

1 **Title page**

2 **Title**

3 VAS: A convenient Web portal for efficient integration of genomic features with
4 millions of genetic variants

5 **Authors**

6 Eric Dun Ho ¹ (dho@cse.cuhk.edu.hk),

7 Qin Cao ¹ (qcao@cse.cuhk.edu.hk),

8 Sau Dan Lee ¹ (sdlee@cse.cuhk.edu.hk),

9 Kevin Y. Yip ^{1,2,3*} (kevinyip@cse.cuhk.edu.hk)

10 **Addresses**

11 ¹Department of Computer Science and Engineering,

12 ²Hong Kong Bioinformatics Centre,

13 ³CUHK-BGI Innovation Institute of Trans-omics,

14 The Chinese University of Hong Kong, Shatin, New Territories, Hong Kong

15 * Corresponding author

VAS: A convenient Web portal for efficient integration of genomic features with millions of genetic variants

Eric Dun Ho¹, Qin Cao¹, Sau Dan Lee¹ and Kevin Y. Yip^{1,2,3*}

Correspondence:
evinyip@cse.cuhk.edu.hk
Department of Computer Science
and Engineering, The Chinese
University of Hong Kong, Shatin,
New Territories, Hong Kong
Full list of author information is
available at the end of the article

Abstract

Background: High-throughput experimental methods have fostered the systematic detection of millions of genetic variants from any human genome. To help explore the potential biological implications of these genetic variants, software tools have been previously developed for integrating various types of information about these genomic regions from multiple data sources. Most of these tools were designed either for studying a small number of variants at a time, or for local execution on powerful machines.

Results: To make exploration of whole lists of genetic variants simple and accessible, we have developed a new Web-based system called VAS (Variant Annotation System, available at <https://yip1ab.cse.cuhk.edu.hk/vas/>). It provides a large variety of information useful for studying both coding and non-coding variants, including whole-genome transcription factor binding, open chromatin and transcription data from the ENCODE consortium. By means of data compression, millions of variants can be uploaded from a client machine to the server in less than 50 megabytes of data. On the server side, our customized data integration algorithms can efficiently link millions of variants with tens of whole-genome datasets. These two enabling technologies make VAS a practical tool for annotating genetic variants from large genomic studies. We demonstrate the use of VAS in annotating genetic variants obtained from a migraine meta-analysis study and multiple data sets from the Personal Genomes Project. We also compare the running time of annotating 6.4 million SNPs of the CEU trio by VAS and another tool, showing that VAS is efficient in handling new variant lists without requiring any pre-computations.

Conclusions: VAS is specially designed to handle annotation tasks with long lists of genetic variants and large numbers of annotating features efficiently. It is complementary to other existing tools with more specific aims such as evaluating the potential impacts of genetic variants in terms of disease risk. We recommend using VAS for a quick first-pass identification of potentially interesting genetic variants, to minimize the time required for other more in-depth downstream analyses.

Keywords: Annotation; Genetic Variants; Genomic Studies; Data Integration

18 **Background**

19 High-density microarrays and massively parallel sequencing have made genome-wide
20 detection of genetic variants from human DNA samples systematic, efficient and
21 inexpensive. In these experiments, it is common to observe hundreds of thousands
22 or even millions of loci in the DNA of a studied sample that differ from the reference
23 genome. To explore possible links between these variants and the phenotypes of the
24 sample, it is necessary to first analyze the potential biological significance of each
25 variant.

26 Early-days analysis methods have focused on the potential impacts of genetic
27 variants in coding regions, the functional consequences of which are usually related
28 to alterations to the corresponding proteins. There have been many successful soft-
29 ware tools for classifying coding variants into those that are synonymous, missense
30 and nonsense, whether they may affect splicing or cause frameshift, and the level
31 of disruption to protein functions and structures [1, 9, 22, 25, 29, 33].

32 On the other hand, it is now well-recognized that many functionally important
33 genetic variants do not change the coding sequences directly but rather perturb
34 gene regulation [11, 13]. For example, a single nucleotide variant (SNV) may hit
35 the binding motif of a transcription factor, which affects the proper binding of it
36 and leads to an expression level change of the regulated gene. Since currently there
37 is not a complete catalog of regulatory regions in the human genome, in order
38 to determine how likely a genetic variant may affect gene regulation, one needs to
39 examine many types of static and cell/tissue-specific features indicative of functional
40 significance. Static features such as evolutionary conservation and sequence motifs
41 help evaluate the possibility for a genomic region to ever play a functional role,
42 while cell/tissue-specific features provide information about regulatory activities in
43 each genomic region in particular cell/tissue types and conditions. Combining both
44 types of features provides a quick and low-cost way to pinpoint the potentially
45 most interesting variants for downstream validation and functional studies. For
46 example, DNase I hypersensitivity and certain histone marks together could identify
47 regulatory regions active in particular cell types that are far away from their target
48 genes [18], while integrating such information with sequence motifs could further
49 predict the transcription factors involved in the gene regulation.

50 A large amount of data containing cell/tissue-specific features have been pro-
51 duced for various human cell types in large-scale studies such as ENCODE [13]
52 and Roadmap Epigenomics [5]. To utilize these data in studying genetic variants, a
53 number of Web tools have been developed for automatic large-scale genomic data
54 integration [3, 6, 7, 16, 20, 21, 23, 26, 28, 34]. Each of them provides a database of
55 genomic features collected from multiple data sources, and a procedure for users to
56 query selected features around their genetic variants. These tools face two common
57 challenges, namely 1) A list of genetic variants in standard Variant Call Format
58 (VCF) could take up hundreds of megabytes and need a long time to upload; and
59 2) Integrating a long variant list with a large number of whole-genome features is
60 time-consuming.

61 Concerning the data uploading issue, some tools restrict the maximum number of
62 genetic variants per job to a small value, while others do not set an explicit limit but
63 practically cannot handle full lists of millions of variants [3, 6, 7, 23]. Some other
64 tools avoid the uploading of large files by allowing local installation and execution,
65 which requires a large amount of genomic features to be downloaded to the user
66 machine [26].

