

GigaScience

SV-plaudit: A cloud-based framework for manually curating thousands of structural variants

--Manuscript Draft--

Manuscript Number:	GIGA-D-18-00103R2	
Full Title:	SV-plaudit: A cloud-based framework for manually curating thousands of structural variants	
Article Type:	Research	
Funding Information:	National Human Genome Research Institute (K99HG009532)	Dr Ryan Layer
	National Human Genome Research Institute (R01HG006693)	Dr Aaron R Quinlan
	National Human Genome Research Institute (R01GM124355)	Dr Aaron R Quinlan
	National Cancer Institute (US) (U24CA209999)	Dr Aaron R Quinlan
Abstract:	SV-plaudit is a framework for rapidly curating structural variant (SVs) predictions. For each SV, we generate an image that visualizes the coverage and alignment signals from a set of samples. Images are uploaded to our cloud framework where users assess the quality of each image using a client-side web application. Reports can then be generated as a tab-delimited file or annotated VCF. As a proof of principle, nine researchers collaborated for one hour to evaluate 1,350 SVs each. We anticipate that SV-plaudit will become a standard step in variant calling pipelines and the crowd-sourced curation of other biological results.	
Corresponding Author:	Ryan Layer UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Jonathan R Belyeu	
First Author Secondary Information:		
Order of Authors:	Jonathan R Belyeu	
	Thomas J Nicholas, PHD	
	Brent S Pedersen, PHD	
	Thomas A Sasani	
	James M Havrilla	
	Stephanie N Kravitz	
	Megan E Conway	
	Brian K Lohman, PHD	
	Aaron R Quinlan, PHD	
	Ryan Layer	
Order of Authors Secondary Information:		

Response to Reviewers:	We added the citation for the data snapshot and RRID.
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends. Have you included all the information requested in your manuscript?	Yes
Resources A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible. Have you included the information requested as detailed in our Minimum Standards Reporting Checklist ?	Yes
Availability of data and materials All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript. Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist ?	Yes

1 **SV-plaudit: A cloud-based framework for manually curating thousands of structural**
2
3
4 **variants**

5
6
7
8
9 Jonathan R. Belyeu^{1,2}, Thomas J. Nicholas^{1,2}, Brent S. Pedersen^{1,2}, Thomas A. Sasani^{1,2}, James M. Havrilla^{1,2},
10
11 Stephanie N. Kravitz^{1,2}, Megan E. Conway¹, Brian K. Lohman^{1,2}, Aaron R. Quinlan^{1,2,3+}, Ryan M. Layer^{1,2+}
12
13
14

15 1. Department of Human Genetics, University of Utah, Salt Lake City, UT
16
17 2. USTAR Center for Genetic Discovery, University of Utah, Salt Lake City, UT
18
19 3. Department of Biomedical Informatics, University of Utah, Salt Lake City, UT
20
21
22

23
24 + *To whom correspondence should be addressed*
25
26
27
28

29 **ABSTRACT**
30
31

32 *SV-plaudit* is a framework for rapidly curating structural variant (SVs) predictions. For each SV, we generate an
33
34 image that visualizes the coverage and alignment signals from a set of samples. Images are uploaded to our
35
36 cloud framework where users assess the quality of each image using a client-side web application. Reports can
37
38 then be generated as a tab-delimited file or annotated VCF. As a proof of principle, nine researchers collaborated
39
40 for one hour to evaluate 1,350 SVs each. We anticipate that *SV-plaudit* will become a standard step in variant
41
42 calling pipelines and the crowd-sourced curation of other biological results.
43
44
45
46

47 Code available at <https://github.com/jbelyeu/SV-plaudit>
48
49

50 Demonstration video available at <https://www.youtube.com/watch?v=ono8kHMKxDs>
51
52
53

54 **KEYWORDS**
55
56

57 Structural variants; Visualization; Manual curation
58
59
60
61
62
63
64
65

1 BACKGROUND

2
3 Large genomic rearrangements, or structural variants (SVs), are an abundant form of genetic variation within the
4 human genome^{1,2}, and they play an important role in both species evolution^{3,4} and human disease phenotypes<sup>5-
5
6
7
8
9</sup>. While many methods have been developed to identify SVs from whole-genome sequencing (WGS) data¹⁰⁻¹⁴,
10 the accuracy of SV prediction remains far below that of single-nucleotide and insertion-deletion variants¹.
11
12 Improvements to SV detection algorithms have, in part, been limited by the availability and applicability of high-
13
14
15 quality truth sets. While the Genome in a Bottle¹⁵ consortium has made considerable progress toward a gold-
16
17 standard variant truth set, the incredibly high quality of the data underlying this project (300X and PCR-free) calls
18
19 into question the generality of the accuracy obtained in typical quality WGS datasets (30X with PCR-
20
21 amplification).
22
23
24
25

