



# THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Genome sequencing reveals Zika virus diversity and spread in the Americas

**Citation for published version:**

Metsky, HC, Matranga, CB, Wohl, S, Schaffner, SF, Freije, CA, Winnicki, S, West, K, Qu, J, Baniecki, ML, Gladden-Young, A, Lin, AE, Christopher, T-T, Ye, SH, Park, DJ, Luo, C, Barnes, KG, Shah, RR, Chak, B, Barbosa-Lima, G, Delatorre, E, Vieira, YR, Paul, LM, Tan, AL, Barcellona, CM, Porcelli, MC, Vasquez, C, Cannons, AC, Cone, MR, Hogan, KN, Kopp IV, EW, Anzinger, JJ, Garcia, KF, Parhap, LA, Gelvez Ramirez, RM, Montoya, M, Rojas, DP, Brown, CM, Hennigan, S, Sabina, B, Scotland, S, Gangavarapu, K, Grubaugh, ND, Oliveira, G, Robles-Sikisaka, R, Rambaut, A, Gehrke, L, Smole, S, Halloran, ME, Villar Centeno, LA, Mattar, S, Lorenzana, I, Cerbino-Neto, J, Valim, C, Degraeve, W, Bozza, PT, Souza, TML, Bosch, I, Yozwiak, NL, MacInnis, BL & Sabeti, PC 2017, 'Genome sequencing reveals Zika virus diversity and spread in the Americas', *Nature*, vol. 546, pp. 411–415. <https://doi.org/10.1038/nature22402>

**Digital Object Identifier (DOI):**

[10.1038/nature22402](https://doi.org/10.1038/nature22402)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

Nature

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Genome sequencing reveals Zika virus diversity and spread in the Americas

Metsky, H.C.<sup>\*1,2</sup>, Matranga, C.B.<sup>\*1</sup>, Wohl, S.<sup>\*1,3</sup>, Schaffner, S.F.<sup>\*1,3,4</sup>, Freije, C.A.<sup>1,3</sup>, Winnicki, S.M.<sup>1</sup>, West, K.<sup>1</sup>, Qu, J.<sup>1</sup>, Baniecki, M.L.<sup>1</sup>, Gladden-Young, A.<sup>1</sup>, Lin, A.E.<sup>1,3</sup>, Tomkins-Tinch, C.H.<sup>1</sup>, Ye, S.H.<sup>1,5</sup>, Park, D.J.<sup>1</sup>, Luo, C.Y.<sup>1,3</sup>, Barnes, K.G.<sup>1,3</sup>, Shah, R.R.<sup>1,6</sup>, Chak, B.<sup>1,3</sup>, Barbosa-Lima, G.<sup>7</sup>, Delatorre, E.<sup>8</sup>, Vieira, Y.R.<sup>7</sup>, Paul, L.M.<sup>9</sup>, Tan, A.L.<sup>9</sup>, Barcellona, C.M.<sup>9</sup>, Porcelli, M.C.<sup>10</sup>, Vasquez, C.<sup>10</sup>, Cannons, A.C.<sup>11</sup>, Cone, M.R.<sup>11</sup>, Hogan, K.N.<sup>11</sup>, Kopp, E.W. IV<sup>11</sup>, Anzinger, J.J.<sup>12</sup>, Garcia, K.F.<sup>13</sup>, Parham, L.A.<sup>13</sup>, Gélvez Ramírez, R.M.<sup>14</sup>, Miranda Montoya, M.C.<sup>14</sup>, Rojas, D.P.<sup>15</sup>, Brown, C.M.<sup>16</sup>, Hennigan, S.<sup>16</sup>, Sabina, B.<sup>16</sup>, Scotland, S.<sup>16</sup>, Gangavarapu, K.<sup>17</sup>, Grubaugh, N.D.<sup>17</sup>, Oliveira, G.<sup>18</sup>, Robles-Sikisaka, R.<sup>17</sup>, Rambaut, A.<sup>19,20</sup>, Gehrke, L.<sup>21,22</sup>, Smole, S.<sup>16</sup>, Halloran, M.E.<sup>23,24</sup>, Villar Centeno, L.A.<sup>14</sup>, Mattar, S.<sup>25</sup>, Lorenzana, I.<sup>13</sup>, Cerbino-Neto, J.<sup>7</sup>, Valim, C.<sup>4,26</sup>, Degraeve, W.<sup>27</sup>, Bozza, P.T.<sup>28</sup>, Gnirke, A.<sup>1</sup>, Andersen, K.G.<sup>†17,18,29</sup>, Isern, S.<sup>†9</sup>, Michael, S.F.<sup>†9</sup>, Bozza, F.A.<sup>†7,30</sup>, Souza, T.M.L.<sup>¶†31,32</sup>, Bosch, I.<sup>†21</sup>, Yozwiak, N.L.<sup>†1,3</sup>, MacInnis, B.L.<sup>¶†1,4</sup>, Sabeti, P.C.<sup>†1,3,4,33</sup>

