson H: **ANNOVAR: functional annotation of genetic variants
568          from high-throughput sequencing data.** *Nucleic Acids Res* 2010, **38:**e164.

569   2.   Shetty AC, Athri P, Mondal K, Horner VL, Steinberg KM, Patel V, Caspary T, Cutler DJ,

570        Zwick ME: **SeqAnt: a web service to rapidly identify and annotate DNA sequence**

571        **variations.** *BMC Bioinformatics* 2010, **11:**471.

572   3.   McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, Flicek P, Cunningham F:

573        **The Ensembl Variant Effect Predictor.** *Genome Biol* 2016, **17:**122.

574   4.   DeFreitas T, Saddiki H, Flaherty P: **GEMINI: a computationally-efficient search**

575        **engine for large gene expression datasets.** *BMC Bioinformatics* 2016, **17:**102.

576   5.   Sandve GK, Nekrutenko A, Taylor J, Hovig E: **Ten simple rules for reproducible**

577        **computational research.** *PLoS Comput Biol* 2013, **9:**e1003285.

578   6.   Genomes Project C, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel

579        JO, Marchini JL, McCarthy S, McVean GA, Abecasis GR: **A global reference for**

580        **human genetic variation.** *Nature* 2015, **526:**68-74.

581        *http://dx.doi.org/10.1038/nature15393*

582   7.   Chang X, Wang K: **wANNOVAR: annotating genetic variants for personal genomes**

583        **via the web.** *J Med Genet* 2012, **49:**433-436.

584   8.   Goecks J, Nekrutenko A, Taylor J, Galaxy T: **Galaxy: a comprehensive approach for**

585        **supporting accessible, reproducible, and transparent computational research in**

586        **the life sciences.** *Genome Biol* 2010, **11:**R86.

587   9.   Yen JL, Garcia S, Montana A, Harris J, Chervitz S, Morra M, West J, Chen R, Church

588        DM: **A variant by any name: quantifying annotation discordance across tools and**

589        **clinical databases.** *Genome Med* 2017, **9:**7. *http://dx.doi.org/10.1186/s13073-016-*

590        *0396-7*

591   10.  Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE,

592        Lunter G, Marth GT, Sherry ST, et al: **The variant call format and VCFtools.**

593        *Bioinformatics* 2011, **27:**2156-2158.

594    11.    Johnston HR, Chopra P, Wingo TS, Patel V, International Consortium on B, Behavior in

595           22q11.2 Deletion S, Epstein MP, Mulle JG, Warren ST, Zwick ME, Cutler DJ: **PEMapper**

596           **and PECaller provide a simplified approach to whole-genome sequencing.** *Proc*

597           *Natl Acad Sci U S A* 2017.

598    12.    O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, Rajput B,

599           Robbertse B, Smith-White B, Ako-Adjei D, et al: **Reference sequence (RefSeq)**

600           **database at NCBI: current status, taxonomic expansion, and functional**

601           **annotation.** *Nucleic Acids Res* 2016, **44:**D733-745.

602    13.    Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K: **dbSNP:**

603           **the NCBI database of genetic variation.** *Nucleic Acids Res* 2001, **29:**308-311.

604    14.    Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A: **Detection of nonneutral**

605           **substitution rates on mammalian phylogenies.** *Genome Res* 2010, **20:**110-121.

606    15.    Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, Gu B, Hart J,

607           Hoffman D, Hoover J, et al: **ClinVar: public archive of interpretations of clinically**

608           **relevant variants.** *Nucleic Acids Res* 2016, **44:**D862-868.

609    16.    Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, O'Donnell-Luria

610           AH, Ware JS, Hill AJ, Cummings BB, et al: **Analysis of protein-coding genetic**

611           **variation in 60,706 humans.** *Nature* 2016, **536:**285-291.

612    17.    Ganna A, Genovese G, Howrigan DP, Byrnes A, Kurki MI, Zekavat SM, Whelan CW,

613           Kals M, Nivard MG, Bloemendal A, et al: **Ultra-rare disruptive and damaging**

614           **mutations influence educational attainment in the general population.** *Nat Neurosci*

615           2016, **19:**1563-1565.

616    18.    Taylor RC: **An overview of the Hadoop/MapReduce/HBase framework and its**

617           **current applications in bioinformatics.** *BMC Bioinformatics* 2010, **11 Suppl 12:**S1.

618   19.   Li H: **A statistical framework for SNP calling, mutation discovery, association**

619         **mapping and population genetical parameter estimation from sequencing data.**

620         *Bioinformatics* 2011, **27:**2987-2993.

621   20.   Kotlar A, Trevino C, Zwick M, Cutler D.J, Wingo T.S. **Bystro: rapid online variant**

622         **annotation and natural-language filtering at whole-genome scale**. Zenodo. 2017.

623         *http://dx.doi.org/10.5281/zenodo.834960*

624

**Figure 1**

**A)**



**B)**

**Figure 2**



1000 Genomes Phase 3 Online Processing Time