1    LAVA: a streamlined visualization tool for longitudinal analysis of viral alleles

2

3    Michelle J. Lin[1,2,¶]          mjlin@uw.edu

4    Ryan C. Shean[1,2,¶]          rcs333@uw.edu

5    Negar Makhsous[1,2]          negarm@uw.edu

6    Alexander L. Greninger[1,2,*]    agrening@uw.edu

7

8    [1]Department of Laboratory Medicine, University of Washington, Seattle, WA, USA

9    [2]Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle,

10   WA, USA

11   [¶]These authors contributed equally

12   [*]Corresponding author

13

14

15

16

17

18

19

20

21

22

23

24 **Abstract**

25 With their small genomes, fast evolutionary rates, and clinical significance, viruses have long

26 been fodder for studies of whole genome evolution. One common need in these studies is the

27 analysis of viral evolution over time through longitudinal sampling. However, there exists no

28 simple tool to automate such analyses. We created a simple command-line visualization tool

29 called LAVA (Longitudinal Analysis of Viral Alleles). LAVA allows dynamic and interactive

30 visualization of viral evolution across the genome and over time. Results are easily shared via a

31 single HTML file that also allows interactive analysis based on read depth and allele frequency.

32 LAVA requires minimal input and runs in minutes for most use cases. LAVA is programmed

33 mainly in Python 3 and is compatible with Mac and Linux machines. LAVA is a user-friendly

34 command-line tool for generating, visualizing, and sharing the results of longitudinal viral

35 genome evolution analysis. Instructions for downloading, installing, and using LAVA can be

36 found at https://github.com/michellejlin/lava.

37 **Keywords:** LAVA, viral evolution, longitudinal, viral allele, antiviral resistance, bioinformatics,

38 visualization, NGS

39

40

41

42

43

44

45

46

47

48    **Introduction**

49        With the rapid and significant advancements in sequencing technologies in recent years,

50    whole-genome sequencing has become more cost-effective, more efficient, and more accurate

51    than ever (1). A common area of bioinformatics research in virology is the comparison of viral

52    evolution in longitudinal samples.  At a basic level, viral genome evolution may be examined

53    over routine passage in cell culture (2,3).  Drug manufacturers routinely check for the

54    development of resistance mutations in response to *in vitro* antiviral pressure (4,5). Clinical

55    researchers want to know how viruses evolve longitudinally in normal or immunocompromised

56    patients, in response to a drug pressure, or in different areas of the body (3,6,7).

57        In order to facilitate these routine analyses of viral evolution, we developed a simple

58    command-line tool called Longitudinal Analysis of Viral Alleles (LAVA) for analyzing and

59    visualizing the evolution of minor variants in viral genomes over time. The basic tenor of these

60    analyses involves the calling of a consensus genome for the initial sample and then using that

61    genome as a reference for downstream samples.  Viral sequence data is plotted both across the

62    genome to show where mutations cluster and over time to show allele frequency changes.  The

63    metadata associated with the experiment may be minimal, consisting simply of sample names

64    and units of time.  The units of time are arbitrary and may be minutes, hours, days, months,

65    years or even different categorical experimental conditions. LAVA also generates interactive

66    HTML files for sequence data analysis.  The HTML files may be manipulated by users without

67    significant bioinformatic experience according to the nature of their biological question,

68    alleviating a significant conundrum for sequencing and bioinformatics groups as demand for

69    their services continues to increase.

70

**Methods**

71

72    LAVA (Longitudinal Analysis of Viral Alleles) can be downloaded at

73    https://github.com/michellejlin/lava. Installation and usage instructions, a folder with example

74    inputs, and the full source code, are also available at this link. The general workflow is shown in

75    Figure 1. A brief explanation is also given here, but a more in depth look into the pipeline is

76    available at the GitHub link, including options and arguments passed to third party tools, the full

77    LAVA source code, as well as an informative readme document.

78    Installation of LAVA and all required dependencies is performed by an install script

79    which is included in the GitHub repository. The install script only requires Python, a Java

80    runtime environment, brew/apt-get for Mac/Linux systems, and an Internet connection. All third

81    party tools except for ANNOVAR (8), which must be manually registered for and downloaded,

82    are also automatically installed. The install script can also be run in 'check mode' and the script

83    will check for all required dependencies and print error messages with instructions for how to fix

84    any missing dependencies. The GitHub readme also contains a walkthrough for manually

85    installing all dependencies and LAVA.