67 Regarding the data integration issue, most tools use a relational database to store
68 the collected data. As a result, a table join between a stored feature and the uploaded
69 genetic variants is often performed by time-inefficient algorithms that make use of
70 standard tree-based indices. Although more efficient linear-time sort-merge join
71 algorithms are available, it could be difficult to instruct the query optimizer to use
72 them. Some tools attempted to solve this problem by pre-computing the results of
73 a large amount of table joins [10, 26], which requires extra disk space for storing the
74 pre-computed results and new pre-computation needs to be performed every time
75 a new genomic feature is added to the database.

76 To overcome these two issues, we have developed VAS (Variant Annotation Sys-
77 tem), a tool for efficient genomic data integration.

78 **Implementation**

79 The overall workflow of VAS is shown in Figure 1. Below we describe its different
80 components in detail.

81 Genomic Features in VAS

82 VAS provides a large variety of genomic features collected from different data sources
83 (Table 1). To help explore genetic variants in non-coding regions, it provides a
84 rich set of whole-genome features about sequence patterns, conservation, chromatin
85 states and expression signals from various experimental and computational data
86 sets. Cell/tissue-specific data based on different cell types studied by the ENCODE
87 Project Consortium and Roadmap Epigenomics are provided for some features. Ad-
88 ditional features are provided for referencing previous findings about known vari-
89 ants and their loci, including previously cataloged SNPs, information about disease
90 SNPs, and Gencode gene annotation, which contains a large number of non-coding
91 RNAs.

92 Feature selection, data compression and data integration

93 A user uploads a list of genetic variants and selects the features to be integrated
94 through a user-friendly Web interface. Multiple data formats are supported for
95 the input list of genetic variants, including VCF and white-space-delimited lists.
96 In our test, uploading 3 million genetic variants involved less than 50 megabytes
97 of data transfer (Figure 2). The enabling technology behind this small uploading
98 data size is a compression procedure that VAS performs on the client side. In a
99 standard VCF file, there is a lot of information not required for the data integration
100 purpose. Our Flash plugin takes the user-supplied variant file, retains only genomic
101 locations, and removes repetitive text such as chromosome names. The resulting file
102 contains compact arrays of chromosomal locations, one for each chromosome. This
103 compression process is transparent to the user in that a user only needs to specify
104 a standard genetic variant file as input and the compression will be automatically
105 performed before the compressed data is transferred to the server.

106 The genomic features to be integrated with the genetic variants are selected from
107 a Web interface that provides a list of the features available. Functions are also pro-
108 vided for searching for particular datasets using their attributes such as cell type
109 (Figure 3a). For each genetic variant, VAS can search for genomic features overlap-
110 ping its exact location or a flanking window of it up to 1Mb, allowing exploration
111 of nearby loci in genetic linkage to the input variants.

112 Upon submitting the input variants and the selected genomic features, the data
113 integration job is added to a queue on the server side. The back-end system adopts
114 a scalable design that allows executing multiple jobs on different computing nodes
115 in parallel. The user is redirected to a waiting page that provides the latest status
116 of the job. Optionally, if an email address is entered, an email notification will be
117 sent to the user when the job is finished.

118 We store data in a customized file format without relying on a relational database,
119 which facilitated our design of linear-time integration algorithms that can efficiently
120 identify overlapping genomic regions in different data files. Specifically, for each
121 feature, the genomic regions containing feature values are sorted according to their
122 genomic locations. Special pointers are added to particular locations (such as the
123 start of each chromosome) in the genome to allow direct access of these locations
124 without a sequential scan of all regions from the beginning of the file.

125 We provide two types of data integration. The first one is identifying genomic fea-
126 tures overlapping exactly the locations of the input genetic variants (exact location
127 for an SNV or insertion, mid-point for a deletion). The second one is identifying
128 genomic features overlapping a flanking window of each input genetic variant. Both
129 types of integration are performed by sort-merge algorithms.

130 For the first type of data integration, we first sort the input variants according
131 to their locations. We then use a pointer to scan through all the genetic variants
132 and all the genomic feature regions sequentially. Whenever a region of the genomic
133 feature is encountered, we add it to a feature queue. Any genetic variant that is
134 then encountered before the end of the region will be annotated with the region
135 and the result is stored in the variant map (see Figure 4 for an example). More
136 specifically, during the scanning process, the algorithm takes one of the following
137 actions whenever a point of the corresponding type is encountered:

- 138 • Location of a variant: Annotate the variant with all the regions currently in
139 the feature queue and store the results in the variant map
- 140 • Starting position of a feature region: Add the region to the feature queue
- 141 • Ending position of a feature region: Remove the region from the feature queue

142 For the second type of data integration, the integration algorithm is similar to
143 the one for the first type, except that now instead of considering a single location
144 of each genetic variant, we consider the starting and ending positions of its flanking

145 window. During the scanning process, the algorithm takes one of the following
146 actions whenever a point of the corresponding type is encountered (see Figure 5 for
147 an example):

- 148 • Starting position of the flanking window of a variant: Add the variant to
149 the variant queue, annotate the variant with all the regions currently in the
150 feature queue and store the results in the variant map
- 151 • Ending position of the flanking window of a variant: Remove the variant from
152 the variant queue
- 153 • Starting position of a feature region: Add the region to the feature queue,
154 annotate all variants currently in the variant queue with the region and store
155 the results in the variant map
- 156 • Ending position of a feature region: Remove the region from the feature queue

157 We have compared the speed efficiency of these data integration algorithms with
158 some alternative methods. For all the methods, we tried to integrate a list of 57,902
159 variants with a genomic feature with 17,524 regions. We tested both types of data
160 integration, with the size of the flanking window set to 100bp in the second type of
161 integration. The time needed for the different methods to perform the integration
162 task is shown in Table 2. Our customized algorithms were found to be the most
163 efficient among the methods in comparison.