26 Given the high false positive rate of SV calls from genome and exome sequencing, manual inspection is a critical
27
28 quality control step, especially in clinical cases. Scrutiny of the evidence supporting an SV is considered to be a
29
30 reliable "dry bench" validation technique, as the human eye can rapidly distinguish true SV signal from alignment
31
32 artifacts. In principle, we could improve the accuracy of SV call sets by visually validating every variant. In
33
34 practice, however, current genomic data visualization methods¹⁶⁻²¹ were designed primarily for spot checking a
35
36 small number of variants and are difficult to scale to the thousands of SVs in typical call sets. Therefore, a curated
37
38 set of SVs requires a new framework that scales to thousands of SVs, minimizes the time needed to adjudicate
39
40 individual variants, and manages the collective judgment of large and often geographically dispersed teams.
41
42
43
44
45

46 Here we present *SV-plaudit*, a fast, highly-scalable framework enabling teams of any size to collaborate on the
47
48 rapid, web-based curation of thousands of SVs. In the web interface, users **consider a curation question for a**
49
50 **series of pre-computed images (Fig 1, Supplementary Fig 1) that contain the coverage, paired-end alignments,**
51
52 **and split-read alignments for the region surrounding a candidate SV for a set of relevant samples (e.g., tumor**
53
54 **and matched normal samples). The curation question is defined by the researcher to match the larger**
55
56 **experimental design (e.g., a cancer study may ask if the variant a somatic variant, a germline variant, or a false**
57
58 **positive). Responses are collected and returned as a report which can be used to identify high-quality variants.**
59
60
61
62
63
64
65

1
2
3 While a team of curators is not required, collecting multiple opinions for each SV allows *SV-plaudit* to report the
4
5 consensus view (i.e., a "curation score") of each variant. This consensus is less susceptible to human error and
6
7 does not require expert users to score variants. With *SV-plaudit*, it is practical to inspect and score every variant
8
9 in a call set, thereby improving the accuracy of SV predictions in individual genomes and allowing curation of
10
11 high quality-truth sets for SV method tuning.
12
13
14
15
16

17 RESULTS

18
19 To assess *SV-plaudit's* utility for curating SVs, nine researchers in the Quinlan laboratory at the University of
20
21 Utah manually inspected and scored the 1,350 SVs (1,310 deletions, 8 duplications, 4 insertions, and 28
22
23 inversions) that the 1000 Genomes Project¹ identified in the NA12878 genome (**Supplemental File 1**). Since we
24
25 expect trio analysis to be a common use case of *SV-plaudit*, we included alignments from NA12878 and her
26
27 parents (NA12891 and NA12892), and participants considered the curation question "The SV in the top sample
28
29 (NA12878) is:" and answers "GOOD", "BAD", or "DE NOVO". In total, the full experiment took less than two hours
30
31 with Amazon costs totaling less than \$0.05. The images (**Supplemental File 2**) were generated in 3 minutes (20
32
33 threads, 2.7 seconds per image) and uploading to S3 required 5 minutes (full command list in **Supplemental**
34
35 **File 3**). The mean time to score all images was 60.1 minutes (2.67 seconds per image) (**Fig 2A**, reports in
36
37 **Supplemental Files 4,5**). In the scoring process, no de novo variants were identified. 40 images did not render
38
39 correctly due to issues in the alignment files (e.g., coverage gaps) and were removed from the subsequent
40
41 analysis (**Supplemental File 6**).
42
43
44
45
46
47
48

49 For this experiment, we use a curation score that mapped "GOOD" and "DE NOVO" to the value one, "BAD" to
50
51 the value zero, and the mean as the aggregation function (**Fig 2B**). Most (70.5%) of variants were scored
52
53 unanimously, with 67.1% being unanimously "GOOD" (score = 1.0, e.g., **Fig 1A**) and 3.4% being unanimously
54
55 "BAD" (score = 0.0, e.g. **Fig 1B**). Since we had nine scores for each variant, we expanded our definition of
56
57 "unambiguous" variants to be those with at most one dissenting vote (score <0.2 or >0.8), which accounts for
58
59 87.1% of the variants. The 12.9% of SVs that were "ambiguous" (more than one dissenting vote, 0.2<= score
60
61
62
63
64
65

1 <=0.8) were generally small (median size of 310.5bp versus 899.5bp for all variants, **Fig 2C**) or contained
2
3 conflicting evidence (e.g., paired-end and split-read evidence indicated an inversion and the read-depth evidence
4
5 indicated a deletion, e.g., **Fig 1C.**)
6

7
8
9
10 Other methods, such as SVTYPER²² and CNVNATOR²³, can independently assess the validity of SV calls.
11
12 SVTYPER genotypes SVs for a given sample by comparing the number of discordant paired-end alignments
13
14 and split-read alignments that support the SV to the number of pairs and reads that support the reference allele.
15
16 CNVNATOR uses sequence coverage to estimate copy number for the region affected by the SV. Both of these
17
18 methods confirm the voting results (**Fig 2D**). Considering the set of “unambiguous” deletions, SVTYPER and
19
20 CNVNATOR agree with the *SV-plaudit* curation score in 92.3% and 81.7% of cases, respectively. Here,
21
22 agreement means that unambiguous false SVs (curation score < 0.2) have a CNVNATOR copy number near
23
24 two (between 1.4 and 2.4) or an SVTYPER genotype of homozygous reference. Unambiguous true SVs (curation
25
26 score > 0.8) have a CNVNATOR copy number near one or zero (less than 1.4), or an SVTYPER genotype of
27
28 non-reference (heterozygous or homozygous alternate).
29
30
31
32
33
34

