

VGEA: A snakemake pipeline for RNA virus genome assembly from next generation sequencing data

Paul E. Oluniyi^{1,2}, **Fehintola V. Ajogbasile**^{1,2}, **Judith U. Oguzie**^{1,2}, **Jessica N. Uwanibe**^{1,2}, **Adeyemi T. Kayode**^{1,2}, **Anise N. Happi**³, **Chinedu A. Ugwu**^{1,2}, **Testimony J. Olumade**^{1,2}, **Olusola Ogunsanya**³, **Philomena E. Eromon**², **Onikepe A. Folarin**^{1,2}, **Simon D.W. Frost**^{4,5}, **Jonathan L. Heeney**⁶, **Christian T. Happi**^{Corresp. 1,2}

¹ Department of Biological Sciences, Faculty of Natural Sciences, Redeemer's University, Ede, Osun, Nigeria

² African Centre of Excellence for Genomics of Infectious Diseases (ACEGID), Redeemer's University, Ede, Osun, Nigeria

³ Department of Veterinary Pathology, Faculty of Veterinary Medicine, University of Ibadan, Ibadan, Oyo, Nigeria

⁴ Microsoft Research, Redmond, 98052, Washington, United States of America

⁵ London School of Hygiene & Tropical Medicine, London, United Kingdom

⁶ Department of Veterinary Medicine, University of Cambridge, Cambridge, United Kingdom

Corresponding Author: Christian T. Happi

Email address: happic@run.edu.ng

Background. Next generation sequencing (NGS)-based studies have vastly increased our understanding of viral diversity. Viral sequence data obtained from NGS experiments are a rich source of information, these data can be used to study their epidemiology, evolution, transmission patterns, and can also inform drug and vaccine design. Viral genomes however represent a great challenge to bioinformatics due to their high mutation rate and forming quasispecies in the same infected host. This has therefore brought about the need to develop/implement advanced bioinformatics tools to assemble genomes well-representative of the viral population circulating in individual patients.

Results. Here we present VGEA (Viral Genomes Easily Assembled), a snakemake workflow for advanced assembly of RNA viral genomes from NGS data. VGEA enables users to split bam files into forward and reverse reads, carry out de novo assembly of forward and reverse reads to generate contigs, pre-process reads for quality and contamination, and map reads to a reference tailored to the sample using corrected contigs supplemented by the user's choice of reference sequences.

Conclusion. VGEA is freely available on GitHub at: <https://github.com/pauloluniyi/VGEA> under the GNU General Public License and also on Zenodo (doi: 10.5281/zenodo.3702287).

1 **VGEA: A snakemake pipeline for RNA virus genome assembly from next generation**
2 **sequencing data**

3

4 Paul E. Oluniyi^{1,2}, Fehintola V. Ajogbasile^{1,2}, Judith U. Oguzie^{1,2}, Jessica N. Uwanibe^{1,2},
5 Adeyemi T. Kayode^{1,2}, Anise N. Happi³, Chinedu A. Ugwu^{1,2}, Testimony J. Olumade^{1,2}, Olusola
6 Ogunsanya³, Philomena E. Eromon², Onikepe A. Folarin^{1,2}, Simon D.W. Frost^{4,5}, Jonathan L.
7 Heeney⁶, Christian T. Happi^{1,2*}

8

9 1. Department of Biological Sciences, Faculty of Natural Sciences, Redeemer's University, Ede,
10 Osun State, Nigeria.

11 2. African Centre of Excellence for Genomics of Infectious Diseases (ACEGID), Redeemer's
12 University, Ede, Osun State, Nigeria.

13 3. Department of Veterinary Pathology, Faculty of Veterinary Medicine, University of Ibadan,
14 Ibadan, Nigeria.

15 4. Microsoft Research, Redmond, 98052, WA, USA.

16 5. London School of Hygiene & Tropical Medicine, London, UK.

17 6. Department of Veterinary Medicine, University of Cambridge, Cambridge, United Kingdom.

18

19

20 *Corresponding Author:

21 Christian Happi^{1,2}

22 Email address: happic@run.edu.ng

23

24 **ABSTRACT**

25

26 **Background.** Next generation sequencing (NGS)-based studies have vastly increased our
27 understanding of viral diversity. Viral sequence data obtained from NGS experiments are a rich
28 source of information, these data can be used to study their epidemiology, evolution,
29 transmission patterns, and can also inform drug and vaccine design. Viral genomes however
30 represent a great challenge to bioinformatics due to their high mutation rate and forming
31 quasispecies in the same infected host. This has therefore brought about the need to
32 develop/implement advanced bioinformatics tools to assemble genomes well-representative of
33 the viral population circulating in individual patients.

34 **Results.** Here we present **VGEA (Viral Genomes Easily Assembled)**, a snakemake workflow for
35 advanced assembly of RNA viral genomes from NGS data. VGEA enables users to split bam
36 files into forward and reverse reads, carry out *de novo* assembly of forward and reverse reads to
37 generate contigs, pre-process reads for quality and contamination, and map reads to a reference
38 tailored to the sample using corrected contigs supplemented by the user's choice of reference
39 sequences.