164 When the data integration is finished, the results are displayed on a Web page
165 that shows information about the selected features around each input variant (Fig-
166 ure 3b). In the case of numeric features, the average feature values around each
167 variant and their percentiles among all genomic regions are also shown. Details
168 of the features can be displayed in a signal-track image generated by the UCSC
169 Genome Browser (Figure 3c). Linking to a corresponding UCSC Genome Browser
170 session is provided for more visualization options and interactive explorations. Inte-
171 gration results can also be downloaded in Microsoft Excel or tab-delimited formats
172 for further analyses.

173 Each data integration job is given a unique 512-bit identifier. The user who issues
174 a job can browse and download the results at a later time by using the provided
175 hyperlink with this identifier embedded. All job files are kept on the server for 30
176 days. Other users without this identifier are unable to access the uploaded data or
177 the corresponding data integration results.

178 Currently there are several related tools providing genome-wide annotation of
179 genetic variants. Each of these tools has its unique features and advantages. We list
180 in Table 3 some of the distinctive properties of VAS.

181 **Results and discussion**

182 **Case studies**

183 As a demonstration of using VAS in exploring potential functional meanings of
184 genetic variants, we used it to analyze two sets of genetic variants with different
185 sets of genomic features.

186 The first set of genetic variants includes the susceptibility loci for migraine identi-
187 fied in a recent study [2]. In that study, a genome-wide meta-analysis was performed
188 on the data from 29 genome-wide association studies, which together involved 23,285
189 individuals with migraine and 95,425 population-matched controls. Twelve loci were
190 identified to be significantly associated with migraine, while 5 loci were found to
191 have significant expression quantitative trait loci (eQTL). We used VAS to retrieve
192 information about various types of static and cell-specific data around these 17 loci.
193 For static features, we considered evolutionary conservation, known variants in db-
194 SNP and GWAS Catalog, protein binding motifs and CpG islands. For cell-specific
195 features, we considered histone modifications, open chromatin and transcription
196 factor binding data from ENCODE sequencing experiments for both normal brain
197 and spinal cord cells (HAc, HA-h, HA-sp and NH-A) and brain cancer lines (BE2_C,
198 Gliobla, Medullo and SK-N-SH_RA).

199 Figure 3b shows part of the annotation results, where the darkness of a table en-
200 try indicates how strong the signal value is. It can be seen that many features have
201 strong signals around the susceptibility loci. As an example, Figure 3c shows the de-
202 tailed view of rs12134493 (marked by the red line), which is at position 115,479,469
203 (hg18)/ 115,677,946 (hg19) of chromosome 1. It is located in an intergenic region
204 downstream of and close to the TSPAN2 gene. In the original study [2], it was found
205 that the susceptibility loci in general had strong open chromatin signals in terms of
206 DNase I hypersensitivity, and they overlapped with some transcription factor bind-
207 ing motifs. Consistent with their findings, VAS was able to find overlaps between
208 the SNP and open chromatin signals in various normal brain cells (Figure 3c i, ii)

209 and the presence of binding motifs for multiple transcription factors around that
210 region (Figure 3c iii).

211 We also made a number of additional interesting observations based on the VAS
212 results. First, the open chromatin signals were found only in normal brain cells but
213 not in the cancer line SK-N-SH_RA. Second, in astrocytes (NH-A), the SNP over-
214 lapped a local region with strong H3K27ac signals (Figure 3c iv), which suggests
215 that the region could be an active enhancer in this cell type. Third, the SNP was
216 inside a region with strong evolutionary conservation among placental mammals
217 and among vertebrates (Figure 3c v), suggesting that the region is under evolu-
218 tionary constraints. Finally, there was active binding of CTCF, RAD21 and YY1
219 in a nearby region a few kilobases away (Figure 3c vi) with corresponding open
220 chromatin signals. Given the closeness of this region and the susceptibility locus, it
221 may be useful to include this region into the study.

222 The second set of genetic variants comes from the Personal Genome Project [8]
223 (<https://my.pgp-hms.org/>). We randomly downloaded 5 lists of genetic variants
224 with at least one variant reported to have high clinical importance according to the
225 report on the Web site (Table 4). We tested if we could identify these variants of
226 potential clinical importance using VAS, by annotating them with the information
227 from GWAS Catalog [35] and the Human Gene Mutation Database [32]. On average,
228 uploading and completing the annotation of each data file took less than 10 minutes.
229 VAS was able to annotate all 21 unique variants reported to be likely pathogenic
230 and rare pathogenic using the information from the two databases, which confirms
231 that VAS can be used to quickly integrate information from diverse sources for more
232 in-depth downstream analyses.

233 Data uploading and integration time

234 To test the speed performance of VAS in handling large data files, we recorded
235 the time required to integrate 6.4 million genetic variants present in the CEU trio
236 obtained from the 1000 Genomes Project with the information of the whole list of
237 SNPs in dbSNP. We compared the performance of VAS with both the reported re-
238 sults and our local execution of GEMINI [26], a tool that allows large-scale genomic
239 data integration by means of local execution and pre-caching of table join results.

240 Both VAS and our local execution of GEMINI were tested on a machine with dual
241 quad core Xeon CPU at 2.4GHz and 64GB of main memory.

242 The resulting time measurements of the two tools (Table 5) show that VAS finished
243 the data integration in around half an hour. As for GEMINI, although our time
244 measurements are different from those reported in the original paper due to the use
245 of different machines, in general a long data loading time (1.5-3 hours) was required
246 for the extensive pre-computation, followed by a very quick data integration phase.
247 This pre-computation step needs to be performed whenever a new set of genetic
248 variants is to be annotated.

249 Since GEMINI was executed locally while VAS is an online system, there was
250 extra data uploading time for VAS. For the data set tested, the data uploading time
251 was negligible as compared to the time needed for data integration. This result is
252 consistent with our above analysis on file size and data uploading time at different
253 numbers of input genetic variants (Figure 2).

254 Overall, VAS is more efficient and flexible in handling new variant lists since it
255 does not require pre-loading of data, while GEMINI works better in situations where
256 the same list of genetic variants is to be repeatedly analyzed by integrating with
257 many different subsets of genomic data.