## Affiliations:

1. Broad Institute of MIT and Harvard, Cambridge MA, USA
2. Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge MA, USA
3. Center for Systems Biology, Department of Organismic and Evolutionary Biology, Harvard University, Cambridge MA, USA
4. Department of Immunology and Infectious Diseases, Harvard T.H. Chan School of Public Health, Harvard University, Boston MA, USA
5. Harvard-MIT Health Sciences and Technology, Massachusetts Institute of Technology, Cambridge MA, USA
6. Harvard University Extension School, Cambridge, MA, USA
7. National Institute of Infectious Diseases Evandro Chagas, Fundação Oswaldo Cruz (FIOCRUZ), Rio de Janeiro RJ, Brazil
8. Laboratório de AIDS e Imunologia Molecular, Instituto Oswaldo Cruz, FIOCRUZ, Rio de Janeiro RJ, Brazil
9. Department of Biological Sciences, College of Arts and Sciences, Florida Gulf Coast University, Fort Myers FL, USA
10. Miami-Dade County Mosquito Control, Miami FL, USA
11. Bureau of Public Health Laboratories, Division of Disease Control and Health Protection, Florida Department of Health, Tampa FL, USA
12. Department of Microbiology, The University of the West Indies, Mona, Kingston, Jamaica
13. Instituto de Investigacion en Microbiologia, Universidad Nacional Autónoma de Honduras, Honduras
14. Grupo de Epidemiología Clínica, Universidad Industrial de Santander, Bucaramanga, Colombia
15. Department of Epidemiology, College of Public Health and Health Professions, University of Florida, Gainesville FL, USA
16. Massachusetts Department of Public Health, Jamaica Plain MA, USA
17. Department of Immunology and Microbial Science, The Scripps Research Institute, La Jolla CA, USA
18. Scripps Translational Science Institute, La Jolla CA, USA
19. University of Edinburgh, Edinburgh, UK
20. Fogarty International Center, National Institutes of Health, Bethesda MD, USA
21. Institute for Medical Engineering and Science, Massachusetts Institute of Technology, Cambridge MA, USA
22. Department of Microbiology and Immunobiology, Harvard Medical School, Boston MA, USA
23. Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle WA, USA
24. Department of Biostatistics, University of Washington, Seattle WA, USA
25. Institute for Tropical Biology Research, Universidad de Córdoba, Colombia
26. Department of Osteopathic Medical Specialties, Michigan State University, East Lansing MI, USA
27. Fiocruz, Instituto Oswaldo Cruz, Laboratório de Genômica Funcional e Bioinformática, Rio de Janeiro RJ, Brazil
28. Laboratório de Imunofarmacologia, Instituto Oswaldo Cruz, Fundação Oswaldo Cruz, Rio de Janeiro RJ, Brazil
29. Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla CA, USA
30. D'Or Institute for Research and Education, Brazil
31. National Institute for Science and Technology on Innovation on Neglected Diseases, Fiocruz, Rio de Janeiro RJ, Brazil
32. Center for Technological Development in Health, Fiocruz, Rio de Janeiro RJ, Brazil
33. Howard Hughes Medical Institute, Chevy Chase MD, USA

\* co-first author

† co-senior author

¶ co-corresponding author: B.L.M. (bronwyn@broadinstitute.org); T.M.L.S. (tmoreno@cchts.fiocruz.br)

55  
56 **Despite great attention given to the recent Zika virus (ZIKV) epidemic in the Americas and its link**  
57 **to birth defects<sup>1,2</sup>, much remains unknown about ZIKV disease epidemiology and ZIKV evolution,**  
58 **in part due to a lack of genomic data. We applied multiple sequencing approaches to generate 110**  
59 **ZIKV genomes from clinical and mosquito samples from 10 countries and territories, greatly**  
60 **expanding the observed viral genetic diversity from this outbreak. We analyzed the timing and**  
61 **patterns of introductions into distinct geographic regions; our phylogenetic evidence suggests rapid**  
62 **expansion of the outbreak in Brazil and multiple introductions of outbreak strains into Puerto Rico,**  
63 **Honduras, Colombia, other Caribbean islands, and the continental US. We find that ZIKV**  
64 **circulated undetected in multiple regions for many months before the first locally transmitted cases**  
65 **were confirmed, highlighting the importance of viral surveillance. We identify mutations with**  
66 **possible functional implications for ZIKV biology and pathogenesis, as well as those potentially**  
67 **relevant to the effectiveness of diagnostic tests.**  
68

69 Since its introduction into the Americas, mosquito-borne ZIKV (Family: *Flaviviridae*) has spread rapidly,  
70 causing hundreds of thousands of cases of ZIKV disease, as well as ZIKV congenital syndrome and likely  
71 other neurological complications<sup>1-3</sup>. Phylogenetic analysis of ZIKV can reveal the trajectory of the  
72 outbreak and detect mutations that may be associated with new disease phenotypes or affect molecular  
73 diagnostics. Despite the 70 years since its discovery and the scale of the recent outbreak, however, fewer  
74 than 100 ZIKV genomes have been sequenced directly from clinical samples. This is due in part to  
75 technical challenges posed by low viral loads (for example, often orders of magnitude lower than in Ebola  
76 virus or dengue virus infection<sup>4-6</sup>), and by loss of RNA integrity in samples collected and stored without  
77 sequencing in mind. Culturing the virus increases the material available for sequencing but can result in  
78 genetic variation that is not representative of the original clinical sample.  
79

80 We sought to gain a deeper understanding of the viral populations underpinning the ZIKV epidemic by  
81 extensive genome sequencing of the virus directly from samples collected as part of ongoing surveillance.  
82 We initially pursued unbiased metagenomic RNA sequencing to capture both ZIKV and other viruses  
83 known to be co-circulating with ZIKV<sup>5</sup>. In most of the 38 samples examined by this approach there  
84 proved to be insufficient ZIKV RNA for genome assembly, but it still proved valuable to verify results  
85 from other methods. Metagenomic data also revealed RNA from other viruses, including 41 likely novel  
86 viral sequence fragments in mosquito pools (**Extended Data Table 1**). In one patient we detected no  
87 ZIKV sequence but did assemble a complete genome from dengue virus (type 1), one of the viruses that  
88 co-circulates with and presents similarly to ZIKV<sup>7</sup>.  
89