40 **Conclusion.** VGEA is freely available on GitHub at: <https://github.com/pauloluniyi/VGEA>
41 under the GNU General Public License and also on Zenodo (doi: 10.5281/zenodo.3702287).

42

43 **Keywords:** VGEA, NGS, Genome, Assembly

44

45

46

47 **INTRODUCTION**

48 The most abundant biological entities on Earth are viruses as they can be found among all
49 cellular forms of life. So far, over four thousand five hundred viral species have been discovered,
50 from which a huge amount of sequence information has been collected by researchers and
51 scientists all over the world (*Pickett et al., 2012; Sharma et al., 2015; Brister et al., 2015*). In
52 recent times (past two decades), a number of these viruses have emerged in the human
53 population causing outbreaks sometimes pandemics., These viruses include mainly: Influenza
54 virus, Severe Acute Respiratory Syndrome (SARS) coronavirus, Middle East Respiratory
55 Syndrome (MERS) coronavirus, Ebola virus, Yellow fever virus, Lassa virus (LASV), Zika virus
56 (*Chan, 2002; Bean et al., 2013; Folarin et al., 2016; Grubaugh et al., 2017; Metsky et al., 2017;*
57 *Siddle et al., 2018; Ajogbasile et al., 2020*) and SARS-CoV-2 (*Chen et al., 2020; Holshue et al.,*
58 *2020; Sohrabi et al., 2020*). During these outbreaks and pandemics, identification of the
59 causative agents and carrying out sequencing to obtain the genomes of the viruses have proved to
60 be critical in helping inform disease surveillance and epidemiology.

61 NGS platforms have been widely accepted as high-throughput, unbiased technologies that have
62 many attractive features compared to conventional diagnostic methods for virus detection and
63 assembly (*Tang & Chiu, 2010; Mokili et al., 2012*). NGS-based studies have vastly increased our
64 understanding of viral diversity (*Reyes et al., 2010; Cantalupo et al., 2011*). Pathogen sequence
65 data obtained from NGS experiments are a rich source of information, these data can be used to
66 study their epidemiology, evolution, transmission patterns, and can also inform drug and vaccine
67 design. The field of genomics, especially pathogen genomics has been transformed by NGS, with
68 costs constantly decreasing, equipment becoming more portable/field deployable during

69 outbreaks and remarkable increase in data availability. The huge amount of data being generated
70 has brought about the need to develop simple and user friendly bioinformatics tools to assemble
71 pathogen genomes well-representative of the pathogen population circulating in individual
72 patients. Here we present **VGEA**, a pipeline for assembly of RNA viral genomes from next
73 generation sequencing data.

74 **MATERIALS AND METHODS**

75 The VGEA pipeline is built on the snakemake workflow management system (*Köster &*
76 *Rahmann, 2012*), a workflow management system that allows the effortless deployment and
77 execution of complex distributed computational workflows in any UNIX-based system, from
78 local machines to high-performance computing clusters. In order to guarantee reproducibility of
79 the results obtained, the VGEA pipeline integrates fixed versions of the tools implemented in the
80 pipeline from conda (<https://docs.conda.io/en/latest/>). Several tools are used to perform different
81 tasks within the pipeline: **Samtools** (*Li et al., 2009*) for splitting of bam files into forward and
82 reverse reads; **IVA** (*Hunt et al., 2015*) for *de novo* assembly to generate contigs; **Shiver**
83 (*Wymant et al., 2018*) to pre-process reads for quality and contamination, then map to a reference
84 tailored to the sample using corrected contigs supplemented with the user's choice of existing
85 reference sequences.

86

87 The **VGEA** pipeline requires the following dependencies:

- 88 Python 3 (www.python.org).
- 89 Samtools (*Li et al., 2009*).
- 90 IVA (*Hunt et al., 2015*).
- 91 Shiver (*Wymant et al., 2018*).

- 92 Fastp (*Chen et al., 2018*).
- 93 Trimmomatic, optional but highly recommended (*Bolger et al., 2014*).
- 94 KMC (*Kokot et al., 2017*).
- 95 MUMmer (*Marçais et al., 2018*).
- 96 SMALT (<https://www.sanger.ac.uk/science/tools/smalt-0>) or BWA (*Li & Durbin, 2009*)
- 97 or BOWTIE (*Langmead, 2010*).
- 98 Fastaq (<https://github.com/sanger-pathogens/Fastaq>).
- 99 Biopython (*Cock et al., 2009*).
- 100 MAFFT (*Katoh et al., 2002*).
- 101 BLAST version 2.2.28 (*Altschul et al., 1990*).
- 102 SPAdes (*Bankevich et al., 2012*).

103

104 We have also made available a singularity recipe file (*Kurtzer et al., 2017*) on the GitHub page:
105 <https://github.com/pauloluniyi/VGEA>. With the provision of the singularity recipe file, users can
106 easily build a local image of the VGEA container that includes all necessary tools, in their fixed
107 versions, and their dependencies using the command below:

108 ***sudo singularity build vgea.simg Singularity***

109 Users can then proceed to run the entire VGEA pipeline with all the dependencies installed from
110 the singularity container. This approach ensures the reproducibility and the tracking of both
111 software code and version, regardless of the operating system used. With the provision of the
112 singularity container, users can easily deploy VGEA to run in the cloud or on high performance
113 computing (HPC) clusters.