35 Despite this consistency, using either SVTYPER or CNVNATOR to validate SVs can lead to false positives or
36
37 false negatives. For example, CNVNATOR reported a copy number loss for 44.2% of the deletions that were
38
39 scored as unanimously BAD, and SVTYPER called 30.7% of the deletions that were unanimously GOOD as
40
41 homozygous reference. Conversely, CNVNATOR had few false negatives (2.4% of unanimously GOOD
42
43 deletions were called as copy neutral), and SVTYPER had few false positives (0.2% of non-reference variants
44
45 were unanimously BAD). **This comparison is meant to demonstrate that different methods have distinct strengths
46
47 and weaknesses; and should not be taken as a direct comparison between SVTYPER and CNVNATOR, since
48
49 CNVNATOR was one of nine methods used by the 1000 Genomes project while SVTYPER was not.**
50
51
52
53
54

55 These results demonstrate that, with *SV-plaudit*, manual curation can be a cost-effective and robust part of the
56
57 SV detection process. While we anticipate that automated SV detection methods will continue to improve, due
58
59 in part to the improved truth sets that *SV-plaudit* will provide, directly viewing SVs will remain an essential
60
61
62
63
64
65

1 validation technique. By extending this validation to full call sets, *SV-plaudit* not only improves specificity but can
2
3 also enhance sensitivity by allowing users to relax quality filters and rapidly screen large sets of calls. Beyond
4
5 demonstrating *SV-plaudit's* utility, our curation of SVs for NA12878 is useful as a high-quality truth set for method
6
7 development and tuning. A VCF of these variants annotated with their curation score is available in
8
9

10 **Supplementary File 5.**

14 **DISCUSSION**

16
17 *SV-plaudit* is an efficient, scalable, and flexible framework for the manual curation of large-scale SV call sets.
18
19 Backed by Amazon S3 and DynamoDB, *SV-plaudit* is easy to deploy and scales to teams of any size. Each
20
21 instantiation of *SV-plaudit* is completely independent and can be deployed locally for private or sensitive
22
23 datasets, or be distributed publicly to maximize participation. By rapidly providing a direct view of the raw data
24
25 underlying candidate SVs, *SV-plaudit* delivers the infrastructure to manually inspect full SV call sets. *SV-plaudit*
26
27 also allows researchers to specify the questions and answers that users consider, to ensure that the curation
28
29 outcome supports the larger experimental design. This functionality is vital to a wide range of WGS experiments,
30
31 from method development to the interpretation of disease genomes. We are actively working on machine
32
33 learning methods that will leverage the curation scores for thousands of SV predictions as training data.
34
35
36
37
38

39 **CONCLUSIONS**

41
42 *SV-plaudit* was designed to judge how well the data in an alignment file corroborate a candidate SV. The question
43
44 of whether a particular SV is a false positive due to artifacts from sequencing or alignment is a broader issue
45
46 that must be answered in the context of other data sources such as mappability and repeat annotations. While
47
48 this second level of analysis is crucial, it is beyond the scope of this paper, and we argue this analysis be
49
50 performed only for those SVs that are fully supported by the alignment data. While *SV-plaudit* combines *samplot*
51
52 and *PlotCritic* to enable the curation of structural variant images, we emphasize that the *PlotCritic* framework
53
54 can be used to score images of any type. Therefore, we anticipate that this framework will facilitate "crowd-
55
56 sourced" curation of many other biological images.
57
58
59
60
61
62
63
64
65

METHODS

Overview. *SV-plaudit* (Fig 3) is based on two software packages: *samplot* for SV image generation, and *PlotCritic* for staging the Amazon cloud environment and managing user input. Once the environment is staged, users log into the system and are presented with a series of SV images in either a random or predetermined order. For each image, the user answers the curation question and responses are logged. Reports on the progress of a project can be quickly generated at any point in the process.