90 To capture sufficient ZIKV content for genome assembly, we turned to two targeted approaches for  
91 enrichment before sequencing: multiplex PCR amplification<sup>8</sup> and hybrid capture<sup>9</sup>. We sequenced and  
92 assembled complete or partial genomes from 110 samples from across the epidemic, out of 229 attempted  
93 (221 clinical samples from confirmed and possible ZIKV disease cases and eight mosquito pools; **Table**  
94 **1, Supplementary Table 1**). This dataset, which we used for further analysis, includes 110 genomes  
95 produced using multiplex PCR amplification (amplicon sequencing) and a subset of 37 genomes  
96 produced using hybrid capture (out of 66 attempted). Because these approaches amplify any contaminant  
97 ZIKV content, we relied heavily on negative controls to detect artefactual sequence, and we established  
98 stringent, method-specific thresholds on coverage and completeness for calling high confidence ZIKV

99 assemblies (**Fig. 1a**). Completeness and coverage for these genomes are shown in **Fig. 1b and c**; the  
100 median fraction of the genome with unambiguous base calls was 93%. Per-base discordance between  
101 genomes produced by the two methods was 0.017% across the genome, 0.15% at polymorphic positions,  
102 and 2.2% for minor allele base calls. Concordance of within-sample variants is shown in more detail in  
103 **Fig. 1d-f**. Patient sample type (urine, serum, or plasma) made no significant difference in sequencing  
104 success in our study (**Extended Data Fig. 1**).

105  
106 To investigate the spread of ZIKV in the Americas we performed a phylogenetic analysis of the 110  
107 genomes from our dataset, together with 64 published genomes available on NCBI GenBank and in our  
108 companion papers<sup>10,11</sup> (**Fig. 2a**). Our reconstructed phylogeny (**Fig. 2b**), which is based on a molecular  
109 clock (**Extended Data Fig. 2**), is consistent with the outbreak originating in Brazil<sup>12</sup>: Brazil ZIKV  
110 genomes appear on all deep branches of the tree, and their most recent common ancestor is the root of the  
111 entire tree. We estimate the date of that common ancestor to have been in early 2014 (95% credible  
112 interval, CI, August 2013 to July 2014). The shape of the tree near the root remains uncertain (i.e. the  
113 nodes have low posterior probabilities) because there are too few mutations to clearly distinguish the  
114 branches. This pattern suggests rapid early spread of the outbreak, consistent with the introduction of a  
115 new virus to an immunologically naive population. ZIKV genomes from Colombia ( $n=10$ ), Honduras  
116 ( $n=18$ ), and Puerto Rico ( $n=3$ ) cluster within distinct, well-supported clades. We also observed a clade  
117 consisting entirely of genomes from patients who contracted ZIKV in one of three Caribbean countries  
118 (the Dominican Republic, Jamaica, and Haiti) or the continental US, containing 30 of 32 genomes from  
119 the Dominican Republic and 19 of 20 from the continental US. We estimated the within-outbreak  
120 substitution rate to be  $1.15 \times 10^{-3}$  substitutions/site/year (95% CI [ $9.78 \times 10^{-4}$ ,  $1.33 \times 10^{-3}$ ]), similar to prior  
121 estimates for this outbreak<sup>12</sup>. This is somewhat higher (1.3x–5x) than reported rates for other  
122 flaviviruses<sup>13</sup>, but is measured over a short sampling period, and therefore may include a higher  
123 proportion of mildly deleterious mutations that have not yet been removed through purifying selection.

124  
125 Determining when ZIKV arrived in specific regions helps elucidate the spread of the outbreak and track  
126 rising incidence of possible complications of ZIKV infection. The majority of the ZIKV genomes from  
127 our study fall into four major clades from different geographic regions, for which we estimated a likely  
128 date for ZIKV arrival. In each case, the date was months earlier than the first confirmed, locally  
129 transmitted case, indicating ongoing local circulation of ZIKV before its detection. In Puerto Rico, the  
130 estimated date was 4.5 months earlier than the first confirmed local case<sup>14</sup>; it was 8 months earlier in  
131 Honduras<sup>15</sup>, 5.5 months earlier in Colombia<sup>16</sup>, and 9 months earlier for the Caribbean/continental US  
132 clade<sup>17</sup>. In each case, the arrival date represents the estimated time to the most recent common ancestor  
133 (tMRCA) for the corresponding clade in our phylogeny (**Fig. 2c**). See **Extended Data Fig. 3** and  
134 **Extended Data Table 2** for details. Similar temporal gaps between the tMRCA of local transmission  
135 chains and the earliest detected cases were seen when chikungunya virus emerged in the Americas<sup>18</sup>. We  
136 also observed evidence for several introductions of ZIKV into the continental US, and found that  
137 sequences from mosquito and human samples collected in Florida cluster together, consistent with the  
138 finding of local ZIKV transmission in Florida in a companion paper<sup>11</sup>.

139  
140 Principal component analysis (PCA) is consistent with the phylogenetic observations (**Fig. 2d**). It shows  
141 tight clustering among ZIKV genomes from the continental US, the Dominican Republic, and Jamaica.  
142 ZIKV genomes from Brazil and Colombia are similar and distinct from genomes sampled in other

143 countries. ZIKV genomes from Honduras form a third cluster that also contains genomes from Guatemala  
144 or El Salvador. The PCA results show no clear stratification of ZIKV within Brazil.