114 The VGEA pipeline consists of three major steps:

115

116 □ **Splitting of BAM files**

117 BAM (and SAM and CRAM) are file formats that contain sequencing reads: either aligned,
118 unaligned, or a combination of the two. One initial step after carrying out next-generation
119 sequencing is often to get rid of host contamination, e.g., by mapping the FASTQ reads obtained
120 from the sequencing machine against the human genome. This will ultimately yield a BAM file.
121 The BAM file obtained from host contaminant removal can then be used as input for the VGEA
122 pipeline. The pipeline has been developed to facilitate splitting of bam files into fastq files of
123 forward and reverse reads using **Samtools** (*Li et al., 2009*). Another reason for having BAM files
124 as starting input for the VGEA pipeline is that scientists, especially in resource-limited settings,
125 usually have BAM files handy. BAM files are usually smaller in size than their corresponding
126 FASTQ files making them easily transferable or uploadable to the cloud.
127 During the splitting step, the pipeline checks in the current directory or within the container (if
128 the user is making use of the singularity container provided) for any file with a ‘.bam’ extension,
129 if it finds any, it splits into FASTQ files of forward and reverse reads.

130

131 □ **Assembly**

132 Following splitting of the BAM files into FASTQ files of forward and reverse reads, the VGEA
133 pipeline carries out *de novo* assembly to generate contigs using **IVA** (*Hunt et al., 2015*). **IVA** is
134 used as our default assembler because it was designed specifically for read pairs sequenced at
135 highly variable depth from RNA virus samples and has been demonstrated to outperform all
136 other virus *de novo* assemblers (*Hunt et al., 2015*).

137

138 □ **Mapping**

139 The VGEA pipeline then uses the **Shiver** software (*Wymant et al., 2018*) for mapping of reads to
140 a reference. First, the contigs generated in the assembly step are compared with reference(s)
141 supplied by the user using BLASTN (*Altschul et al., 1990*), this is to remove contaminants and
142 low-quality contigs. Using **MAFFT** (*Katoh et al., 2002*), the processed contigs are added to the
143 alignment of existing references initially supplied by the user from which **Shiver** identifies the
144 closest existing reference by comparison with all of the contigs. Using contig sequences and the
145 closest existing reference to fill in gaps between contigs, if any exists, **Shiver** creates a reference
146 for mapping. User-supplied raw reads are then mapped to the shiver-created reference to
147 generate a consensus sequence well-representative of the viral population in the patient sample.
148 Before mapping however, the reads are trimmed using Trimmomatic (*Bolger et al., 2014*) and
149 Fastaq (<https://github.com/sanger-pathogens/Fastaq>) in order to remove low-quality bases,
150 adapters and primer sequences. Adapter and Primer sequences are provided by the user.

151

152 **RESULTS**

153

154 We demonstrated the usage and performance of the VGEA pipeline by applying it to generate
155 consensus whole genomes from Lassa virus and SARS-COV-2 datasets sequenced on the
156 illumina MiSeq and illumina FGx sequencing machines in our lab at the African Centre of
157 Excellence for Genomics of Infectious Diseases (ACEGID), Redeemer's University, Nigeria.
158 We also applied the pipeline to generate whole genomes from HIV-1 datasets sequenced on the
159 illumina HiSeq 2500 obtained from NCBI Sequence Read Archive (SRA). We made use of 60
160 test datasets (Lassa Virus (20), SARS-CoV-2 (20) and HIV-1 (20)) for the validation of the

161 VGEA pipeline. All our test datasets are available on figshare
162 (<https://doi.org/10.6084/m9.figshare.13009997>). The performance of our pipeline was consistent
163 irrespective of whether the samples were sequenced with the illumina MiSeq, FGx or HiSeq
164 2500. We compared our genomes obtained with **VGEA** to genomes obtained with **viral-ngs**
165 (<https://github.com/broadinstitute/viral-ngs>) (*Park et al., 2015*) which is a suite of genomic
166 analysis pipelines for viral sequencing and is one of the most widely used resources for whole
167 genome assembly of viruses (Tables 1,2,3&4).
168 We compared the means of our genome lengths using **VGEA** and **viral-ngs** to determine the
169 significance of the differences in genome lengths from the two pipelines; using the **R**
170 **programming language**, we carried out a student's t-test and found that genome length with
171 **VGEA** are significantly longer than genome length with **viral-ngs** (p-value = 0.002962, C.I =
172 95% for Lassa virus (S) dataset; p-value = 0.01748, C.I = 95% for Lassa virus (L) dataset; p-
173 value = 0.04283, C.I = 95% for SARS-CoV-2 dataset; p-value = 0.001286, C.I = 95% for HIV-1
174 datasets).
175 Using the **VGEA** pipeline, we also obtained Lassa virus partial genomes from three rodent
176 samples, however running the same samples through **viral-ngs** we couldn't obtain any genome
177 despite making the pipeline as less stringent as possible.