258 **Conclusion**

259 In this paper, we have described VAS, a new Web tool that can efficiently integrate
260 millions of genetic variants with tens of whole-genome data sets in a single inte-
261 gration task. The client-side data compression procedure and the customized data
262 store allowed fast uploading and integrating whole lists of genetic variants obtained
263 from genomic studies, making VAS a practical tool for routine first-step annotation
264 of genetic variants.

265 When analyzing large-scale genomic data, the main bottleneck is usually inspect-
266 ing long lists of results, pinpointing the most biologically or medically significant
267 parts, and making correct interpretations of them. The time spent on data inte-
268 gration is usually relatively unimportant. However, the time difference between a
269 standard data integration method and a customized one could become large when
270 the numbers of input genetic variants and integrating genomic features are large. In
271 addition, since VAS can accept multiple job requests from different users simultane-

272 ously, having an efficient data integration method can also shorten the time spent
273 on waiting for other earlier jobs in the queue to complete.

274 Currently VAS supports job-level parallelization, which means multiple jobs can
275 be run at the same time in parallel on different computing units. In the future, we
276 plan to extend VAS to support sub-job-level parallelization, which means a single
277 job can be divided into sub-tasks simultaneously performed on different computing
278 units. As the integration of each genetic variant is independent of the other variants,
279 high-level distributed computing frameworks such as MapReduce should be readily
280 applicable. An additional advantage of adopting such a framework is the distribution
281 of data to multiple machines, which allows for better scalability.

282 VAS is currently implemented as an online system, which enjoys the advantage
283 of requiring no local installation or downloading of genomic features by the user.
284 We have ensured data integrity and confidentiality by providing encrypted network
285 connections and assigning task IDs that are only made known to the users who
286 submit the tasks. However, there are situations in which some private data can
287 only be analyzed locally. Theoretically a user can install a local version of VAS on
288 his/her own machine to perform the analysis offline, but that would also require
289 downloading a large amount of stored data features. We will investigate ways to
290 facilitate data integration in these situations, such as allowing users to easily down-
291 load a selected subset of features or dynamically download data features at the time
292 needed, and developing privacy-preserving distributed data integration algorithms.

293 In the case study we have demonstrated that with the data currently loaded into
294 VAS, one could already use it to obtain some interesting patterns around each
295 genetic variant. As more and more cell/tissue-specific data are being produced,
296 we will keep updating the data repository of VAS to cover more cell/tissue types
297 and more data for each cell/tissue type. We also plan on supporting the GRCh38
298 human reference genome when most data files in our database have a CRCh38
299 version available.

300 **Availability and requirements**

301 **Project name:** Variant Annotation System (VAS)

302 **Project home page:** <https://yiplab.cse.cuhk.edu.hk/vas/>

303 **Operating system:** VAS can be accessed from any platform by using one of the
304 listed Web browsers

305 **Programming languages:** PHP, Python

306 **Other requirements:** We recommend accessing VAS by using Google Chrome
307 (version 35 or higher), Microsoft Internet Explorer (version 10 or higher), or Mozilla
308 Firefox (version 24 or higher), with JavaScript enabled and a minimum screen res-
309 olution of 1024 pixels x 768 pixels

310 **Any restrictions to use by non-academics:** Nil

311 **Competing interests**

312 The authors declare that they have no competing interests.

313 **Author's contributions**

314 KYY conceived the study. EDH, SDL and KYY designed the system. EDH and QC collected the data and
315 implemented the system. EDH, QC, SDL and KYY tested the system. EDH and KYY wrote the manuscript.

316 **Acknowledgements**

317 SDL is partially supported by the HKRGC Theme-based Research Scheme T12-401/13-R. KYY is partially
318 supported by the HKRGC Early Career Scheme 419612.

319 **Author details**

320 ¹Department of Computer Science and Engineering, The Chinese University of Hong Kong, Shatin, New Territories,
321 Hong Kong. ²Hong Kong Bioinformatics Centre, The Chinese University of Hong Kong, Shatin, New Territories,
322 Hong Kong. ³CUHK-BGI Innovation Institute of Trans-omics, The Chinese University of Hong Kong, Shatin, New
323 Territories, Hong Kong.

324 **References**

- 325 1. Ivan A Adzhubei, Steffen Schmidt, Leonid Peshkin, Vasily E Ramensky, Anna Gerasimova, Peer Bork, Alexey S
326 Kondrashov, and Shamil R Sunyaev. A method and server for predicting damaging missense mutations. *Nature*
327 *Methods*, 7(4):248–249, 2010.
- 328 2. Verner Anttila, Bendik S Winsvold, Padhraig Gormley, Tobias Kurth, Francesco Bettella, George McMahon,
329 Mikko Kallela, Rainer Malik, Boukje de Vries, Gisela Terwindt, Sarah E Medland, Unda Todt, Wendy L
330 McArdle, Lydia Quaye, Markku Koiranen, M Arfan Ikram, Terho Lehtimäki, Anine H Stam, Lannie Ligthart,
331 Juho Wedenoja, Ian Dunham, Benjamin M Neale, Priit Palta, Eija Hamalainen, Markus Schurks, Lynda M
332 Rose, Julie E Buring, Paul M Ridker, Stacy Steinberg, Hreinn Stefansson, Finnbogi Jakobsson, Debbie A
333 Lawlor, David M Evans, Susan M Ring, Markus Färkkilä, Ville Artto, Mari A Kaunisto, Tobias Freilinger, Jean
334 Schoenen, Rune R Frants, Nadine Pelzer, Claudia M Weller, Ronald Zielman, Andrew C Heath, Pamela A F
335 Madden, Grant W Montgomery, Nicholas G Martin, Guntram Borck, Hartmut Göbel, Axel Heinze, Katja
336 Heinze-Kuhn, Frances M K Williams, Anna-Liisa Hartikainen, Anneli Pouta, Joyce van den Ende, Andre G
337 Uitterlinden, Albert Hofman, Najaf Amin, Jouke-Jan Hottenga, Jacqueline M Vink, Kauko Heikkilä, Michael
338 Alexander, Bertram Muller-Myhsok, Stefan Schreiber, Thomas Meitinger, Heinz Erich Wichmann, Arpo
339 Aromaa, Johan G Eriksson, Bryan J Traynor, Daniah Trabzuni, "North American Brain Expression
340 Consortium", "UK Brain Expression Consortium", Elizabeth Rossin, Kasper Lage, Suzanne B R Jacobs,
341 J Raphael Gibbs, Ewan Birney, Jaakko Kaprio, Brenda W Penninx, Dorret I Boomsma, Cornelia van Duijn, Olli
342 Raitakari, Marjo-Riitta Jarvelin, John-Anker Zwart, Lynn Cherkas, David P Strachan, Christian Kubisch,
343 Michel D Ferrari, Arn M J M van den Maagdenberg, Martin Dichgans, Maija Wessman, George Davey Smith,
344 Kari Stefansson, Mark J Daly, Dale R Nyholt, Daniel I Chasman, and Aarno Palotie. Genome-wide
345 meta-analysis identifies new susceptibility loci for migraine. *Nature Genetics*, 45(8):912–917, 2013.