145  
146 Genetic variation can provide important clues to understanding ZIKV biology and pathogenesis and can  
147 reveal potentially functional changes in the virus. We observed 1030 single nucleotide polymorphisms  
148 (SNPs) in the complete dataset, well distributed across the genome (**Fig. 3a**). Any effect of these  
149 mutations cannot be determined from these data; however, the most likely candidates for functional  
150 mutations would be among the 202 nonsynonymous SNPs (**Supplementary Table 2**) and the 32 SNPs in  
151 the 5' and 3' untranslated regions (UTRs). Adaptive mutations are more likely to be found at high  
152 frequency or to be seen multiple times, although both effects can also occur by chance. We observed five  
153 positions with nonsynonymous mutations at >5% minor allele frequency that occur on two or more  
154 branches of the tree (**Fig. 3b**); two of these (at 4287 and 8991) occur together and might represent  
155 incorrect placement of a Brazil branch in the tree. The remaining three are more likely to represent  
156 multiple nonsynonymous mutations; one (at 9240) appears to involve nonsynonymous mutations to two  
157 different alleles.

158  
159 To assess the possible biological significance of these mutations, we looked for evidence of selection in  
160 the ZIKV genome. Viral surface glycoproteins are known targets of positive selection, and mutations in  
161 these proteins can confer adaptation to new vectors<sup>19</sup> or aid immune escape<sup>20,21</sup>. We therefore searched for  
162 an excess of nonsynonymous mutations in the ZIKV envelope glycoprotein (E). However, the  
163 nonsynonymous substitution rate in E proved to be similar to that in the rest of the coding region (**Fig. 3c**,  
164 left); moreover, amino acid changes were significantly more conservative in that region than elsewhere  
165 (**Fig. 3c**, middle and right). Any diversifying selection occurring in the surface protein thus appears to be  
166 operating under selective constraint. We also found evidence for purifying selection in the ZIKV 3' UTR  
167 (**Fig. 3d, Supplementary Table 3**), a region important for viral replication<sup>22</sup>.

168  
169 While the transition-to-transversion ratio (6.98) was within the range seen in other viruses<sup>23</sup>, we observed  
170 a significantly higher frequency of C-to-T and T-to-C substitutions than other transitions (**Fig. 3d**,  
171 **Extended Data Fig. 4, Supplementary Table 3**). This enrichment is apparent both in the genome as a  
172 whole and at 4-fold degenerate sites, where selection pressure is minimal. Many processes may contribute  
173 to this conspicuous mutation pattern, including mutational bias of the ZIKV RNA-dependent RNA  
174 polymerase, host RNA editing enzymes (e.g. APOBECs, ADARs) acting upon viral RNA, and chemical  
175 deamination, but further investigation is required to determine the cause of this phenomenon.

176  
177 Mismatches between PCR assays and viral sequence are a potential source of poor diagnostic  
178 performance in this outbreak<sup>24</sup>. To assess the potential impact of ongoing viral evolution on diagnostic  
179 function, we compared eight published qRT-PCR-based primer/probe sets to our data. We found  
180 numerous sites where the probe or primer did not match an allele found among the 174 ZIKV genomes  
181 from the current dataset (**Fig. 3e**). In most cases, the discordant allele was shared by all outbreak samples,  
182 presumably because it was present in the Asian lineage that entered the Americas. These mismatches  
183 could affect all uses of the diagnostic assay in the outbreak. We also found mismatches from new  
184 mutations that occurred following ZIKV entry into the Americas. Most of these were present in less than  
185 10% of samples, although one was seen in 29%. These observations suggest that genome evolution has

186 not caused widespread degradation of diagnostic performance during the course of the outbreak, but that  
187 mutations continue to accumulate and ongoing monitoring is needed.

188  
189 Analysis of within-host viral genetic diversity can reveal important information for understanding virus-  
190 host interactions and viral transmission. However, accurately identifying these variants in low-titer  
191 clinical samples is challenging, and further complicated by potential artefacts associated with enrichment  
192 prior to sequencing. To investigate whether we could reliably detect within-host ZIKV variants in our  
193 data, we identified within-host variants in a cultured ZIKV isolate used as a positive control throughout  
194 our study, and found that both amplicon sequencing and hybrid capture data produced concordant and  
195 replicable variant calls (**Fig. 1d**). In clinical samples, hybrid capture within-host variants were noisier but  
196 contained a reliable subset: although most variants were not validated by the other sequencing method or  
197 by a technical replicate, those at high frequency were always replicable, as were those that passed a  
198 previously described filter<sup>25</sup> (**Fig. 1e-f, Extended Data Table 3**). Within this high confidence set we  
199 looked for variants shared between samples as a clue to transmission patterns, but there were too few  
200 variants to draw any meaningful conclusions. By contrast, within-host variants identified in amplicon  
201 sequencing data were unreliable at all frequencies (**Fig. 1f, Extended Data Table 3**), suggesting that  
202 further technical development is needed before amplicon sequencing can be used to study within-host  
203 variation in ZIKV and other clinical samples with low viral titer.