Samplot. *Samplot* is a Python program that uses *pysam*²⁴ to extract alignment data from a set of BAM or CRAM files, and *matplotlib*²⁵ to visualize the raw data for the genomic region surrounding a candidate SV (Fig 3A). For each alignment file, *samplot* renders the depth of sequencing coverage, paired-end alignments, and split-read alignments where paired-end and split-read alignments are color-coded based by the type of SV they support (e.g., black for deletion, red for a duplication, etc.) (Fig 1 Supplementary Figure 2, which considers variants at different sequencing coverages, and Supplementary Figure 3, which depicts variants supported by long-read sequencing).^{26,27} Alignments are positioned along the x-axis by genomic location and along the left y-axis by the distance between the ends (insert size), which helps users to differentiate normal alignments from discordant alignments that support an SV. Depth of sequencing coverage is also displayed on the right y-axis to allow users to inspect whether putative copy number changes are supported by the expected changes in coverage. To improve performance for large events, we downsample “normal” paired-end alignments (a +/- orientation and an insert size range that is within Z standard deviations from the mean; by default Z = 4). Plots for each alignment file are stacked and share a common x-axis that reports the chromosomal position. By convention, the sample of interest (e.g., proband or tumor) is displayed as the top track, followed by the set of related reference genomes tracks (e.g., parents and siblings, matched normal sample). Users may specify the exact order by using command line parameters to *samplot*. A visualization of genome annotations and genes and exons within the locus is displayed below the alignment plots to provide context for assessing the SV's relevance to phenotypes. Rendering time depends on the number of samples, *sequence coverage*, and the size of the SV, but most images will require less than 5 seconds, and *samplot* rendering can be parallelizable by SV call.

1 **PlotCritic.** *PlotCritic* (**Fig 3B**) provides a simple web interface for scoring images and viewing reports that
2
3 summarize the results from multiple users and SV images. *PlotCritic* is both highly scalable and easy to deploy.
4
5 Images are stored on Amazon Web Services (AWS) S3 and DynamoDB tables store project configuration
6
7 metadata and user responses. These AWS services allow *PlotCritic* to dynamically scale to any number of users.
8
9
10 It also precludes the need for hosting a dedicated server, thereby facilitating deployment.

11
12
13
14
15 After *samplot* generates the SV images, *PlotCritic* manages their transfer to S3 and configures tables in
16
17 DynamoDB based on a JSON configuration file (`config.json` file in **Fig 3B**). In this configuration file, one
18
19 defines the curation questions posed to reviewers, as well as the allowed answers and associated keyboard
20
21 bindings to allow faster responses (`curationQandA` field in **Fig 3B**). In turn, these dictate the text and buttons
22
23 that appear on the resulting web interface. As such, it allows the interface to be easily customized to support a
24
25 wide variety of curation scenarios. For example, a cancer experiment may display a tumor sample and matched
26
27 normal sample and ask users if the SV appears in both samples (i.e., a germline variant) or just in the tumor
28
29 sample (i.e., a somatic variant). To accomplish this, the curation question (`question` field in **Fig 3B**) could be
30
31 “In which samples does the SV appear?”, and the answer options (`answers` field in **Fig 3B**) could be “TUMOR”,
32
33 “BOTH”, “NORMAL”, “NEITHER”. Alternatively, in the case of a rare disease, the interface could display a
34
35 proband and parents and ask if the SV is only in the proband (i.e., de novo) or if it is also in a parent (i.e.,
36
37 inherited). Since there is no limit to the length of a questions or number of answers options, *PlotCritic* can support
38
39 more complex experimental scenarios.
40
41
42
43
44
45

46
47 Once results are collected, *PlotCritic* can generate a tab-delimited report or annotated VCF that, for each SV
48
49 image, details the number of times the image was scored and the full set of answers it received. Additionally, a
50
51 curation score can be calculated for each image by providing a value for each answer option and an aggregation
52
53 function (e.g., mean, median, mode, standard deviation, min, max). For example, consider the cancer example
54
55 from above where the values three, two, one, and zero mapped to the answers “TUMOR”, “BOTH”, “NORMAL”,
56
57 and “NEITHER”, respectively. If “mode” were selected as the curation function, then the curation score would
58
59 reflect the opinion of a plurality of users. The mean would reflect the consensus among all users, and the
60
61
62
63
64
65

1 standard deviation would capture the level of disagreement about each image. While we expect mean, median,
2
3 mode, standard deviation, min, and max to satisfy most use cases, users can implement custom scores by
4
5 operating on the tab-delimited reported.
6
7
8
9

10 Each *PlotCritic* project is protected by AWS Cognito user authentication, which securely restricts access to the
11
12 project website to authenticated users. A project manager is the only authorized user at startup and can
13
14 authenticate other users using Cognito's secure services. The website can be further secured using HTTPS and
15
16 additional controls, such as IP restrictions, can be put in place by configuring AWS IAM access controls directly
17
18 for S3 and DynamoDB.
19
20
21
22

23 **AVAILABILITY OF SOURCE CODE AND REQUIREMENTS**

24 Project name: SV-Plaudit

25
26 Project home page: <https://github.com/jbelyeu/SV-plaudit>

27
28 Operating systems: Mac OS and Linux

29
30 Programming language: Python, bash

31
32 License: MIT

33
34
35
36
37 **Research Resource Initiative Identification ID: SCR_016285**

38 39 40 41 **AVAILABILITY OF SUPPORTING DATA AND MATERIAL**

42
43 The datasets generated and/or analyzed during the current study are available in the 1000 Genomes Project
44
45 repository, <ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp/phase3/data/>

46
47
48
49
50
51 **All test data used or generated during this study, and a snapshot of the code, are available in the GigaScience**
52
53 **GigaDB repository.²⁸**

54
55
56
57
58
59
60
61
62
63
64
65

1 **DECLARATIONS**

2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

List of abbreviations

SV: Structural Variant

Ethics approval and consent to participate

Not applicable

Consent for publication

Not applicable

Competing interests

The authors declare that they have no competing interests.