86    Before execution, LAVA requires a reference genome for sequencing read alignment,

87    which can be provided as either an NCBI GenBank Accession number, or a local nucleotide

88    FASTA file along with a GFF file to provide gene and protein annotations (9).  Sequencing reads

89    for all samples are input as adapter and quality-trimmed FASTQ files. LAVA currently does not

90    perform adapter or quality trimming so this needs to be done beforehand as required. A two-

91    column CSV metadata file is also required to match sample with its longitudinal temporal

92    information such as passage number or day number.  Other available options include removing

93    PCR duplicates for metagenomic sequencing, manually specifying the name of the output

94    folder, analyzing variants by nucleotide changes instead of amino acid changes, saving all

95    intermediate files (which LAVA by default removes), automatically specifying an allele frequency

96    cutoff (which is by default set at 1%), and saving output as PNG files.

97      There are two different methods of selecting annotations for the reference sequence:

98      automatic and manual. In automatic mode LAVA begins with searching GenBank for a user

99      provided accession number corresponding to the viral species to be analyzed (10). This record

100     is then downloaded both as a nucleotide FASTA file and the complete GenBank record. LAVA

101     aligns the first FASTQ file provided to the downloaded reference (11,12) and calls a majority

102     consensus sequence based off this alignment using samtools (13–15). Then coding sequence

103     annotations are pulled from the GenBank record and transferred to the new majority consensus

104     FASTA using a MAFFT alignment (16). In manual mode the user specifies a reference FASTA

105     and a GFF file containing protein annotations for this reference sequence. LAVA assumes that

106     the FASTA is the majority consensus for the first sample and the GFF is a correct annotation of

107     the reference nucleotide sequence. The result of both the automatic and manual processes is

108     an annotated majority consensus of the first sample.

109     LAVA then aligns each of the FASTQ files specified in the metadata CSV to this newly

110     generated file, using bwa-mem (17). By default, LAVA does not remove PCR duplicates given

111     the common use of AmpliSeq-like approaches for viral genome sequencing; however, this

112     option can be added in cases where removing PCR duplicates would give a more accurate

113     representation of the data, such as the analysis of metagenomic samples.  Variants are called

114     for every position in the genome for every sample using VarScan and saved as standard VCF

115     files (18). These files are removed from the output folder during cleanup to keep disk usage low,

116     but can be saved using the --save option. Variants for all bases are annotated using Annovar

117     and GATK as nonsynonymous, synonymous, complex, stop-gain, or stop-loss, along with the

118     coverage and allele percentage at each base (8,19). The main text file generated by the pipeline

119     is a table called merged.csv containing all the samples, their metadata, and all the amino acid

120     changes.  This file, along with reads.csv, the individual .bam files, and the individual

121     .genomecov files, is used for generating the interactive visualization but can also be manually

122     parsed and examined for more in-depth or non-standard analysis. Reads.csv provides read

123    mapping information for each sample, such as total number of reads in sample and percentage

124    of reads mapping to reference. A .bam file is generated for each sample during the alignment

125    process, and these can be viewed for understanding the alignments and how the reads were

126    mapped. Genome coverage for each base in each sample is parsed and extracted into a file

127    with extension .genomecov, so genome-wide depth can be examined and analyzed.

128        LAVA then visualizes this information with the Bokeh Python module (20–22), allowing

129    for an easily readable and interactive data visualization. The output for this step is an HTML file

130    containing two interactive plots (Fig 2). The first plot depicts allele frequency changes for each

131    variant across the genome for each sample. Tabs at the top of the plot allow easy switching

132    between samples.  Sliders to the right of the main plot allow the user to dynamically change the

133    visibility of variants by depth and by coverage. To the right side of the plot, a line graph shows

134    per-base coverage across the whole genome to help inform the user of reasonable coverage

135    thresholds.  The second plot shows allele frequency across the samples over time. Given that

136    the samples will be representing different time points (passages, cultures, days past infection,

137    etc.) of a single virus, this plot shows the longitudinal evolution of amino acid changes,

138    separated by protein. Tabs at the top of the plot allow the user to specify which protein they

139    want to examine. Allele frequencies for all changes in the selected protein are plotted over time.

140    Here, variants can also be filtered by depth and coverage. Both plots support zooming and

141    panning and each mutational change has an associated tooltip which can be viewed by

142    hovering with the mouse over the associated data point to display locus-associated metadata.

143    Data can be filtered by type of mutation (synonymous, non-synonymous, stopgains/stoplosses,

144    and complex mutations), as well as if the same mutation occurs across multiple samples.

145    Another available option in the command line is to show nucleotide changes such as

146    transversions and transitions rather than amino acid changes, which may be relevant to cases

147    when examining nucleoside-analog antivirals directed against viral polymerases, base editors

148    such as APOBEC or ADAR proteins, or other aspects of viral epigenetics (23–25).

149    LAVA outputs all these files in HTML format (Fig 2), which are readily interpretable in

150    any web browser by groups without significant bioinformatics experience. Once generated by

151    LAVA, all these graphs can be sent and shared as standalone files. Additionally, LAVA also has

152    an option to generate static PNG images of the results for situations where interactive

153    visualization is not appropriate such as publications or presentations.