204  
205 Sequencing low titer viruses like ZIKV directly from clinical samples presents several challenges that  
206 have likely contributed to the paucity of genomes available from the current outbreak. While development  
207 of technical and analytical methods will surely continue, we note that factors upstream in the process,  
208 including collection site and cohort, were strong predictors of sequencing success in our study (**Extended**  
209 **Data Fig. 1**). This highlights the importance of continuing development and implementation of best  
210 practices for sample handling, without disrupting standard clinical workflows, for wider adoption of  
211 genome surveillance during outbreaks. Additional sequencing, however challenging, remains critical to  
212 ongoing investigation of ZIKV biology and pathogenesis. Together with two companion studies<sup>10,11</sup>, this  
213 effort advances both technological and collaborative strategies for genome surveillance in the face of  
214 unexpected outbreak challenges.

215

216 **Author Contributions**  
217 C.B.M., S.W., C.A.F., S.M.W., K.W., J.Q., M.L.B., A.G.-Y., C.Y.L., R.R.S., G.B.-L., Y.R.V., L.M.P., A.L.T., C.M.B., M.C.P.,  
218 C.Vasquez., A.C.C., M.R.C., K.N.H., E.W.K.IV, J.J.A., K.F.G., L.A.P., R.M.G.R., M.C.M.M., C.M.B., S.H., B.S., S.Scotland.,  
219 K.G., G.O., R.R.-S., and I.B. performed laboratory experiments and prepared samples for sequencing. H.C.M., C.B.M., C.A.F.,  
220 S.M.W., K.W., J.Q., M.L.B., C.Y.L., A.G.-Y., N.G.D., A.G., and K.G.A. developed methods for ZIKV detection, targeted  
221 enrichment, and/or sequencing library preparation. H.C.M., C.B.M., S.W., S.F.S., M.L.B., A.E.L., C.H.T.-T., S.Y., D.J.P., E.D.,  
222 A.R., T.M.L.S., I.B., and B.L.M. performed sequence assembly, curation, and/or data analyses. S.Smole., L.A.V.C., S.M., I.L.,  
223 S.I., S.F.M., and F.A.B. led clinical studies and/or study sites. K.G.B., B.C., D.P.R., N.D.G., L.G., M.E.H., A.R., A.G., J.C.-N.,  
224 C.Valim., W.G., P.T.B., A.G., K.G.A., S.I., S.F.M., F.A.B., T.M.L.S., and I.B. provided critical insights and guidance. H.C.M.,  
225 C.B.M., T.M.L.S., N.L.Y., B.L.M., and P.C.S. oversaw study design and management. H.C.M., C.B.M., S.W., S.F.S., A.E.L.,  
226 N.L.Y., B.L.M. and P.C.S. drafted the manuscript. All authors reviewed the manuscript.

227  
228 **Acknowledgements**  
229 We are grateful for the vision and support of Marc and Lynne Benioff, and for the support and guidance of Liliana Brown, Eun  
230 Mi Lee, and Maria Giovanni (NIAID), and Justine Levin-Allerhand and Eric S. Lander (Broad Institute). We thank Molly  
231 Schleicher, Emily Lipscomb, August Felix, Andrea Saltzman, and Stacey Donnelly (Broad Institute) for assistance with IRB and  
232 Ethics processes; Eva Mair and Liisa Nogelo (Broad Institute), and Erika Carmean (HHMI) for legal counsel; Tamara Mason and  
233 the Broad Institute Genomics Platform for sequencing support; Ashley Matthews, Sinéad Chapman, Daniel Neafsey, Bruce  
234 Birren (Broad Institute) for management and guidance; Oliver Pybus (Oxford University) and ZiBRA Project colleagues for  
235 sharing data prior to publication; Daniel Olson and Edwin Asturias (Children’s Hospital Colorado), Marc Salit (NIST), and  
236 Etienne Simon-Loriere (Pasteur Institute) for sharing samples and reagents; and Edward Holmes (University of Sydney), Gonzalo  
237 Bello (Fiocruz), Ryan Tewhey, Anne Piantadosi, Chris Edwards and the Sabeti Lab (Broad Institute) for helpful discussions and  
238 critical reading of the manuscript. We are indebted to Zika patients and clinical teams for making this work possible.

239  
240 **Funding**  
241 Funding was provided by: Marc and Lynne Benioff (P.C.S.); NIH NIAID U19AI110818 (Broad Institute); Howard Hughes  
242 Medical Institute (P.C.S.); Broad Institute BroadNext10 program (A.G. and P.C.S.); AWS Cloud Credits for Research (P.C.S.);  
243 Conselho Nacional de Desenvolvimento Científico e Tecnológico (440909/2016-3) and Fundação de Amparo a Pesquisa do  
244 Estado do Rio de Janeiro (E-26/201.320/2016, E-26/201.332/2016, E-26/010.000194/2015) (P.T.B. and F.A.B.); NIH NIAID  
245 1R01AI099210 (S.I. and S.F.M.); MIDAS-National Institute of General Medical Sciences U54GM111274 (M.E.H.); NIH NIAID  
246 AI100190 (I.B. and L.G.); AEDES Network (I.B.) and Colombian Science, Technology and Innovation Fund of Sistema General  
247 de Regalías-BPIN 2013000100011 (L.A.V., R.M.G., M.C.M.M., and I.B.); ASTMH Shope Fellowship (K.G.B.); NSF DGE  
248 1144152 (A.E.L.); PNPD/CAPES Postdoctoral Fellowship (E.D.); Fulbright-Colciencias Doctoral Scholarship (D.P.R.); NIH  
249 training grant 5T32AI007244-33 (N.D.G.); EU under grant agreements 278433-PREDEMICS and 643476-COMPARE (A.R.);  
250 and NIH NCATS CTSA UL1TR001114, NIH NIAID contract HHSN272201400048C, The Ray Thomas Foundation, and Pew  
251 Biomedical Scholarship (K.G.A.).