- 346 3. Maxim Barrenboim and Thomas Manke. ChroMoS: An integrated web tool for SNP classification, prioritization
347 and functional interpretation. *Bioinformatics*, 29(17):2197–2198, 2013.
- 348 4. Gary Benson. Tandem repeats finder: A program to analyze DNA sequences. *Nucleic Acids Research*,
349 27:573–580, 1999.
- 350 5. Bradley E Bernstein, John A Stamatoyannopoulos, Joseph F Costello, Bing Ren, Aleksandar Milosavljevic,
351 Alexander Meissner, Manolis Kellis, Marco A Marra, Arthur L Beaudet, Joseph R Ecker, Peggy J Farnham,
352 Martin Hirst, Eric S Lander, Tarjei S Mikkelsen, and James A Thomson. The NIH roadmap epigenomics
353 mapping consortium. *Nature Biotechnology*, 28(10):1045–1048, 2010.
- 354 6. Alan P Boyle, Eurie L Hong, Manoj Hariharan, Yong Cheng, Marc A Schaub, Maya Kasowski, Konrad J
355 Karczewski, Julie Park, Benjamin C Hitz, Shuai Weng, J Michael Cherry, and Michael Snyder. Annotation of
356 functional variation in personal genomes using RegulomeDB. *Genome Research*, 22:1790–1797, 2012.
- 357 7. Yu-Chang Cheng, Fang-Chih Hsiao, Erh-Chan Yeh, Wan-Jia Lin, Cheng-Yang Louis Tang, Huan-Chin Tseng,
358 Hsing-Tsung Wu, Chuan-Kun Liu, Chih-Cheng Chen, Yuan-Tsong Chen, and Adam Yao. VarioWatch:
359 Providing large-scale and comprehensive annotations on human genomic variants in the next generation
360 sequencing era. *Nucleic Acids Research*, 40:W76–W81, 2012.
- 361 8. G M Church. The personal genome project. *Molecular Systems Biology*, 1(2005.0030), 2005.
- 362 9. Pablo Cingolani, Adrian Platts, Le Lily Wang, Melissa Coon, Tung Nguyen, Luan Wang, Susan J Land,
363 Douglas M Ruden, and Xiangyi Lu. A program for annotating and predicting the effects of single nucleotide
364 polymorphisms, SnpEff: SNPs in the genome of drosophila melanogaster strain w1118; iso-2; iso-3. *Fly*,
365 6(2):80–92, 2012.
- 366 10. Sergio Contrino, Richard N. Smith, Daniela Butano, Adrian Carr, Fengyuan Hu, Rachel Lyne, Kim Rutherford,
367 Alex Kalderimis, Julie Sullivan, Seth Carbon, Ellen T. Kephart, Paul Lloyd, E. O. Stinson, Nicole L.
368 Washington, Marc D. Perry, Peter Ruzanov, Zheng Zha, Suzanna E. Lewis, Lincoln D. Stein, and Gos Micklem.
369 modMine: Flexible access to modENCODE data. *Nucleic Acids Research*, 40:D1082–D1088, 2012.
- 370 11. Gregory M. Cooper and Jay Shendure. Needles in stacks of needles: Finding disease-causal variants in a wealth
371 of genomic data. *Nature Reviews Genetics*, 12(9):628–640, 2011.
- 372 12. Thomas Derrien, Jordi Estellé, Santiago Marco Sola, David G. Knowles, Emanuele Raineri, Roderic Guigó, and
373 Paolo Ribeca. Fast computation and applications of genome mappability. *PLOS ONE*, 7(e30377), 2012.
- 374 13. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature*,
375 489(7414):57–74, 2012.
- 376 14. Jason Ernst and Manolis Kellis. ChromHMM: Automating chromatin-state discovery and characterization.
377 *Nature Methods*, 9(3):215–216, 2012.
- 378 15. Paul Flicek, M. Ridwan Amode, Daniel Barrell, Kathryn Beal, Konstantinos Billis, Simon Brent, Denise
379 Carvalho-Silva, Peter Clapham, Guy Coates, Stephen Fitzgerald, Laurent Gil, Carlos Garcia Giron, Leo Gordon,
380 Thibaut Hourlier, Nathan Hunt, Sarah Johnson, Thomas Juettemann, Andreas K. Kahari, Stephen Keenan,
381 Eugene Kulesha, Fergal J. Martin, Thomas Maurel, William M. McLaren, Daniel N. Murphy, Rishi Nag, Bert
382 Overduin, Miguel Pignatelli, Bethan Pritchard, Emily Pritchard, Harpreet S. Riat, Magali Ruffier, Daniel
383 Sheppard, Kieron Taylor, Anja Thormann, Stephen J. Trevanion, Alessandro Vullo, Steven P. Wilder, Mark
384 Wilson, Amonida Zadissa, Bronwen L. Aken, Ewan Birney, Fiona Cunningham, Jennifer Harrow, Javier Herrero,
385 Tim J.P. Hubbard, Rhoda Kinsella, Matthieu Muffato, Anne Parker, Giulietta Spudich, Andy Yates, Daniel R.
386 Zerbino, and Stephen M.J. Searle. Ensembl 2014. *Nucleic Acids Research*, 42:D749–D755, 2014.
- 387 16. Lukas Habegger, Suganthi Balasubramanian, David Z. Chen, Ekta Khurana, Andrea Sboner, Arif Harmanci,
388 Joel Rozowsky, Declan Clarke, Michael Snyder, and Mark Gerstein. VAT: A computational framework to
389 functionally annotate variants in personal genomes within a cloud-computing environment. *Bioinformatics*,
390 28(17):2267–2269, 2010.
- 391 17. Jennifer Harrow, Adam Frankish, Jose M. Gonzalez, Electra Tapanari, Mark Diekhans, Felix Kokocinski,
392 Bronwen L. Aken, Daniel Barrell, Amonida Zadissa, Stephen Searle, If Barnes, Alexandra Bignell, Veronika
393 Boychenko, Toby Hunt, Mike Kay, Gaurab Mukherjee, Jeena Rajan, Gloria Despacio-Reyes, Gary Saunders,
394 Charles Steward, Rachel Harte, Michael Lin, Cédric Howald, Andrea Tanzer, Thomas Derrien, Jacqueline
395 Chrast, Nathalie Walters, Suganthi Balasubramanian, Baikang Pei, Michael Tress, Jose Manuel Rodriguez,
396 Iakes Ezkurdia, Jeltje van Baren, Michael Brent, David Haussler, Manolis Kellis, Alfonso Valencia, Alexandre