154

155    **Results and Discussion**

156    To demonstrate the intended use cases of LAVA and demonstrate why it represents a

157    new and useful tool, we illustrate two real world examples from our own lab.

158

159    *Case Study 1 – Evolution of human parainfluenza virus 3 in culture*

160    The provided examples (https://github.com/michellejlin/lava/tree/master/example), which

161    are included with the software, illustrate the automation of a task which the authors first

162    performed manually. For case study 1, these example files are truncated versions of the real

163    data analyzed in Iketani et al. (26), and are named Example 1 in LAVA. Example 1 illustrates

164    how to use LAVA to rapidly perform whole genome analysis on matched samples to understand

165    how a unique selective pressure (i.e. culture exposure) affects viral evolution.

166    Briefly, paired human parainfluenza virus type 3 (HPIV3) samples were sequenced

167    directly from nasal sampling and after isolation in culture. The study aimed to examine how

168    HPIV3 adapts to brief exposure to culture. Sequencing reads were adapter and quality trimmed

169    using cutadapt, producing the FASTQ files available on the GitHub (27). A simple metadata

170    sheet called Example1_metadata.csv was created containing file names and 'passage

171    numbers'. In this case, because we only use two samples, we put the first sample (nasal swab

172    SC332) occurring at passage 0 and the second sample (cultured CUL 332) at passage 1. We

173    have also provided a manually generated GFF/FASTA reference pair containing the protein

174    locations for all HPIV3 proteins except C and D (which are created through RNA editing and

175    thus do not automatically translate correctly).  To reduce the file size, the example files

176    uploaded to GitHub contain only the first 20,000 original reads that correctly mapped to HPIV3 –

177    full sequencing read files are available from BioProject PRJNA338014.  Example 1 shows how

178    rapid adaptations to culture can be discovered using LAVA as two non-synonymous mutations

179    (S554G and P241L) appear in the sample after very brief growth in culture. This example also

180    shows off the utility of the depth and allele frequency sliders which can be used to quickly filter

181    low-level sequencing artifacts and mapping errors out of the data, allowing the user to focus on

182    the most relevant points of data.

183

184    *Case Study 2*

185    We have also included data for a case study which fully highlights the longitudinal

186    analysis nature of LAVA. In this study, norovirus samples were recovered from a >250 day

187    infection over 11 time points from a single patient (28). The fundamental question in this

188    analysis is what whole genome changes accrue as norovirus adapts to the

189    immunocompromised host over almost a yearlong period.

190    Samples were sequenced and reads were adapter and quality trimmed using cutadapt

191    as part of our routine metagenomics analysis pipeline (27). As in Example 1, we selected a few

192    samples: ST107, ST283, and ST709 (all available on BioProject PRJNA338014). Reads were

193    trimmed to reduce file size to upload onto GitHub. A two-column metadata sheet called

194    Example2_metadata.csv was created mapping samples to collection day. The analysis was run

195    with the one-line command "lava.py -q MH260507 ST107.fastq Example2_metadata.csv -o

196    norovirus_output" (MH260507 is the GenBank Accession number for the actual day 0

197    consensus of these samples). This command showcases the alternative method of generating

198    reference files: using the -q flag to automatically download a GenBank reference and transfer

199    annotations. This example also highlights the utility of the protein plots, which show how the

200    allele frequency of all variants for each protein changes over time. Instead of using passages as

201     in Example 1, these plots demonstrate the evolution over number of days of infection. Using

202     these plots, one can see how the entire norovirus genome accumulates fixed mutational

203     changes over a long-term infection with an increased rate of fixed mutational changes in VP1,

204     the capsid protein and main antigenic determinant of norovirus (29).

205

206     *Comparisons*

207         While there are many programs that process and visualize somatic mutations, LAVA is

208     unique in its focus on monitoring minor variant alleles in viruses (30–32). With both its

209     component parts of pipeline and visualizer, LAVA fills an important need in the viral

210     bioinformatics community. The Broad Institute, for example, has several well-documented

211     workflows for both germline and somatic variant discovery: HaplotypeCaller and MuTect2.

212     These tools are excellent for their intended use cases and LAVA uses a workflow inspired by

213     these tools. However, HaplotypeCaller is not well suited for whole genome analysis of viral

214     genomes, as the tool is focused on germline SNPs and does not handle the extreme allelic

215     variance found in viral genomes. MuTect2, the Broad Institute's somatic SNP and indel caller,

216     performs well for its intended use but does not emit all bases of a genome, which is vital

217     information for viral whole genome analysis. Both of these tools are excellent for their intended

218     purposes but would have to be significantly modified to reproduce the analysis of LAVA.