178

179

180

181

182

183 **DISCUSSION**

184 Using the workflow management system, **Snakemake**, we have developed a user-friendly
185 pipeline for advanced assembly of viral genomes from NGS data. Its main features are: (i)
186 splitting bam files into forward and reverse reads, (ii) *de novo* assembly to generate contigs, (iii)
187 pre-processing of reads for quality and contamination, (iv) mapping reads to a sample-tailored
188 reference. This significantly improves the quality of the genomes obtained from NGS data and
189 generates genomes well-representative of the pathogen population circulating in individual
190 patients and communities. VGEA is freely available on GitHub at:
191 <https://github.com/pauloluniyi/VGEA> under the GNU General Public License and also on
192 Zenodo (doi: 10.5281/zenodo.3702287). All test datasets used for the validation of the pipeline
193 are available on NCBI and also on figshare (<https://doi.org/10.6084/m9.figshare.13009997>).

194

195 CONCLUSION

196 VGEA was built primarily by biologists and in a manner that is easy to be employed by users
197 without significant computational background. As new and innovative tools for viral genome
198 analysis and assembly are increasingly being developed, these can easily be incorporated into the
199 VGEA pipeline. We hope that other scientists can build upon and improve VGEA as a tool to
200 extract more qualitative and quantitative information from viral genomes.

201

202

203

204

205

206 ACKNOWLEDGEMENTS

207 We appreciate the continuous support of ACEGID staff and the management of Redeemer's
208 University.

209

210 **Abbreviations**

211 **VGEA** Viral Genomes Easily Assembled

212 **NGS** Next generation sequencing

213 **RNA** Ribonucleic acid

214 **SARS** Severe Acute Respiratory Syndrome

215 **MERS** Middle East Respiratory Syndrome

216 **IVA** Iterative Virus Assembler

217 **SHIVER** Sequences from HIV Easily Reconstructed

218 **HPC** High Performance Computing

219

220 **ADDITIONAL INFORMATION AND DECLARATIONS**

221

222 **Funding**

223 This work was supported by grants from the National Institute of Allergy and Infectious Diseases

224 (<https://www.niaid.nih.gov>), NIH-H3Africa (<https://h3africa.org>) (U01HG007480 and

225 U54HG007480 to C.T.H), the World Bank grant (worldbank.org) (project ACE019 to C.T.H.)

226 and the Wellcome Trust grant (<https://wellcome.ac.uk>) (216619/Z/19/Z to C.T.H and J.L.H). The

227 funders had no role in study design, data collection and analysis, decision to publish, or

228 preparation of the manuscript.

229

230 Grant Disclosures

231 The following grant information was disclosed by the authors:

232 National Institute of Allergy and Infectious Diseases, NIH-H3Africa: U01HG007480 and

233 U54HG007480.

234 The World Bank: ACE019.

235 The Wellcome Trust: 216619/Z/19/Z.

236

237 Competing Interests

238 Simon D.W. Frost is employed by Microsoft Research. All other authors have declared that no

239 competing interests exist.

240

241 Author Contributions

242 ● Paul E. Oluniyi developed/implemented the pipeline and wrote the manuscript.

243 ● Paul E. Oluniyi, Fehintola V. Ajogbasile, Judith U. Oguzie, Jessica N. Uwanibe,

244 Adeyemi T. Kayode, Philomena E. Eromon, Anise N. Happi, Testimony J. Olumade and

245 Olusola Ogunsanya performed sequencing and metagenomics analysis.

246 ● Anise N. Happi, Chinedu A. Ugwu, Testimony J. Olumade and Olusola Ogunsanya

247 collected and processed rodent samples for sequencing.

248 ● Onikepe A. Folarin, Simon D.W. Frost, Jonathan L. Heeney and Christian T. Happi

249 supervised the project.

250 ● Christian T. Happi reviewed and corrected the manuscript and approved the final draft.

251

252 Data Availability

253 The following information was supplied regarding data availability:

254 VGEA is freely available on GitHub at: <https://github.com/pauloluniyi/VGEA> under the
255 GNU General Public License.

256

257 All primary test datasets used for the validation of the VGEA pipeline are available on figshare
258 (<https://doi.org/10.6084/m9.figshare.13009997>). All SARS-CoV-2 and Lassa virus test datasets

259 have been submitted to NCBI SRA (BioProject accession numbers, PRJNA666685 and