Funding

This research was supported by a US National Human Genome Research Institute awards to RML (NIH K99HG009532) and ARQ (NIH R01HG006693 and NIH R01GM124355), as well as a US National Cancer Institute award to ARQ (NIH U24CA209999).

Authors' contributions

JRP and RML developed the software. JRB, TJN, BSP, TAS, JMH, SNK, MEC, BKL, and RML scored variants for the experiment. JRP, ARQ, and RML wrote the manuscript. ARQ and RML conceived the study.

1 REFERENCES

- 2
3
4 1. Sudmant, P. H. *et al.* An integrated map of structural variation in 2,504 human genomes. *Nature* **526**, 75–
5 81 (2015).
6
- 7
8 2. Redon, R. *et al.* Global variation in copy number in the human genome. *Nature* **444**, 444–454 (2006).
9
- 10 3. Newman, T. L. *et al.* A genome-wide survey of structural variation between human and chimpanzee.
11
12 *Genome Res.* **15**, 1344–1356 (2005).
13
- 14 4. Bailey, J. A. & Eichler, E. E. Primate segmental duplications: crucibles of evolution, diversity and disease.
15
16 *Nat. Rev. Genet.* **7**, 552–564 (2006).
17
- 18 5. Payer, L. M. *et al.* Structural variants caused by Alu insertions are associated with risks for many human
19
20 diseases. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E3984–E3992 (2017).
21
22
- 23 6. Schubert, C. The genomic basis of the Williams-Beuren syndrome. *Cell. Mol. Life Sci.* **66**, 1178–1197
24
25 (2009).
26
27
- 28 7. Pleasance, E. D. *et al.* A comprehensive catalogue of somatic mutations from a human cancer genome.
29
30 *Nature* **463**, 191–196 (2010).
31
32
- 33 8. Venkitaraman, A. R. Cancer susceptibility and the functions of BRCA1 and BRCA2. *Cell* **108**, 171–182
34
35 (2002).
36
37
- 38 9. Zhang, F., Gu, W., Hurles, M. E. & Lupski, J. R. Copy number variation in human health, disease, and
39
40 evolution. *Annu. Rev. Genomics Hum. Genet.* **10**, 451–481 (2009).
41
- 42 10. Ye, K., Schulz, M. H., Long, Q., Apweiler, R. & Ning, Z. Pindel: a pattern growth approach to detect break
43
44 points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics* **25**,
45
46 2865–2871 (2009).
47
48
- 49 11. Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end and split-read analysis.
50
51 *Bioinformatics* **28**, i333–i339 (2012).
52
- 53 12. Handsaker, R. E., Korn, J. M., Nemesh, J. & McCarroll, S. A. Discovery and genotyping of genome
54
55 structural polymorphism by sequencing on a population scale. *Nat. Genet.* **43**, 269–276 (2011).
56
57
- 58 13. Kronenberg, Z. N. *et al.* Wham: Identifying Structural Variants of Biological Consequence. *PLoS Comput.*
59
60 *Biol.* **11**, e1004572 (2015).
61
62
63
64
65

- 1 14. Layer, R. M., Chiang, C., Quinlan, A. R. & Hall, I. M. LUMPY: a probabilistic framework for structural
2 variant discovery. *Genome Biol.* **15**, R84 (2014).
- 3
4
- 5 15. Zook, J. M. *et al.* Integrating human sequence data sets provides a resource of benchmark SNP and indel
6 genotype calls. *Nat. Biotechnol.* **32**, 246–251 (2014).
- 7
8
- 9 16. Thorvaldsdóttir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV): high-performance
10 genomics data visualization and exploration. *Brief. Bioinform.* **14**, 178–192 (2013).
- 11
12
- 13 17. Fiume, M., Williams, V., Brook, A. & Brudno, M. Savant: genome browser for high-throughput sequencing
14 data. *Bioinformatics* **26**, 1938–1944 (2010).
- 15
16
- 17 18. Munro, J. E., Dunwoodie, S. L. & Giannoulatou, E. SVPV: a structural variant prediction viewer for paired-
18 end sequencing datasets. *Bioinformatics* **33**, 2032–2033 (2017).
- 19
20
- 21 19. O'Brien, T. M., Ritz, A. M., Raphael, B. J. & Laidlaw, D. H. Gremlin: an interactive visualization model for
22 analyzing genomic rearrangements. *IEEE Trans. Vis. Comput. Graph.* **16**, 918–926 (2010).
- 23
24
- 25 20. Wyczalkowski, M. A. *et al.* BreakPoint Surveyor: a pipeline for structural variant visualization.
26
27
- 28 21. Spies, N., Zook, J. M., Salit, M. & Sidow, A. svviz: a read viewer for validating structural variants.
29
30
- 31 22. Chiang, C. *et al.* SpeedSeq: ultra-fast personal genome analysis and interpretation. *Nat. Methods* **12**,
32 966–968 (2015).
- 33
34
- 35 23. Abyzov, A., Urban, A. E., Snyder, M. & Gerstein, M. CNVnator: an approach to discover, genotype, and
36 characterize typical and atypical CNVs from family and population genome sequencing. *Genome Res.* **21**,
37 974–984 (2011).
- 38
39
- 40 24. [PDF]pysam documentation - Read the Docs.
41
42
- 43 25. Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Computing in Science Engineering* **9**, 90–95 (2007).
- 44
45
- 46 26. Zook, J. M. *et al.* Data Descriptor: Extensive sequencing of seven human genomes to characterize
47 benchmark reference materials. (2016). doi:10.1038/sdata.2016.25
- 48
49
- 50 27. Daniel Kortschak, R., S Pedersen, B. & L Adelson, D. bíogo/hts: high throughput sequence handling for
51 the Go language. *JOSS* **2**, 168 (2017).
- 52
53
- 54
55
- 56
57
- 58
59
- 60
61
- 62
63
- 64
65