252  
253 **Competing financial interests**  
254 The authors declare no competing financial interests.

255 **Figures**

256

Country or territory	Samples	Samples with metagenomic data	Amplicon sequencing genomes	Hybrid capture genomes	Total genomes
Brazil	53	12	27	7	<b>27</b>
Colombia	20	0	4	2	<b>4</b>
Dominican Republic	45	7	30	9	<b>30</b>
Guatemala/El Salvador	3	0	1	0	<b>1</b>
Haiti	4	0	1	0	<b>1</b>
Honduras	20	6	18	8	<b>18</b>
Jamaica	20	0	5	0	<b>5</b>
Martinique	3	0	1	0	<b>1</b>
Puerto Rico	15	0	3	1	<b>3</b>
Continental US	36	12	20	10	<b>20</b>
Other	10	1	0	0	<b>0</b>
<b>Total</b>	<b>229</b>	<b>38</b>	<b>110</b>	<b>37</b>	<b>110</b>

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

**Table 1 | Samples and genomes by region.** Sample source information and sequencing results for 229 clinical and mosquito pool samples. Continental US includes 8 mosquito pool samples; all others are clinical samples. In the final column, genomes generated by both methods are counted only once. “Other” includes regions without a ZIKV genome included in downstream analysis.

**Figure 1 | Sequence data from clinical and mosquito samples.** (a) Thresholds used to select samples for downstream analysis. Each point is a replicate. Red and blue shading: regions of accepted amplicon sequencing and hybrid capture genome assemblies, respectively. Not shown: hybrid capture positive controls with depth >10,000x. (b) Amplicon sequencing coverage by sample (row) across the ZIKV genome. Red: sequencing depth  $\geq 100x$ ; heatmap (bottom) sums coverage across all samples. White horizontal lines: amplicon locations. (c) Relative sequencing depth across hybrid capture genomes. (d) Within-sample variants for a single cultured isolate (PE243) across seven technical replicates. Each point is a variant in a replicate identified using amplicon sequencing (red) or hybrid capture (blue). Variants are plotted if the pooled frequency across replicates by either method is  $\geq 1\%$ . (e) Within-sample variant frequencies across methods. Each point is a variant in an individual sample and points are plotted on a log-log scale. Green points: “verified” variants detected by hybrid capture that pass strand bias and frequency filters. (f) Counts of within-sample variants across two replicate libraries, for each method. Variants are plotted in the frequency bin corresponding to the higher of the two detected frequencies. In (e-f), frequencies <1% are shown at 0%.

**Figure 2 | Zika virus spread throughout the Americas.** (a) Samples were collected in each of the colored countries/territories. Specific state, department, or province of origin for samples in this study are highlighted if known. (b) Maximum clade credibility tree. Dotted tips: genomes generated in this study. Node labels are posterior probabilities indicating support for the node. Violin plots denote probability distributions for the tMRCA of four highlighted clades. (c) Time elapsed between estimated tMRCA and date of first confirmed, locally transmitted case. Color: distributions based on relaxed clock model (also shown in (b)); grey: strict clock. “Caribbean” includes the continental US. (d) Principal component analysis of variants. Circles: data generated in this study; diamonds: other publicly available genomes from this outbreak. Percentage of variance explained by each component is indicated on axis.

**Figure 3 | Geographic and genomic distribution of Zika virus variation.** (a) Location of variants in the ZIKV genome. The minor allele frequency is the proportion of the 174 genomes from this outbreak that share a variant. Dotted bars: <25% of samples had a base call at that position. (b) Phylogenetic distribution of nonsynonymous variants with minor allele frequency  $\geq 5\%$ , shown on the branch where the mutation most likely occurred. Grey outline:



294 variant might be on next-most ancestral branch (in two cases, 2 branches upstream), but exact location is unclear  
295 because of missing data. Red circles: variants occurring at more than one location in the tree. **(c)** Conservation of the  
296 ZIKV envelope (E) region. Left: nonsynonymous variants per amino acid for the E region (dark grey) and the rest of  
297 the coding region (light grey). Middle: proportion of nonsynonymous variants resulting in negative BLOSUM62 scores,  
298 which indicate unlikely or extreme substitutions ( $p < 0.039$ ,  $\chi^2$  test). Right: average of BLOSUM62 scores for  
299 nonsynonymous variants ( $p < 0.037$ , 2-sample  $t$ -test). **(d)** Constraint in the ZIKV 3' UTR and observed transition rates  
300 over the ZIKV genome. **(e)** ZIKV diversity in diagnostic primer and probe regions. Top: locations of published probes  
301 (dark blue) and primers (cyan)<sup>26-31</sup> on the ZIKV genome. Bottom: each column represents a nucleotide position in the  
302 probe or primer. Colors in the column indicate the fraction of ZIKV genomes (out of 174) that match the probe/primer  
303 sequence (grey), differ from it (red), or have no data for that position (white).