260 PRJNA666664). All HIV-1 test datasets are available on NCBI SRA (accession numbers:

261 ERR3953696, ERR3953853, ERR3953893, ERR3953891, ERR3953866, ERR3953846,

262 ERR3953756, ERR3953877, ERR3953876, ERR3953750, ERR3953741, ERR3953697,

263 ERR3953699, ERR3953706, ERR3953708, ERR3953710, ERR3953712, ERR3953716,

264 ERR3953295, ERR3953693).

265

266

267

268

269

270

271

272

273

274

275 **REFERENCES**

- 276 **Ajogbasile FV, Oguzie JU, Oluniyi PE, Eromon PE, Uwanibe JN, Mehta SB, Siddle KJ,**
277 **Odia I, Winnicki SM, Akpede N, Akpede G, Okogbenin S, Ogbaini-Emovon E,**
278 **MacInnis BL, Folarin OA, Modjarrad K, Schaffner SF, Tomori O, Ihekweazu C,**
279 **Sabeti PC, Happi CT. 2020.** Real-time Metagenomic Analysis of Undiagnosed Fever
280 Cases Unveils a Yellow Fever Outbreak in Edo State, Nigeria. *Scientific reports* 10:3180.
281 DOI: 10.1038/s41598-020-59880-w.
- 282
283 **Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990.** Basic local alignment search
284 tool. *Journal of molecular biology* 215:403–410. DOI: 10.1016/S0022-2836(05)80360-2.
- 285
286 **Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,**
287 **Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler**
288 **G, Alekseyev MA, Pevzner PA. 2012.** SPAdes: a new genome assembly algorithm and
289 its applications to single-cell sequencing. *Journal of computational biology: a journal of*
290 *computational molecular cell biology* 19:455–477. DOI: 10.1089/cmb.2012.0021.
- 291
292 **Bean AGD, Baker ML, Stewart CR, Cowled C, Deffrasnes C, Wang L-F, Lowenthal JW.**
293 **2013.** Studying immunity to zoonotic diseases in the natural host - keeping it real. *Nature*
294 *reviews. Immunology* 13:851–861. DOI: 10.1038/nri3551.
- 295
296 **Bolger AM, Lohse M, Usadel B. 2014.** Trimmomatic: a flexible trimmer for Illumina sequence
297 data. *Bioinformatics* 30:2114–2120. DOI: 10.1093/bioinformatics/btu170.
- 298
299 **Brister JR, Ako-Adjei D, Bao Y, Blinkova O. 2015.** NCBI viral genomes resource. *Nucleic*
300 *acids research* 43:D571–7. DOI: 10.1093/nar/gku1207.
- 301
302 **Cantalupo PG, Calgua B, Zhao G, Hundesa A, Wier AD, Katz JP, Grabe M, Hendrix RW,**
303 **Girones R, Wang D, Pipas JM. 2011.** Raw sewage harbors diverse viral populations.
304 *mBio* 2. DOI: 10.1128/mBio.00180-11.
- 305
306 **Chan PKS. 2002.** Outbreak of avian influenza A(H5N1) virus infection in Hong Kong in 1997.
307 *Clinical infectious diseases: an official publication of the Infectious Diseases Society of*
308 *America* 34 Suppl 2:S58–64. DOI: 10.1086/338820.
- 309
310 **Chen S, Zhou Y, Chen Y, Gu J. 2018.** fastp: an ultra-fast all-in-one FASTQ preprocessor.
311 *Bioinformatics* 34:i884–i890. DOI: 10.1093/bioinformatics/bty560.
- 312
313 **Chen N, Zhou M, Dong X, Qu J, Gong F, Han Y, Qiu Y, Wang J, Liu Y, Wei Y, Xia J ’an,**