219         The Broad Institute's viral-ngs suite, pipelines designed specifically for the analysis of

220     viral genomics, takes paired-end reads and calls intrahost variants (iSNVs). Taxonomic read

221     identification is also visualized with Krona. For variant calling in viral genomes, viral-ngs is an

222     excellent tool and we recommend using it over LAVA. However, LAVA was created specifically

223     to automatically compare longitudinal data, which is not a built-in feature of viral-ngs. LAVA also

224     has a visualization tool to easily see and compare minor allele variants across the genome and

225     across time. In these use cases, LAVA adds functionality over other bioinformatics programs.

226        Two other bioinformatics pipelines exist that perform similar tasks as LAVA. SMuPFi is a

227    pipeline that, like LAVA, analyzes NGS data to provide a graphical representation of SNPs and

228    works well for viral analysis (33). However, due to its nature as a tool designed to better

229    understand viral escape mechanisms, SMuPFi operates in the area of co-occurring mutations,

230    and works best with only two co-occurring mutations at the same time due to the complex

231    statistical analysis involved.

232        Another pipeline that serves to identify variant sites is ViVan (34). ViVan takes similar

233    input as LAVA and has a very easy to use, albeit size limit restricted, web interface. It also

234    detects more sensitive variant alleles than LAVA does—it claims to identify variant alleles with a

235    frequency of >0.1%, with a slightly higher rate of false positives, whereas LAVA by default both

236    filters out any minor allele variants below 1% frequency (though this can be adjusted using the -

237    af argument), and allows dynamic filtration in its visualization to suit the user's purpose. ViVan

238    searches for variants within each sample individually and currently provides no built-in feature

239    for comparisons between samples.

240        LAVA combines many of the gold-standard bioinformatics tools into a single pipeline to

241    annotate minor allele variants in viruses and adds a truly unique functionality with its interactive

242    visualization. The plots that LAVA outputs present easily understandable comparisons between

243    longitudinal samples, illustrating complex relationships in a simple format that makes patterns

244    like evolution of minor allele variants across samples, nucleotide change frequency in different

245    proteins, and synonymous vs. nonsynonymous mutations in the genome evident. By allowing

246    dynamic filtering of data by allele frequency and coverage depth, these plots can be adjusted to

247    suit the individual needs of the user.

248        Additionally, the inherently shareable nature of the HTML plots that LAVA creates as

249    output is an advantage. The small size, ability to be viewed on any web browser, and lack of

250    dependencies allow data to be shared quickly and extensively through email or any other

251    means, especially with collaborators who are not comfortable filtering BAM and VCF files.

252

253     *Limitations*

254          LAVA is a powerful tool for analyzing a diverse variety of viral datasets, yet it is not

255     without its limitations. While stopgains and stoplosses are handled correctly and included in the

256     plots, LAVA is currently unable to handle complex mutations, wherein two neighboring

257     nucleotide variants occur within a single codon. Multiple nucleotide changes within the same

258     codon are each treated individually as separate amino acid changes. However, LAVA

259     automatically detects this situation, and both prints a warning to the console and colors points

260     corresponding to complex mutations distinctly. Sequence variations such as copy number

261     changes, recombination, or large deletions and insertions that escape the bwa-mem aligner

262     may also be missed (35). Due to the nature of its visualization, LAVA also does not display

263     overlapping genes properly and instead shows them side-by-side. However, LAVA does print a

264     warning message to the console if overlapping proteins are detected, directing users to the

265     README which contains directions for how to manually prepare a GFF file without overlapping

266     proteins.  LAVA also does not correctly analyze proteins with RNA editing or ribosomal slippage.

267     Many of these limitations can be fixed by editing the GFF file accordingly.

268          Another limitation of LAVA is that web browsers can fail to render the output plots if there

269     are an extremely large number of variants (>5,000). This does not impact the actual analysis,

270     only the visualization, and the merged.csv output file will still contain all relevant data. This could

271     create problems if LAVA was used to analyze bacterial genomes or other extremely large

272     genomes. LAVA will print a warning message if there are greater than 5,000 variants. The

273     nature of the merged.csv output file is such that manual analysis could easily be performed in

274     an environment better suited to visualizing extremely large data sets such as R.

275

276     **Conclusions**

277        LAVA allows users to go from sequencing data to dynamically interactive plots

278    illustrating longitudinal changes in their samples. The only required inputs are 1) FASTQ files

279    with sequences for analysis, 2) either a GFF file and reference FASTA or a Genbank accession

280    number, and 3) a simple metadata.csv file containing information about sample name and

281    passage number. LAVA cuts down the time and effort significantly for data analysis of

282    longitudinal samples, and provides an intuitive and interactive visualization that can be easily

283    shared among collaborators.