## 304 **Methods**

305

### 306 **Ethics statement**

307 The clinical studies from which samples were obtained were evaluated and approved by relevant Institutional  
308 Review Boards/Ethics Review Committees at: Hospital General de la Plaza de la Salud (Santo Domingo, Dominican  
309 Republic), University of the West Indies (Kingston, Jamaica), Universidad Nacional Autónoma de Honduras  
310 (Tegucigalpa, Honduras), Oswaldo Cruz Foundation (Rio de Janeiro, Brazil), Centro de Investigaciones  
311 Epidemiológicas - Universidad Industrial de Santander (Bucaramanga, Colombia), Massachusetts Department of  
312 Public Health (Jamaica Plain, Massachusetts), and Florida Department of Health (Tallahassee, Florida). Informed  
313 consent was obtained from all participants enrolled in studies at Hospital General de la Plaza de la Salud,  
314 Universidad Nacional Autónoma de Honduras, Oswaldo Cruz Foundation, and Universidad Industrial de Santander.  
315 IRBs at the University of West Indies, Massachusetts Department of Public Health, and Florida Department of  
316 Health granted waivers of consent given this research with leftover clinical diagnostic samples involved no more  
317 than minimal risk. Harvard University and Massachusetts Institute of Technology (MIT) Institutional Review  
318 Boards/Ethics Review Committees provided approval for sequencing and secondary analysis of samples collected  
319 by the aforementioned institutions.

320

### 321 **Sample collections and study subjects**

322 Suspected ZIKV cases (including high-risk travelers) were enrolled through study protocols at multiple  
323 aforementioned collection sites. Clinical samples (including blood, urine, cerebrospinal fluid, and saliva) were  
324 obtained from suspected or confirmed ZIKV cases and from high-risk travelers. De-identified information about  
325 study participants and other sample metadata are reported in **Supplementary Table 1**.

326

### 327 **Viral RNA isolation**

328 RNA was isolated following manufacturer's standard operating protocol for 0.14 mL up to 1 mL samples<sup>32</sup> using the  
329 QIAamp Viral RNA Minikit (Qiagen), except that in some cases 0.1 M final concentration of  $\beta$ -mercaptoethanol (as  
330 a reducing agent) or 40  $\mu$ g/mL final concentration of linear acrylamide (Ambion) (as a carrier) were added to AVL  
331 buffer prior to inactivation. Extracted RNA was resuspended in AVE buffer or nuclease-free water. In some cases,  
332 viral samples were concentrated using Vivaspin-500 centrifugal concentrators (Sigma-Aldrich) prior to inactivation  
333 and extraction. In these cases, 0.84 mL of sample was concentrated to 0.14 mL by passing through a 30 kDa filter  
334 and discarding the flow through.

335

### 336 **Carrier RNA and host rRNA depletion**

337 In a subset of human samples, carrier poly(rA) RNA and host rRNA were depleted from RNA samples using RNase  
338 H selective depletion<sup>9,33</sup>. Briefly, oligo d(T) (40 nt long) and/or DNA probes complementary to human rRNA were  
339 hybridized to the sample RNA. The sample was then treated with 15 units of Hybridase Thermostable RNase H  
340 (Epicentre) for 30 minutes at 45°C. The complementary DNA probes were removed by treating each reaction with  
341 an RNase-free DNase (Qiagen) according to the manufacturer's protocol. Following depletion, samples were  
342 purified using 1.8x volume AMPure RNAClean beads (Beckman Coulter Genomics) and eluted into 10  $\mu$ l water for  
343 cDNA synthesis.

344

### 345 **Illumina library construction and sequencing**

346 cDNA synthesis was performed as described in previously published RNA-seq methods<sup>9</sup>. To track potential cross-  
347 contamination, 50 fg of synthetic RNA (gift from M. Salit, NIST) was spiked into samples using unique RNA for  
348 each individual ZIKV sample. ZIKV negative control cDNA libraries were prepared from water, human K-562 total  
349 RNA (Ambion), or EBOV (KY425633.1) seed stock; ZIKV positive controls were prepared from ZIKV Senegal  
350 (isolate HD78788) or ZIKV Pernambuco (isolate PE243; KX197192.1) seed stock. The dual index Accel-NGS® 2S  
351 Plus DNA Library Kit (Swift Biosciences) was used for library preparation. Approximately half of the cDNA

352 product was used for library construction, and indexed libraries were generated using 18 cycles of PCR. Each  
353 individual sample was indexed with a unique barcode. Libraries were pooled at equal molarity and sequenced on the  
354 Illumina HiSeq 2500 or MiSeq (paired-end reads) platforms.

355

### 356 **Amplicon-based cDNA synthesis and library construction**

357 ZIKV amplicons were prepared as described<sup>8,11</sup>, similarly to “RNA jackhammering” for preparing low input viral  
358 samples for sequencing<sup>34</sup>, with slight modifications. After PCR amplification, each amplicon pool was quantified on  
359 a 2200 TapeStation (Agilent Technologies) using High Sensitivity D1000 ScreenTape (Agilent Technologies). 2  $\mu$ L  
360 of a 1:10 dilution of the amplicon cDNA was loaded and the concentration of the 350-550 bp fragments was  
361 calculated. The cDNA concentration, as reported by the TapeStation, was highly predictive of sequencing outcome  
362 (i.e. whether a sample passes genome assembly thresholds) (**Extended Data Fig. 5**). cDNA from each of the two  
363 amplicon pools were mixed equally (10-25 ng each) and libraries were prepared using the dual index Accel-NGS®  
364 2S Plus DNA Library Kit (Swift Biosciences) according to manufacturer's protocol. Libraries were indexed with a  
365 unique barcode using 7 cycles of PCR, pooled equally and sequenced on the Illumina MiSeq (250 bp paired-end  
366 reads) platform. Primer sequences were removed by hard trimming the first 30 bases for each insert read prior to  
367 analysis.