- 314 **Yu T, Zhang X, Zhang L. 2020.** Epidemiological and clinical characteristics of 99 cases
315 of 2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study. *The Lancet*
316 395:507–513. DOI: 10.1016/S0140-6736(20)30211-7.
317
- 318 **Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T,**
319 **Kauff F, Wilczynski B, de Hoon MJL. 2009.** Biopython: freely available Python tools
320 for computational molecular biology and bioinformatics. *Bioinformatics* 25:1422–1423.
321 DOI: 10.1093/bioinformatics/btp163.
322
- 323 **Folarin OA, Ehichioya D, Schaffner SF, Winnicki SM, Wohl S, Eromon P, West KL,**
324 **Gladden-Young A, Oyejide NE, Matranga CB, Deme AB, James A, Tomkins-Tinch**
325 **C, Onyewurunwa K, Ladner JT, Palacios G, Nosamiefan I, Andersen KG, Omilabu**
326 **S, Park DJ, Yozwiak NL, Nasidi A, Garry RF, Tomori O, Sabeti PC, Happi CT.**
327 **2016.** Ebola Virus Epidemiology and Evolution in Nigeria. *The Journal of infectious*
328 *diseases* 214:S102–S109. DOI: 10.1093/infdis/jiw190.
329
- 330 **Grubaugh ND, Ladner JT, Kraemer MUG, Dudas G, Tan AL, Gangavarapu K, Wiley**
331 **MR, White S, Thézé J, Magnani DM, Prieto K, Reyes D, Bingham AM, Paul LM,**
332 **Robles-Sikisaka R, Oliveira G, Pronty D, Barcellona CM, Metsky HC, Baniecki**
333 **ML, Barnes KG, Chak B, Freije CA, Gladden-Young A, Gnirke A, Luo C,**
334 **MacInnis B, Matranga CB, Park DJ, Qu J, Schaffner SF, Tomkins-Tinch C, West**
335 **KL, Winnicki SM, Wohl S, Yozwiak NL, Quick J, Fauver JR, Khan K, Brent SE,**
336 **Reiner RC Jr, Lichtenberger PN, Ricciardi MJ, Bailey VK, Watkins DI, Cone MR,**
337 **Kopp EW 4th, Hogan KN, Cannons AC, Jean R, Monaghan AJ, Garry RF, Loman**
338 **NJ, Faria NR, Porcelli MC, Vasquez C, Nagle ER, Cummings DAT, Stanek D,**
339 **Rambaut A, Sanchez-Lockhart M, Sabeti PC, Gillis LD, Michael SF, Bedford T,**
340 **Pybus OG, Isern S, Palacios G, Andersen KG. 2017.** Genomic epidemiology reveals
341 multiple introductions of Zika virus into the United States. *Nature* 546:401–405. DOI:
342 10.1038/nature22400.
343
- 344 **Holshue ML, DeBolt C, Lindquist S, Lofy KH, Wiesman J, Bruce H, Spitters C, Ericson**
345 **K, Wilkerson S, Tural A, Diaz G, Cohn A, Fox L, Patel A, Gerber SI, Kim L, Tong**
346 **S, Lu X, Lindstrom S, Pallansch MA, Weldon WC, Biggs HM, Uyeki TM, Pillai SK,**
347 **Washington State 2019-nCoV Case Investigation Team. 2020.** First Case of 2019
348 Novel Coronavirus in the United States. *The New England journal of medicine.* DOI:
349 10.1056/NEJMoa2001191.
350
- 351 **Hunt M, Gall A, Ong SH, Brener J, Ferns B, Goulder P, Nastouli E, Keane JA, Kellam P,**
352 **Otto TD. 2015.** IVA: accurate de novo assembly of RNA virus genomes. *Bioinformatics*
353 31:2374–2376. DOI: 10.1093/bioinformatics/btv120.
354
- 355 **Katoh K, Misawa K, Kuma K-I, Miyata T. 2002.** MAFFT: a novel method for rapid multiple

- 356 sequence alignment based on fast Fourier transform. *Nucleic acids research* 30:3059–
357 3066. DOI: 10.1093/nar/gkf436.
- 358
- 359 **Kokot M, Dlugosz M, Deorowicz S. 2017.** KMC 3: counting and manipulating k-mer
360 statistics. *Bioinformatics* 33:2759–2761. DOI: 10.1093/bioinformatics/btx304.
- 361
- 362 **Köster J, Rahmann S. 2012.** Snakemake--a scalable bioinformatics workflow engine.
363 *Bioinformatics* 28:2520–2522. DOI: 10.1093/bioinformatics/bts480.
- 364
- 365 **Kurtzer GM, Sochat V, Bauer MW. 2017.** Singularity: Scientific containers for mobility of
366 compute. *PLoS one* 12:e0177459. DOI: 10.1371/journal.pone.0177459.
- 367
- 368 **Langmead B. 2010.** Aligning short sequencing reads with Bowtie. *Current protocols in*
369 *bioinformatics / editorial board, Andreas D. Baxevanis... [et al.]* Chapter 11:Unit 11.7.
370 DOI: 10.1002/0471250953.bi1107s32.
- 371
- 372 **Li H, Durbin R. 2009.** Fast and accurate short read alignment with Burrows-Wheeler
373 transform. *Bioinformatics* 25:1754–1760. DOI: 10.1093/bioinformatics/btp324.
- 374
- 375 **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,**
376 **Durbin R, 1000 Genome Project Data Processing Subgroup. 2009.** The Sequence
377 Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079. DOI:
378 10.1093/bioinformatics/btp352.
- 379
- 380 **Marçais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A. 2018.** MUMmer4:
381 A fast and versatile genome alignment system. *PLoS computational biology*
382 14:e1005944. DOI: 10.1371/journal.pcbi.1005944.
- 383
- 384 **Metsky HC, Matranga CB, Wohl S, Schaffner SF, Freije CA, Winnicki SM, West K, Qu J,**
385 **Baniecki ML, Gladden-Young A, Lin AE, Tomkins-Tinch CH, Ye SH, Park DJ, Luo**
386 **CY, Barnes KG, Shah RR, Chak B, Barbosa-Lima G, Delatorre E, Vieira YR, Paul**
387 **LM, Tan AL, Barcellona CM, Porcelli MC, Vasquez C, Cannons AC, Cone MR,**
388 **Hogan KN, Kopp EW, Anzinger JJ, Garcia KF, Parham LA, Ramirez RMG,**
389 **Montoya MCM, Rojas DP, Brown CM, Hennigan S, Sabina B, Scotland S,**
390 **Gangavarapu K, Grubaugh ND, Oliveira G, Robles-Sikisaka R, Rambaut A,**
391 **Gehrke L, Smole S, Halloran ME, Villar L, Mattar S, Lorenzana I, Cerbino-Neto J,**
392 **Valim C, Degraeve W, Bozza PT, Gnirke A, Andersen KG, Isern S, Michael SF,**
393 **Bozza FA, Souza TML, Bosch I, Yozwiak NL, MacInnis BL, Sabeti PC. 2017.** Zika
394 virus evolution and spread in the Americas. *Nature* 546:411–415. DOI:
395 10.1038/nature22402.