28. Belyeu, J, R; Nicholas, T, J; Pedersen, B, S; Sasani, T, A; Havrilla, J, M; Kravitz, S, N; Conway, M, E; Lohman, B, K; Quinlan, A, R; Layer, R, M (2018): Supporting data for "SV-plaudit: A cloud-based framework for manually curating thousands of structural variants" GigaScience Database. <http://dx.doi.org/10.5524/100450>

FIGURE LEGENDS

Figure 1. Example *samplot* images of putative deletion calls that were scored as **A)** unanimously GOOD, **B)** unanimously BAD, and **C)** ambiguous with a mix of GOOD and BAD scores with respect to the top sample (NA12878) in each plot. The black bar at the top of the figure indicates the genomic position of the predicted SV, and the following subfigures visualize the alignments and sequence coverage of each sample. Subplots report paired-end (square-ends connected by a solid line, annotated as concordant and discordant paired-end reads in **A)** and split-read (circle-ends connected by a dashed line, annotated in **A)** alignments by their genomic position (x-axis) and the distance between mapped ends (insert size, left y-axis). Colors indicate the type of event the alignment supports (black for deletion, red for duplication, and blue and green for inversion) and intensity indicates the concentration of alignments. The grey filled shapes report the sequence coverage distribution in the locus for each sample (right y-axis, annotated in **A)**. The samples in each panel are a trio of father (NA12891), mother (NA12892), and daughter (NA12878).

Figure 2. A) The distribution of the time between when an image was presented and when it was scored. B) The distribution of curation scores. C) The SV size distribution for all, unanimous (score 0 or 1), unambiguous (score <0.2 or >0.8) and ambiguous (score ≥ 0.2 and ≤ 0.8) variants. D) A comparison of predictions for deletions between CNVNATOR copy number calls (y-axis), SVTYPER genotypes (color, "Ref." is homozygous reference and "Non-ref." is heterozygous or homozygous alternate), and curation scores (x-axis). This demonstrates a general agreement between all methods with a concentration of reference genotypes and copy number two (no evidence for a deletion) at curation score less than 0.2, and non-reference and copy number one or zero events (evidence for a deletion) at curation score greater than 0.8. There are also false positives for CNVNATOR (copy number less than 2 at score = 0), and false negatives for SVTYPER (reference genotype at score = 1).

1 **Figure 3.** The *SV-Plaudit* process. **A)** *Samplot* generates an image for each SV from VCF considering a set of
2
3 alignment (BAM or CRAM) files. **B)** *PlotCritic* uploads the images to an Amazon S3 bucket and prepares
4
5 DynamoDB tables. Users select a curation answer (“GOOD”, “BAD”, or “DE NOVO”) for each SV image.
6
7 DynamoDB logs user responses and generates reports. Within a report, a curation score function can be
8
9 specified by mapping answer options to values and selecting an aggregation function. Here “GOOD” and “DE
10
11 NOVO” were mapped to one, “BAD” to zero, and the mean was used. One useful output option for a report is a
12
13 VCF annotated with the curation scores (shown here in bold as a **SVP**).
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