- 397 Reymond, Mark Gerstein, Roderic Guigó, and Tim J. Hubbard. GENCODE: The reference human genome
398 annotation for the ENCODE project. *Genome Research*, 22:1760–1774, 2012.
- 399 18. Nathaniel D. Heintzman, Gary C. Hon, R. David Hawkins, Pouya Kheradpour, Alexander Stark, Lindsey F.
400 Harp, Zhen Ye, Leonard K. Lee, Rhona K. Stuart, Christina W. Ching, Keith A. Ching, Jessica E.
401 Antosiewicz-Bourget, Hui Liu, Xinmin Zhang, Roland D. Green, Victor V. Lobanenkov, Ron Stewart, James A.
402 Thomson, Gregory E. Crawford, Manolis Kellis, and Bing Ren. Histone modifications at human enhancers
403 reflect global cell-type-specific gene expression. *Nature*, 459(7243):108–112, 2009.
- 404 19. Donna Karolchik, Galt P. Barber, Jonathan Casper, Hiram Clawson, Melissa S. Cline, Mark Diekhans,
405 Timothy R. Dreszer, Pauline A. Fujita, Luvina Guruvadoo, Maximilian Haeussler, Rachel A. Harte, Steve
406 Heitner, Angie S. Hinrichs, Katrina Learned, Brian T. Lee, Chin H. Li, Brian J. Raney, Brooke Rhead, Kate R.
407 Rosenbloom, Cricket A. Sloan, Matthew L. Speir, Ann S. Zweig, David Haussler, Robert M. Kuhn, and
408 W. James Kent. The UCSC genome browser database: 2014 update. *Nucleic Acids Research*, 42:764–770, 2014.
- 409 20. Ekta Khurana, Yao Fu, Vincenza Colonna, Xinmeng Jasmine Mu, Hyun Min Kang, Tuuli Lappalainen, Andrea
410 Sboner, Lucas Lochovsky, Jieming Chen, Arif Harmanci, Jishnu Das, Alexej Abyzov, Suganthi
411 Balasubramanian, Kathryn Beal, Dimple Chakravarty, Daniel Challis, Yuan Chen, Declan Clarke, Laura Clarke,
412 Fiona Cunningham, Uday S. Evani, Paul Flicek, Robert Fragoza, Erik Garrison, Richard Gibbs, Zeynep H.
413 Gümüş, Javier Herrero, Naoki Kitabayashi, Yong Kong, Kasper Lage, Vaja Liluashvili, Steven M. Lipkin,
414 Daniel G. MacArthur, Gabor Marth, Donna Muzny, Tune H. Pers, Graham R. S. Ritchie, Jeffrey A. Rosenfeld,
415 Cristina Sisu, Xiaomu Wei, Michael Wilson, Yali Xue, Fuli Yu, "1000 Genomes Project Consortium",
416 Emmanouil T. Dermitzakis, Haiyuan Yu, Mark A. Rubin, Chris Tyler-Smith, and Mark Gerstein. Integrative
417 annotation of variants from 1092 humans: Application to cancer genomics. *Science*, 342(6154):1235587, 2013.
- 418 21. Martin Kircher, Daniela M Witten, Preti Jain, Brian J R'Roak, Gregory M Cooper, and Jay Shendure. A
419 general framework for estimating the relative pathogenicity of human genetic variants. *Nature Genetics*,
420 46(3):310–315, 2014.
- 421 22. Prateek Kumar, Steven Henikoff, and Ng Pauline C. Predicting the effects of coding non-synonymous variants
422 on protein function using the SIFT algorithm. *Nature Protocols*, 4(7):1073–1081, 2009.
- 423 23. Mulin Jun Li, Panwen Wang, Xiaorong Liu, Ee Lyn Lim, Zhangyong Wang, Meredith Yeager, Maria P Wong,
424 Pak Chung Sham, Stephen J Chanock, and Junwen Wang. GWASdb: A database for human genetic variants
425 identified by genome-wide association studies. *Nucleic Acids Research*, 40:D1047–D1054, 2011.
- 426 24. V. Matys, E. Fricke, R. Geffers, E. Göbbling, M. Haubrock, R. Hehl, K. Hornischer, D. Karas, A. E. Kel, O. V.
427 Kel-Margoulis, D.-U. Kloos, S. Land, B. Lewicki-Potapov, H. Michael, R. Münch, I. Reuter, S. Rotert,
428 H. Saxel, M. Scheer, S. Thiele, and E. Wingender. TRANSFAC: Transcriptional regulation, from patterns to
429 profiles. *Nucleic Acids Research*, 31:374–378, 2003.
- 430 25. William McLaren, Bethan Pritchard, Daniel Rios, Yuan Chen, Paul Flicek, and Fiona Cunningham. Deriving the
431 consequences of genomic variants with the ensembl API and SNP effect predictor. *Bioinformatics*,
432 26(16):2069–2070, 2010.
- 433 26. Umadevi Paila, Brad A Chapman, Rory Kirchner, and Aaron R Quinlan. GEMINI: Integrative exploration of
434 genetic variation and genome annotations. *PLOS Computational Biology*, 9(e1003153), 2013.
- 435 27. Katherine S. Pollard, Melissa J. Hubisz, Kate R. Rosenbloom, and Adam Siepel. Detection of nonneutral
436 substitution rates on mammalian phylogenies. *Genome Research*, 20:110–121, 2010.
- 437 28. Graham R S Ritchie, Ian Dunham, Eleftheria Zeggini, and Paul Flicek. Functional annotation of noncoding
438 sequence variants. *Nature Methods*, 11(3):294–296, 2014.
- 439 29. Christian Schaefer, Alice Meier, Burkhard Rost, and Yana Bromberg. SNPdbe: Constructing and nsSNP
440 functional impacts database. *Bioinformatics*, 28(4):601–602, 2011.
- 441 30. S. T. Sherry, M.-H. Ward, J. Baker, Kholodov, L. Phan, E.M. Smigielski, and K. Sirotkin. dbSNP: the NCBI
442 database of genetic variation. *Nucleic Acids Research*, 29(1):308–311, 2001.
- 443 31. Adam Siepel, Gill Bejerano, Jakob S. Pedersen, Angie S. Hinrichs, Minmei Hou, Kate Rosenbloom, Hiram
444 Clawson, John Spieth, LaDeana W. Hillier, Stephen Richards, George M. Weinstock, Richard K. Wilson,
445 Richard A. Gibbs, W. James Kent, Webb Miller, and David Haussler. Evolutionarily conserved elements in
446 vertebrate, insect, worm, and yeast genomes. *Genome Research*, 15:1034–1050, 2005.
- 447 32. Peter D. Stenson, Matthew Mort, Edward V. Ball, Katy Shaw, Andrew D. Phillips, and David N. Cooper. The