- 396
397 **Mokili JL, Rohwer F, Dutilh BE. 2012.** Metagenomics and future perspectives in virus
398 discovery. *Current opinion in virology* 2:63–77. DOI: 10.1016/j.coviro.2011.12.004.
399
- 400 **Park DJ, Dudas G, Wohl S, Goba A, Whitmer SLM, Andersen KG, Sealfon RS, Ladner
401 JT, Kugelman JR, Matranga CB, Winnicki SM, Qu J, Gire SK, Gladden-Young A,
402 Jalloh S, Nosamiefan D, Yozwiak NL, Moses LM, Jiang P-P, Lin AE, Schaffner SF,
403 Bird B, Towner J, Mamoh M, Gbakie M, Kanneh L, Kargbo D, Massally JLB,
404 Kamara FK, Konuwa E, Sellu J, Jalloh AA, Mustapha I, Foday M, Yillah M,
405 Erickson BR, Sealy T, Blau D, Paddock C, Brault A, Amman B, Basile J, Bearden S,
406 Belser J, Bergeron E, Campbell S, Chakrabarti A, Dodd K, Flint M, Gibbons A,
407 Goodman C, Klena J, McMullan L, Morgan L, Russell B, Salzer J, Sanchez A,
408 Wang D, Jungreis I, Tomkins-Tinch C, Kislyuk A, Lin MF, Chapman S, MacInnis
409 B, Matthews A, Bochicchio J, Hensley LE, Kuhn JH, Nusbaum C, Schieffelin JS,
410 Birren BW, Forget M, Nichol ST, Palacios GF, Ndiaye D, Happi C, Gevao SM,
411 Vandi MA, Kargbo B, Holmes EC, Bedford T, Gnirke A, Ströher U, Rambaut A,
412 Garry RF, Sabeti PC. 2015.** Ebola Virus Epidemiology, Transmission, and Evolution
413 during Seven Months in Sierra Leone. *Cell* 161:1516–1526. DOI:
414 10.1016/j.cell.2015.06.007.
- 415
- 416 **Pickett BE, Sadat EL, Zhang Y, Noronha JM, Squires RB, Hunt V, Liu M, Kumar S,
417 Zaremba S, Gu Z, Zhou L, Larson CN, Dietrich J, Klem EB, Scheuermann RH.
418 2012.** ViPR: an open bioinformatics database and analysis resource for virology research.
419 *Nucleic acids research* 40:D593–8. DOI: 10.1093/nar/gkr859.
- 420
- 421 **Reyes A, Haynes M, Hanson N, Angly FE, Heath AC, Rohwer F, Gordon JI. 2010.** Viruses
422 in the faecal microbiota of monozygotic twins and their mothers. *Nature* 466:334–338.
423 DOI: 10.1038/nature09199.
- 424
- 425 **Sharma D, Priyadarshini P, Vрати S. 2015.** Unraveling the web of viroinformatics:
426 computational tools and databases in virus research. *Journal of virology* 89:1489–1501.
427 DOI: 10.1128/JVI.02027-14.
428
- 429 **Siddle KJ, Eromon P, Barnes KG, Mehta S, Oguzie JU, Odia I, Schaffner SF, Winnicki
430 SM, Shah RR, Qu J, Wohl S, Brehio P, Iruolagbe C, Aiyepada J, Uyigue E,
431 Akhilomen P, Okonofua G, Ye S, Kayode T, Ajogbasile F, Uwanibe J, Gaye A,
432 Momoh M, Chak B, Kotliar D, Carter A, Gladden-Young A, Freije CA, Omoregie
433 O, Osiemi B, Muoebonam EB, Airende M, Enigbe R, Ebo B, Nosamiefan I, Oluniyi
434 P, Nekoui M, Ogbaini-Emovon E, Garry RF, Andersen KG, Park DJ, Yozwiak NL,
435 Akpede G, Ihekweazu C, Tomori O, Okogbenin S, Folarin OA, Okokhere PO,
436 MacInnis BL, Sabeti PC, Happi CT. 2018.** Genomic Analysis of Lassa Virus during an

- 437 Increase in Cases in Nigeria in 2018. *The New England journal of medicine* 379:1745–
438 1753. DOI: 10.1056/NEJMoa1804498.
- 439
- 440 **Sohrabi C, Alsafi Z, O’Neill N, Khan M, Kerwan A, Al-Jabir A, Iosifidis C, Agha R. 2020.**
441 World Health Organization declares Global Emergency: A review of the 2019 Novel
442 Coronavirus (COVID-19). *International journal of surgery* . DOI:
443 10.1016/j.ijssu.2020.02.034.
- 444
- 445 **Tang P, Chiu C. 2010.** Metagenomics for the discovery of novel human viruses. *Future*
446 *microbiology* 5:177–189. DOI: 10.2217/fmb.09.120.
- 447
- 448 **Wymant C, Blanquart F, Golubchik T, Gall A, Bakker M, Bezemer D, Croucher NJ, Hall**
449 **M, Hillebregt M, Ong SH, Ratmann O, Albert J, Bannert N, Fellay J, Fransen K,**
450 **Gourlay A, Grabowski MK, Günsenheimer-Bartmeyer B, Günthard HF, Kivelä P,**
451 **Kouyos R, Laeyendecker O, Liitsola K, Meyer L, Porter K, Ristola M, van Sighem**
452 **A, Berkhout B, Cornelissen M, Kellam P, Reiss P, Fraser C, BEEHIVE**
453 **Collaboration. 2018.** Easy and accurate reconstruction of whole HIV genomes from
454 short-read sequence data with shiver. *Virus evolution* 4:vey007. DOI: 10.1093/ve/vey007.