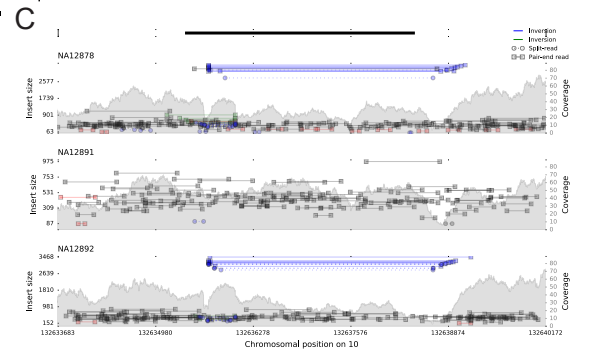
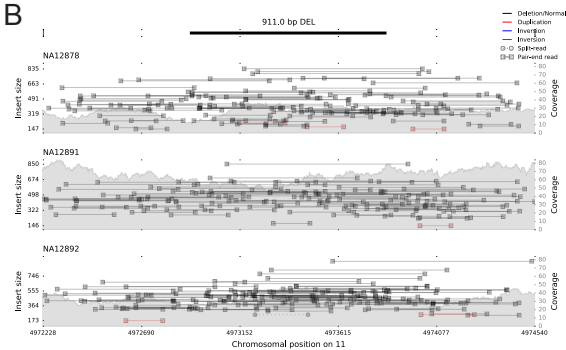
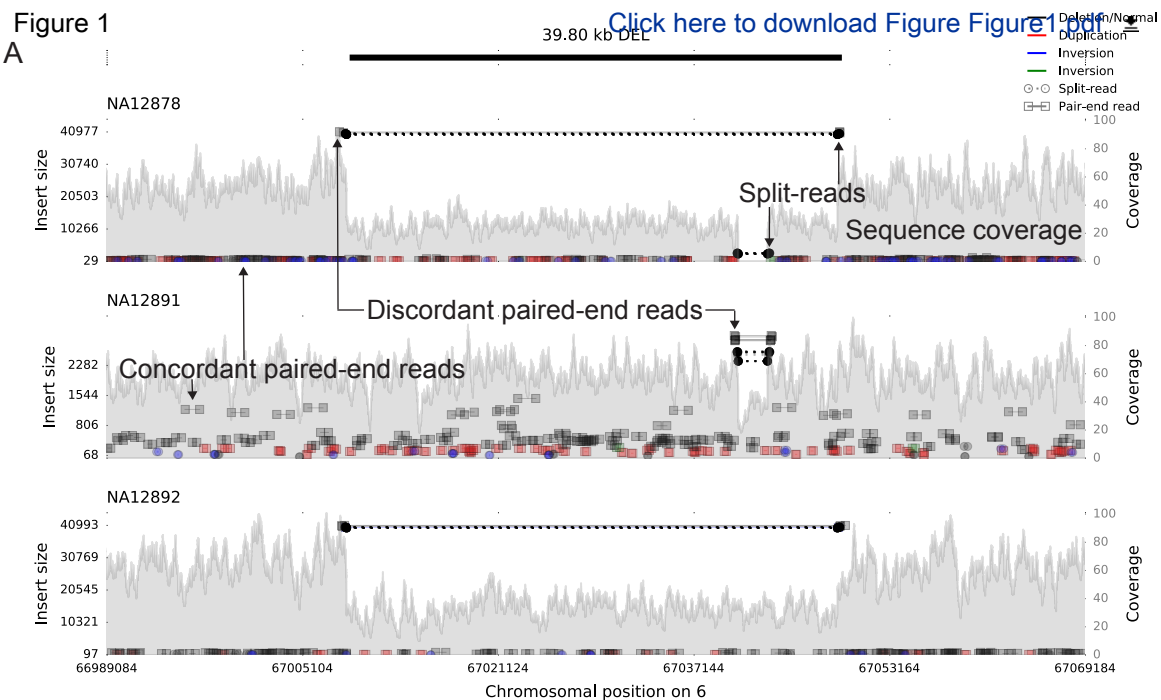
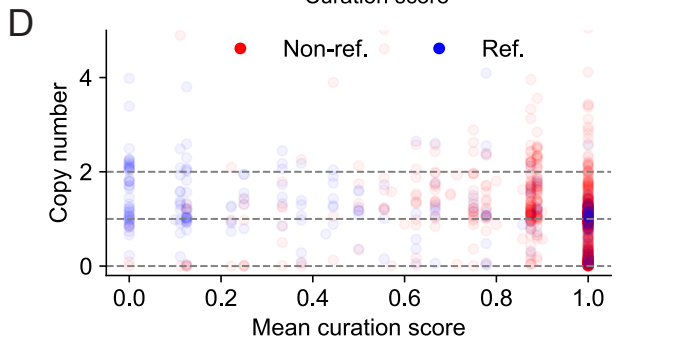
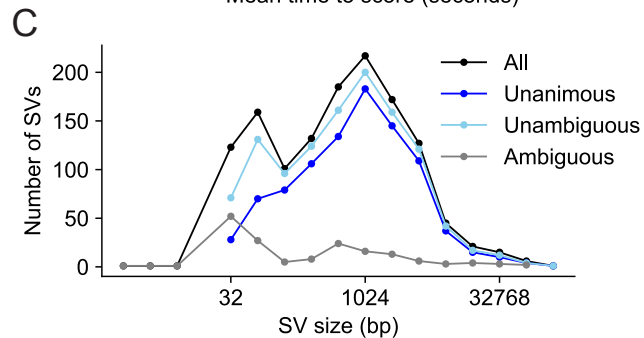
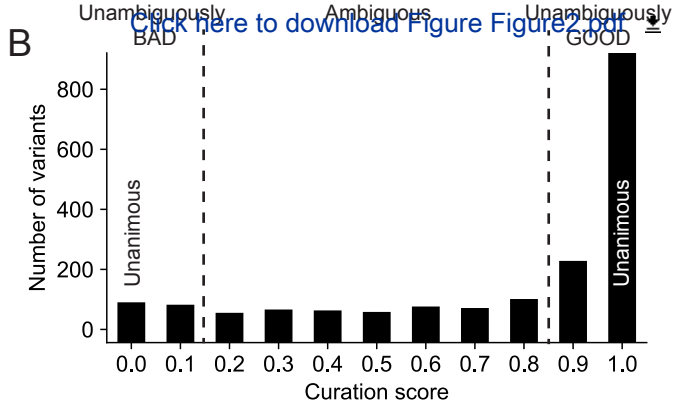
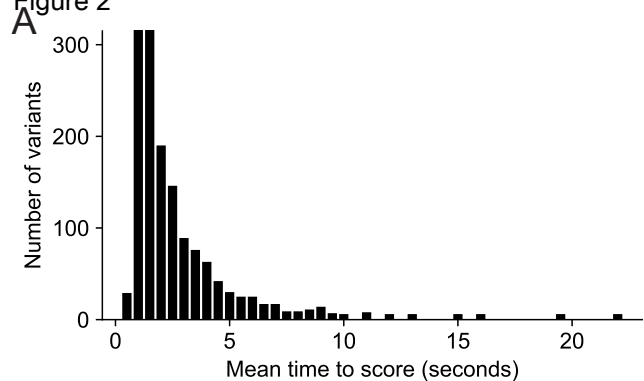
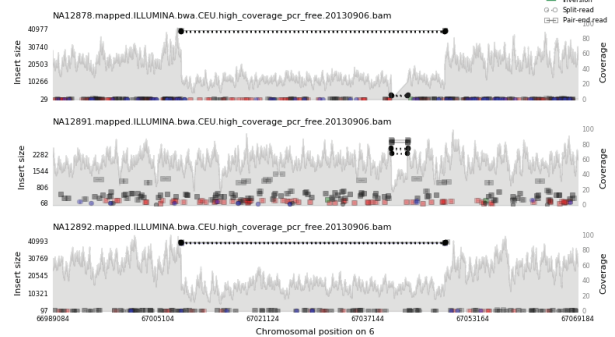