- 448 human gene mutation database: Building a comprehensive mutation repository for clinical and molecular
 449 genetics, diagnostic testing and personalized genomic medicine. *Human Genetics*, 133:1–9, 2014.
- 450 33. Kai Wang, Mingyao Li, and Hakon Hakonarson. ANNOVAR: Functional annotation of genetic variants from
 451 high-throughput sequencing data. *Nucleic Acids Research*, 38:e164, 2010.
- 452 34. Lucas D. Ward and Manolis Kellis. HaploReg: A resource for exploring chromatin states, conservation, and
 453 regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Research*, 40:D930–D934,
 454 2012.
- 455 35. Danielle Welter, Jacqueline MacArthur, Joannella Morales, Tony Burdett, Peggy Hall, Heather Junkins, Alan
 456 Klemm, Paul Flicek, Teri Manolio, Lucia Hindorf, and Helen Parkinson. The NHGRI GWAS catalog, a curated
 457 resource of SNP-trait associations. *Nucleic Acids Research*, 42:D1001–D1006, 2014.
- 458 36. Kevin Y Yip, Chao Cheng, Nitin Bhardwaj, James B Brown, Jing Leng, Anshul Kundaje, Joel Rozowsky, Ewan
 459 Birney, Peter Bickel, Michael Snyder, and Mark Gerstein. Classification of human genomic regions based on
 460 experimentally determined binding sites of more than 100 transcription-related factors. *Genome Biology*,
 461 13(R48), 2012.

462 Figures

Figure 1 Schematic illustration of the VAS workflow. Genomic features are pre-sorted and stored in data files with pointers for direct access to particular genomic locations. A user supplies the list of genetic variants and selects the genomic features to integrate with the variants at the client side. The variants extractor produces a compressed form of the input variants. The task is then sent to the backend and put into a waiting queue, and the user is shown a waiting page. When an execution daemon becomes available, it fetches the next task in the queue and uses the customized algorithms to perform data integration. The integration results are stored in a tab-delimited file. The user will then be shown a summary page of the integration results. An email notification will also be sent, with a link for a user to retrieve the summary page later. The user can then view the integration details of each input variant, perform interactive analysis on the UCSC Genome Browser, or download the annotation results in tab-delimited or Excel format.

Figure 2 Amount of data upload and uploading time required at various sizes of the input list of genetic variants in our simulation study, before and after client-side data compression. The data uploading time for the uncompressed case was estimated based on the file size and the data transfer rate when transferring the compressed version of the same files.

Figure 3 Usage of VAS. (a) Selecting genomic features to be integrated with the genetic variants. (b) Summary of the annotation results. Genomic features identified around each genetic variant (within a 10kb window in this case) are shown, where a darker color indicates a stronger signal value. (c) Detailed view of a genetic variant, with an embedded UCSC Genome Browser image in which each genomic feature is shown as a signal track.

Figure 4 An example of point-to-region data integration using our algorithm.

Figure 5 An example of region-to-region data integration using our algorithm.

463 Tables

Table 1 List of genomic features provided by VAS

Type	Genomic features
Chromatin	ENCODE open chromatin, histone modifications, protein-DNA binding [13], Roadmap Epigenomics DNA methylation [5]
Genomic states	ChromHMM segmentation [14], supervised genomic region classification [36]
Expression	ENCODE RNA-seq [13]
Sequence	UCSC [19] conservation scores [31, 27], transcription factor binding motifs [24], sequence uniqueness [12], repeats[4], GC content
Annotation	Gencode [17]
Variations	dbSNP [30]
Diseases	GWAS Catalog [35], The Human Gene Mutation Database [32]

Table 2 Data integration time of different methods. For BigBed reader and interval tree, we used the implementation of bxpython. For relational database, we tried several indexing methods including standard B-tree index and spatial index, and report here the shortest time among these approaches. Tabix was called using the pytabix library in Python.