Table 1 (on next page)

Comparison between Lassa virus whole genomes (S segment) obtained with VGEA and viral-ngs

- 1 **Table 1:** Comparison between Lassa virus whole genomes (S segment) obtained with **viral-ngs**
- 2 and **VGEA**

S/N	Sample ID	Lassa Virus Genome Segment	Genome length with viral-ngs	Genome length with VGEA
1.	758	S	3392	3413
2.	852	S	3391	3413
3.	934	S	3394	3413
4.	1004	S	3392	3413
5.	1078	S	3382	3414
6.	1126	S	3394	3413
7.	A4	S	3360	3413
8.	A7	S	3380	3413
9.	J1	S	3374	3413
10.	K7	S	3389	3412
11.	0202C	S	3148	3413
12.	1177	S	3389	3413
13.	1801	S	3379	3412
14.	1880	S	3386	3413
15.	540	S	3294	3406
16.	D2	S	3394	3413
17.	F8	S	3391	3413
18.	J4	S	3386	3413
19.	O1	S	3393	3413

20.	O2	S	3381	3413
-----	----	---	------	------

3
4
5
6
7
8

Table 2 (on next page)

Comparison between Lassa virus whole genomes (L segment) obtained with VGEA and viral-ngs

- 1 **Table 2:** Comparison between Lassa virus whole genomes (L segment) obtained with **viral-ngs**
 2 and **VGEA**
 3

S/N	Sample ID	Lassa Virus Genome Segment	Genome length with viral-ngs	Genome length with VGEA
1.	758	L	7225	7271
2.	852	L	7239	7268
3.	934	L	7245	7272
4.	1004	L	7235	7269
5.	1078	L	7225	7272
6.	1126	L	7234	7270
7.	A4	L	7150	7271
8.	A7	L	7214	7273
9.	J1	L	7154	7271
10.	K7	L	7187	7272
11.	0202C	L	6418	7272
12.	1177	L	7121	7271
13.	1801	L	7221	7269
14.	1880	L	7195	7272
15.	540	L	7020	7225
16.	D2	L	7237	7271
17.	F8	L	7230	7272
18.	J4	L	7224	7273

19.	O1	L	7246	7273
20.	O2	L	7224	7275

4

Table 3 (on next page)

Comparison between SARS-CoV-2 whole genomes obtained with VGEA and viral-ngs

- 1 **Table 3:** Comparison between SARS-CoV-2 whole genomes obtained with **viral-ngs** and
 2 **VGEA**
 3

S/N	Sample ID	Genome length with viral-ngs	Genome length with VGEA
1.	CV18	29858	29903
2.	CV29	29865	29986
3.	CV43	29875	29903
4.	CV45	29897	30002
5.	CV47	29900	29903
6.	CV48	29859	30145
7.	CV50	29897	29903
8.	CV55	29898	30636
9.	CV57	29895	30129
10.	CV115	29859	30180
11.	CV145	28898	30250
12.	CV153	29848	29903
13.	CV155	29870	29903
14.	CV156	29856	29903
15.	CV163	29864	30356
16.	CV165	29872	29958
17.	CV167	29894	29900
18.	CV170	29897	29903
19.	CV185	29870	29903
20.	CV192	29897	29903

4

Table 4 (on next page)

Comparison between HIV-1 whole genomes obtained with VGEA and viral-ngs

1 **Table 4:** Comparison between HIV-1 whole genomes obtained with **viral-ngs** and **VGEA**

2

S/N	SRA Accession Number	Genome length with viral-ngs	Genome length with VGEA
1.	ERR3953696	8927	9877
2.	ERR3953853	9054	9872
3.	ERR3953893	9083	9880
4.	ERR3953891	7827	9818
5.	ERR3953866	9088	9872
6.	ERR3953846	9148	9917
7.	ERR3953756	8963	9913
8.	ERR3953877	9035	9802
9.	ERR3953876	7317	9824
10.	ERR3953750	5532	9860
11.	ERR3953741	7376	9822
12.	ERR3953697	6728	9798
13.	ERR3953699	9054	9837
14.	ERR3953706	7317	9826
15.	ERR3953708	7393	9839
16.	ERR3953710	6705	9808
17.	ERR3953712	9108	9846
18.	ERR3953716	9134	9835
19.	ERR3953295	8676	9882
20.	ERR3953693	8885	9857

3