284

285    *Web resources*

286    LAVA can be found at https://github.com/michellejlin/lava and is programmed in Python.

287

288    *Acknowledgements*

289    The authors would like to acknowledge the Broad institute for Picard, as well as the developers

290    and maintainers of UCSC Genome Browser for gff3ToGenePred. We would also like to thank

291    the entire open source bioinformatics community for their commitment to producing freely

292    available and useful tools for everyone.

293

## References

1.    Muir P, Li S, Lou S, Wang D, Spakowicz DJ, Salichos L, et al. The real cost of sequencing: scaling computation to keep pace with data generation. Genome Biol [Internet]. 2016 Dec [cited 2019 Mar 18];17(1). Available from: http://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-0917-0

2.    Iketani S, Shean RC, Ferren M, Makhous N, Aquino DB, Georges A, et al. Viral Entry Properties Required for Fitness in Humans Are Lost Through Rapid Genomic Change during Viral Isolation. mBio. 2018;9(4):e00898-18.

3.    Xue KS, Stevens-Ayers T, Campbell AP, Englund JA, Pergam SA, Boeckh M, et al. Parallel evolution of influenza across multiple spatiotemporal scales. eLife [Internet]. 2017 Jun 27 [cited 2019 Mar 18];6. Available from: https://elifesciences.org/articles/26875

4.    Toots M, Yoon J-J, Cox RM, Hart M, Sticher ZM, Makhsous N, et al. Characterization of orally efficacious influenza drug with high resistance barrier in ferrets and human airway epithelia. Sci Transl Med. 2019 Oct 23;11(515).

5.    Yoon J-J, Toots M, Lee S, Lee M-E, Ludeke B, Luczo JM, et al. Orally Efficacious Broad-Spectrum Ribonucleoside Analog Inhibitor of Influenza and Respiratory Syncytial Viruses. Antimicrob Agents Chemother. 2018 Aug;62(8).

6.    Debbink K, McCrone JT, Petrie JG, Truscon R, Johnson E, Mantlo EK, et al. Vaccination has minimal impact on the intrahost diversity of H3N2 influenza viruses. PLoS Pathog. 2017;13(1):e1006194.

7.    McCrone JT, Woods RJ, Martin ET, Malosh RE, Monto AS, Lauring AS. Stochastic processes constrain the within and between host evolution of influenza virus. eLife [Internet]. 2018 May 3 [cited 2019 Mar 18];7. Available from: https://elifesciences.org/articles/35962

8.    Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. Nucleic Acids Res. 2010 Sep 1;38(16):e164–e164.

9.    Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank. Nucleic Acids Res. 2016 Jan 4;44(D1):D67–72.

10.    Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. BMC Bioinformatics. 2009;10(1):421.

11.    Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform. Bioinforma Oxf Engl. 2010 Mar 1;26(5):589–95.

12.    Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics. 2009 Jul 15;25(14):1754–60.

13.    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format and SAMtools. Bioinforma Oxf Engl. 2009 Aug 15;25(16):2078–9.

14.    Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. Bioinforma Oxf Engl. 2011 Aug 1;27(15):2156–8.

15.    Li H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. Bioinforma Oxf Engl. 2011 Nov 1;27(21):2987–93.

16.    Katoh K, Misawa K, Kuma K, Miyata T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. Nucleic Acids Res. 2002 Jul 15;30(14):3059–66.

337    17.    Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
338    ArXiv13033997 Q-Bio [Internet]. 2013 Mar 16 [cited 2019 Mar 18]; Available from:
339    http://arxiv.org/abs/1303.3997

340    18.    Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, et al. VarScan 2:
341    Somatic mutation and copy number alteration discovery in cancer by exome sequencing.
342    Genome Res. 2012 Mar 1;22(3):568–76.

343    19.    McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
344    Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA
345    sequencing data. Genome Res. 2010 Sep 1;20(9):1297–303.

346    20.    Bokeh Development Team. Bokeh: Python library for interactive visualization [Internet].
347    2014 [cited 2018 Oct 31]. Available from: http://www.bokeh.pydata.org

348    21.    Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely
349    available Python tools for computational molecular biology and bioinformatics. Bioinformatics.
350    2009 Jun 1;25(11):1422–3.

351    22.    Jones E, Oliphant E, Peterson P, et al. SciPy: Open Source Scientific Tools for Python
352    [Internet]. 2001 [cited 2018 Oct 31]. Available from: https://www.scipy.org/

353    23.    Yoon J-J, Toots M, Lee S, Lee M-E, Ludeke B, Luczo JM, et al. Orally Efficacious Broad-
354    Spectrum Ribonucleoside Analog Inhibitor of Influenza and Respiratory Syncytial Viruses.
355    Antimicrob Agents Chemother [Internet]. 2018 Jun 11 [cited 2019 Mar 18];62(8). Available from:
356    http://aac.asm.org/lookup/doi/10.1128/AAC.00766-18

357    24.    Smith HC, Bennett RP, Kizilyer A, McDougall WM, Prohaska KM. Functions and
358    regulation of the APOBEC family of proteins. Semin Cell Dev Biol. 2012 May;23(3):258–68.