Figure 2

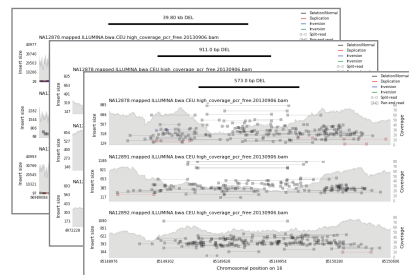


#CHROM	POS	INFO
6	67009228	SVTYPE=DEL;END=67049033 → samplot →
11	4972926	SVTYPE=DEL;END=4973937
16	85149501	SVTYPE=DEL;END=85150074

BAMs
 NA12878.bam
 NA12891.bam
 NA12892.bam

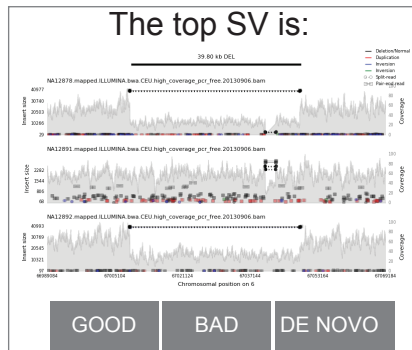


B



config.json

```
{ "curationVariables" : {
  "curationQandA" : {
    "question" :
      "The top SV is:",
    "answers" : {
      "g", "GOOD",
      "b", "BAD",
      "d", "DE NOVO" } },
  "AWSValues" : { ... } }
```



User input

Type: DEL
 Chrom: 6
 Start: 67009228
 End: 67049033
 User: jon@utah.edu
 Result: GOOD



S3

DynamoDB

PlotCritic


VCF report

#CHROM	POS	INFO
6	67009228	SVTYPE=DEL;END=67049033; SVP=1.0
11	4972926	SVTYPE=DEL;END=4973937; SVP=0.0
16	85149501	SVTYPE=DEL;END=85150074; SVP=0.5



Click here to access/download
Supplementary Material
Supplemental_File_1.vcf





Click here to access/download
Supplementary Material
Supplemental_File_3.sh



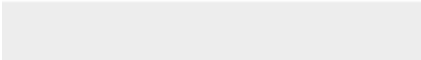



Click here to access/download
Supplementary Material
Supplemental_File_4.csv





Click here to access/download
Supplementary Material
Supplemental_File_5.vcf





Click here to access/download
Supplementary Material
Supplemental_File_6.txt