Method	Integrating variant locations (second)	Integration variant flanking windows (second)
BigBed	277.90	275.63
Interval tree	0.41	0.60
Relational database	8.05	736.23
Tabix	8.87	8.88
Our algorithms	0.21	0.52

Table 3 Some distinctive features of VAS as compared to some related tools. For GWAVA and RegulomeDB, the maximum number of input variants allowed is based on our own tests of the system. Properties of the tools are based on their versions on 8th September 2014.

Tool	CADD [21]	GEMINI [26]	GWASdb [23]	GWAVA [28]	HaploReg [34]	RegulomeDB [6]	VAS
Client-side data compression	No	(local)	N/A	No	No	No	Yes
Input variants allowed	~100,000	(Unlimited)	1	>10,000	10,000	~5,000	3,000,000
Genomic features/aggregated features provided	63	(User defined)	37	14	10	1,012	1,000+
			(5 categories)		(6 categories)	(13 categories)	(16 categories)
Data storage and integration	(Not described)	Relational DB	Relational DB	(Not described)	Relational DB	Relational DB	Customized
Searching flanking regions	No	No	Yes	No	No	No	Yes
Asynchronous access of results	Yes	(local)	No	No	No	No	Yes
Linkout to genome browser	No	No	UCSC [19]	Ensembl [15]	No	UCSC	UCSC

Table 4 Lists of genetic variants from the Personal Genome Project tested on VAS. The variants listed in the “PGP variants” column include likely pathogenic and rare (<2.5%) pathogenic variants according to the reports available on the Personal Genome Project Web site. The information in the “Chromosomal location”, “dbSNP ID” and “Clinical importance” columns was all obtained from these reports.

Sample	Total number of variants	PGP variants	Chromosomal location	dbSNP ID	Clinical importance	Found by VAS
hu47A9D1	960,613	APOA5-S19W	chr11:116662407/chr11:116167616	rs3135506	Low	Yes
		APOE-C130R	chr19:45411941/chr19:50103780	rs429358	High	Yes
		MBL2-G54D	chr10:54531235/chr10:54201240	rs1800450	Low	Yes
		MBL2-R52C	chr10:54531242/chr10:54201247	rs5030737	Low	Yes
		MTRR-I49M	chr5:7870973/chr5:7923972	rs1801394	Low	Yes
		MYO7A-R302H	chr11:76869378/chr11:76547025	rs41298135	High	Yes
		rs5186	chr3:148459988/chr3:149942677	rs5186	Low	Yes
hu7DA960	960,613	AMPD1-Q12X	chr11:115236057/chr11:115037579	rs17602729	Low	Yes
		KCNE1-D85N	chr21:35821680/chr21:34743549	N/A	High	Yes
		KRT5-G138E	chr12:52913668/chr12:51199934	rs11170164	Low	Yes
		MBL2-G54D	chr10:54531235/chr10:54201240	rs1800450	Low	Yes
		rs5186	chr3:148459988/chr3:149942677	rs5186	Low	Yes
hu8D40D6	598,897	APOE-C130R	chr19:45411941/chr19:50103780	rs429358	High	Yes
		HFE-S65C	chr6:26091185	N/A	Low	Yes
		MTRR-I49M	chr5:7870973/chr5:7923972	rs1801394	Low	Yes
		PRPH-D141Y	chr12:49689404	rs58599399	High	Yes
		RPF1-A91V	chr10:72360387/chr10:72030392	rs35947132	Low	Yes
		SERPINA1-E288V	chr14:94847262/chr14:93917014	rs17580	Low	Yes
hu998A3D	960,613	BTD-D444H	chr3:15686693/chr3:15661696	rs13078881	Low	Yes
		C3-R102G	chr19:6718387/chr19:6669386	rs2230199	Moderate	Yes
		COL4A1-Q1334H	chr13:110818598/chr13:109616598	rs3742207	Low	Yes
		HFE-S65C	chr6:26091185	N/A	Low	Yes
		MTRR-I49M	chr5:7870973/chr5:7923972	rs1801394	Low	Yes
		rs5186	chr3:148459988/chr3:149942677	rs5186	Low	Yes
		SERPINA1-E366K	chr14:94844947/chr14:93914699	rs28929474	High	Yes
hgD53911	612,647	COL4A1-Q1334H	chr13:110818598/chr13:109616598	rs3742207	Low	Yes
		MTRR-I49M	chr5:7870973/chr5:7923972	rs1801394	Low	Yes
		PKD1-R4276W	chr16:2139814/chr16:2079814	rs114251396	High	Yes
		rs5186	chr3:148459988/chr3:149942677	rs5186	Low	Yes
		SCNN1G-E197K	chr16:23200963/chr16:23108463	rs5738	Low	Yes
		VWF-R854Q	chr12:6143978/chr12:6014238	rs41276738	Moderate	Yes

Table 5 Time measurement of GEMINI and VAS

Tool		Data loading/uploading (s)*	Data integration (s)	Total (s)
GEMINI (as reported in [26])	Average	5,050.0	24.0	5,064.0
GEMINI (our testing results)	Trial 1	9,944.6	154.1	10,098.6
	Trial 2	9,960.5	155.5	10,116.1
	Trial 3	10,182.4	156.9	10,339.3
	Trial 4	10,182.3	162.8	10,345.1
	Trial 5	10,053.2	169.1	10,222.2
	Average	10,064.6	159.7	10,224.3
	Std. dev.	115.2	6.2	117.6
VAS	Trial 1	9.9	1,711.1	1,721.1
	Trial 2	10.4	1,772.3	1,782.7
	Trial 3	9.7	1,552.5	1,562.1
	Trial 4	9.2	1,541.6	1,550.8
	Trial 5	9.6	1,580.9	1,590.5
	Average	9.8	1,631.7	1,641.4
	Std. dev.	0.4	103.7	104.1