359    25.    Nishikura K. Functions and regulation of RNA editing by ADAR deaminases. Annu Rev
360    Biochem. 2010;79:321–49.

361    26.    Iketani S, Shean RC, Ferren M, Makhsous N, Aquino DB, des Georges A, et al. Viral
362    Entry Properties Required for Fitness in Humans Are Lost through Rapid Genomic Change
363    during Viral Isolation. mBio. 2018 Jul 3;9(4).

364    27.    Martin M. Cutadapt removes adapter sequences from high-throughput sequencing
365    reads. EMBnet.journal. 2011 May 2;17(1):10.

366    28.    Casto AM, Adler AL, Makhsous N, Crawford K, Qin X, Kuypers JM, et al. Prospective,
367    Real-time Metagenomic Sequencing During Norovirus Outbreak Reveals Discrete Transmission
368    Clusters. Clin Infect Dis Off Publ Infect Dis Soc Am. 2019 Aug 30;69(6):941–8.

369    29.    Mahar JE, Donker NC, Bok K, Talbo GH, Green KY, Kirkwood CD. Identification and
370    Characterization of Antibody-Binding Epitopes on the Norovirus GII.3 Capsid. J Virol. 2014 Feb
371    15;88(4):1942–52.

372    30.    Ardin M, Cahais V, Castells X, Bouaoun L, Byrnes G, Herceg Z, et al. MutSpec: a
373    Galaxy toolbox for streamlined analyses of somatic mutation spectra in human and mouse
374    cancer genomes. BMC Bioinformatics [Internet]. 2016 Dec [cited 2019 Mar 18];17(1). Available
375    from: http://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-016-1011-z

376    31.    Lee J, Lee AJ, Lee J-K, Park J, Kwon Y, Park S, et al. Mutalisk: a web-based somatic
377    MUTation AnaLyIS toolKit for genomic, transcriptional and epigenomic signatures. Nucleic Acids
378    Res. 2018 Jul 2;46(W1):W102–8.

379    32.    Cario CL, Witte JS. Orchid: a novel management, annotation and machine learning
380    framework for analyzing cancer mutations. Hancock J, editor. Bioinformatics. 2018 Mar
381    15;34(6):936–42.
382    33.    Leung P, Bull R, Lloyd A, Luciani F. A Bioinformatics Pipeline for the Analyses of Viral
383    Escape Dynamics and Host Immune Responses during an Infection. BioMed Res Int.
384    2014;2014:1–12.
385    34.    Isakov O, Bordería AV, Golan D, Hamenahem A, Celniker G, Yoffe L, et al. Deep
386    sequencing analysis of viral infection and evolution allows rapid and detailed characterization of
387    viral mutant spectrum. Bioinformatics. 2015 Jul 1;31(13):2141–50.
388    35.    Greninger AL, Roychoudhury P, Makhsous N, Hanson D, Chase J, Krueger G, et al.
389    Copy Number Heterogeneity, Large Origin Tandem Repeats, and Interspecies Recombination
390    in Human Herpesvirus 6A (HHV-6A) and HHV-6B Reference Strains. Longnecker RM, editor. J
391    Virol [Internet]. 2018 Feb 28 [cited 2019 Mar 18];92(10). Available from:
392    http://jvi.asm.org/lookup/doi/10.1128/JVI.00135-18
393

394    **Figure Legends**

395

396    **Figure 1** - General workflow of the LAVA pipeline is depicted to offer a high-level overview of

397    program execution. Dashed arrows represent optional steps. Input are shown boxed in green,

398    output in blue, and the main lava program is circled in pink. For input, either a GenBank

399    Accession number or a FASTA/GFF pair is required. If a GenBank Accession number is

400    provided, LAVA generates a FASTA/GFF pair following the outlined steps. The linked chain

401    symbol between the metadata.csv input and the FASTQ reads is meant to emphasize that the

402    metadata.csv must contain all the file names that you wish to include in your analysis. General

403    steps are given with tools used during that specific step listed to the side or underneath each

404    step in parentheses. The final output is given as HTML files that contain the interactive plots.

405    For exactly what is passed to each of the other programs and information about parameters and

406    optional arguments (such as mapping parameters), the source code is available on GitHub.

407

408    **Figure 2 -** Example LAVA output is shown, this figure shows the results from running Example 1

409    (All files are available on GitHub and a more in-depth coverage of this data is provided in *Case*

410    *study 1.*) This example is a screenshot of a Chrome browser displaying the final HTML created

411    by LAVA. The plot on the top of the page shows all amino acid changes across the whole

412    genome for each sample. You can switch between the samples using the tabs highlighted in a

413    red box. The bottom plot shows changes in by-protein allele frequencies over time. You can use

414    tabs once again to switch between proteins. All changes meeting display requirements are

415    plotted over time (or whatever your numerical metadata was). For example, this example shows

416    the hemagglutinin-neuraminidase protein for HPIV3 undergoing changes during the culturing

417    process. All output can be filtered by depth, allele frequency and type of mutation using the

418    sliders boxed in red to the right of each main plot. A small plot is displayed next to the whole

419 genome graph providing a visual representation of the per-base coverage of reads mapping to
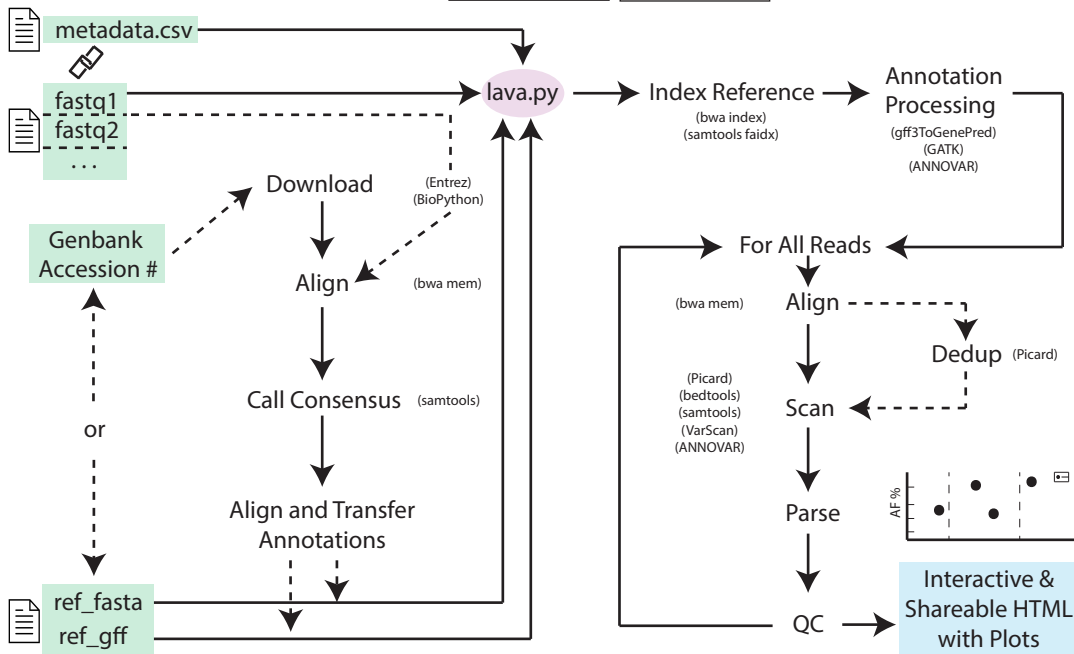
420 the consensus.

421

lava.py -f example.fasta -g example.gff   fastq1.fastq   metadata.csv   -o output_folder -save

Reference fasta/GFF pair to align first sample to. Can replace with -q
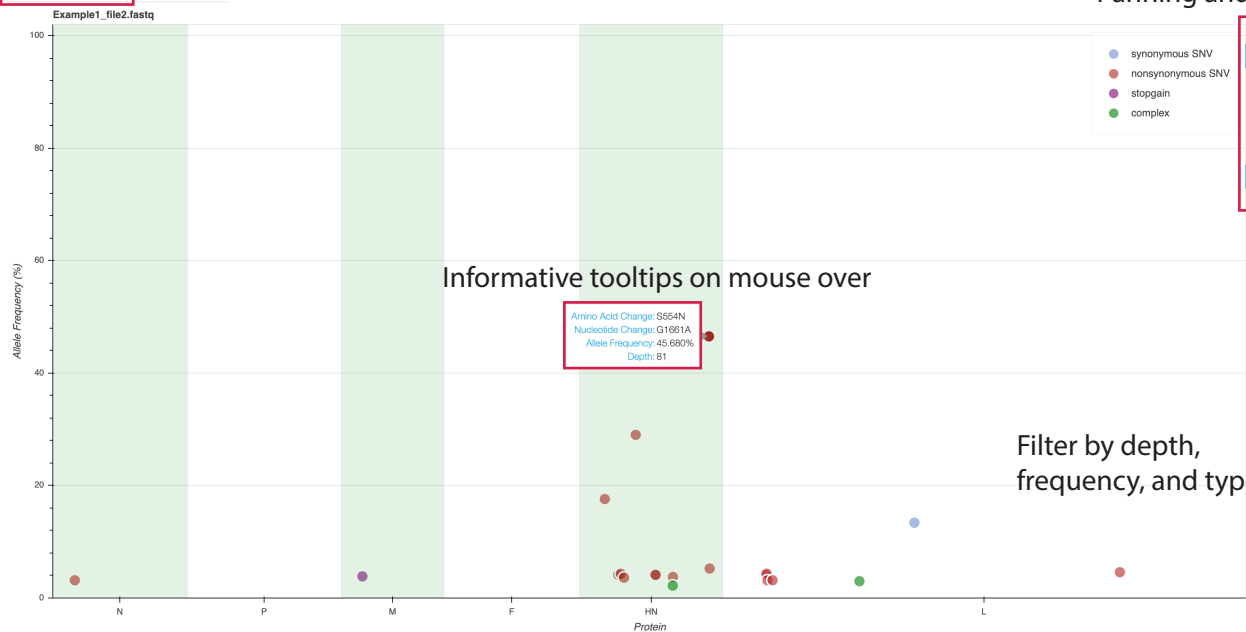
Set of reads to build reference majority consensus from

Metadata file with names and time values

Optional arguments for different output style

metadata.csv

fastq1
fastq2
. . .

lava.py → Index Reference
(bwa index)
(samtools faidx)
→ Annotation Processing
(gff3ToGenePred)
(GATK)
(ANNOVAR)

Download
(Entrez)
(BioPython)

Genbank
Accession #

Align
(bwa mem)

Call Consensus   (samtools)

Align and Transfer
Annotations

or

ref_fasta
ref_gff

For All Reads

Align   (bwa mem)   Dedup   (Picard)

Scan
(Picard)
(bedtools)
(samtools)
(VarStan)
(ANNOVAR)

Parse

AF %

QC → Interactive &
Shareable HTML
with Plots

Switch between samples

Panning and zooming utilities

Coverage at each position

Informative tooltips on mouse over

Amino Acid Change: S554N
Nucleotide Change: G1661A
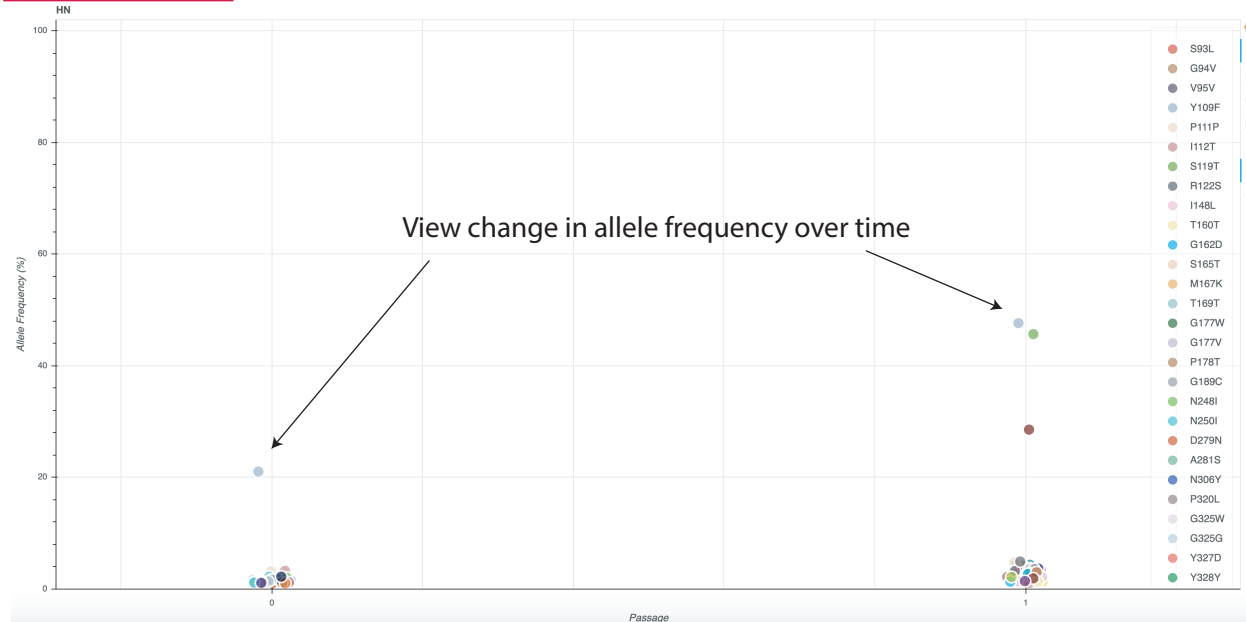Allele Frequency: 45.680%
Depth: 81

Filter by depth, frequency, and type

Whole genome information by protein

Read mapping information

Select proteins

View change in allele frequency over time

Filter and customize output