368

### 369 **Zika virus hybrid capture**

370 Viral hybrid capture was performed as previously described<sup>9</sup>. Probes were created to target ZIKV and chikungunya  
371 virus (CHIKV). Candidate probes were created by tiling across publicly available sequences for ZIKV and CHIKV  
372 on NCBI GenBank<sup>35</sup>. Probes were selected from among these candidate probes to minimize the number used while  
373 maintaining coverage of the observed diversity of the viruses. Alternating universal adapters were added to allow  
374 two separate PCR amplifications, each consisting of non-overlapping probes. (To download probe sequences, see  
375 Supplementary Information.)

376

377 The probes were synthesized on a 12k array (CustomArray). The synthesized oligos were amplified by two separate  
378 emulsion PCR reactions with primers containing T7 RNA polymerase promoter. Biotinylated baits were in vitro  
379 transcribed (MEGAscript, Ambion) and added to prepared ZIKV libraries. The baits and libraries were  
380 hybridized overnight (~16 hrs), captured on streptavidin beads, washed, and re-amplified by PCR using the Illumina  
381 adapter sequences. Capture libraries were then pooled and sequenced. In some cases, a second round of hybrid  
382 capture was performed on PCR-amplified capture libraries to further enrich the ZIKV content of sequencing  
383 libraries (**Extended Data Fig. 6**). In the main text, “hybrid capture” refers to a combination of hybrid capture  
384 sequencing data and data from the same libraries without capture (unbiased), unless explicitly distinguished.

385

### 386 **Genome assembly**

387 We assembled reads from all sequencing methods into genomes using viral-ngs v1.13.3<sup>36,37</sup>. We taxonomically  
388 filtered reads from amplicon sequencing against a ZIKV reference, KU321639.1. We filtered reads from other  
389 approaches against the list of accessions provided in Supplementary Information. To compute results on individual  
390 replicates, we *de novo* assembled these and scaffolded against KU321639.1. To obtain final genomes for analysis,  
391 we pooled data from multiple replicates of a sample, *de novo* assembled, and scaffolded against KX197192.1. For  
392 all assemblies, we set the viral-ngs ‘assembly\_min\_length\_fraction\_of\_reference’ and ‘assembly\_min\_unambig’  
393 parameters to 0.01. For amplicon sequencing data, unambiguous base calls required at least 90% of reads to agree in  
394 order to call that allele (‘major\_cutoff’ = 0.9); for hybrid capture data, we used the default threshold of 50%. We  
395 modified viral-ngs so that calls to GATK’s UnifiedGenotyper set ‘min\_indel\_count\_for\_genotyping’ to 2.

396

397 At 3 sites with insertions or deletions (indels) in the consensus genome CDS, we corrected the genome using Sanger  
398 sequencing of the RT-PCR product (namely, at 3447 in the genome for sample DOM\_2016\_BB-0085-SER; at 5469  
399 in BRA\_2016\_FC-DQ12D1-PLA; and at 6516-6564 in BRA\_2016\_FC-DQ107D1-URI, with coordinates in  
400 KX197192.1). At other indels in the consensus genome CDS, we replaced the indel with ambiguity.

401  
402 Depth of coverage values from amplicon sequencing include read duplicates. In all other cases, we removed  
403 duplicates with viral-ngs.

#### 405 **Identification of non-ZIKV viruses in samples by unbiased sequencing**

406 Using Kraken v0.10.6<sup>38</sup> in viral-ngs, we built a database that includes its default “full” database (which incorporates  
407 all bacterial and viral whole genomes from RefSeq<sup>39</sup> as of October 2015). Additionally, we included the whole  
408 human genome (hg38), genomes from PlasmoDB<sup>40</sup>, sequences covering mosquito genomes (*Aedes aegypti*, *Aedes*  
409 *albopictus*, *Anopheles albimanus*, *Anopheles quadrimaculatus*, *Culex quinquefasciatus*, and the outgroup  
410 *Drosophila melanogaster*) from GenBank<sup>35</sup>, protozoa and fungi whole genomes from RefSeq, SILVA LTP 16s  
411 rRNA sequences<sup>41</sup>, and all sequences from NCBI’s viral accession list<sup>42</sup> (as of October 2015) for viral taxa that have  
412 human as a host. (To download database, see Supplementary Information.)

413  
414 For each sample, we ran Kraken on data from unbiased sequencing replicates (not including hybrid capture data)  
415 and searched its output reports for viral taxa with more than 100 reported reads. We manually filtered the results,  
416 removing ZIKV, bacteriophages, and known lab contaminants. For each sample and its associated taxa, we  
417 assembled genomes using viral-ngs as described above; results are in **Extended Data Table 1a**. We used the  
418 following genomes for taxonomically filtering reads and as the reference for assembly: KJ741267.1 (cell fusing  
419 agent virus), AY292384.1 (deformed wing virus), NC\_001477.1 (dengue virus type 1), LC164349.1 (JC  
420 polyomavirus). When reporting sequence identity of an assembly to its taxon, we used BLASTN<sup>43</sup> to determine the  
421 identity between the sequence and the reference used for its assembly.

422  
423 To focus on metagenomics of mosquito pools (**Extended Data Table 1b**), we considered unbiased sequencing data  
424 from 8 mosquito pools (not including hybrid capture data). We first ran the depletion pipeline of viral-ngs on raw  
425 data and then ran the viral-ngs Trinity<sup>44</sup> assembly pipeline on the depleted reads to assemble them into contigs. We  
426 pooled contigs from all mosquito pool samples and identified all duplicate contigs with sequence identity >95%  
427 using CD-HIT<sup>45</sup>. Additionally, we used predicted coding sequences from Prodigal 2.6.3<sup>46</sup> to identify duplicate  
428 protein sequences at >95% identity. We classified contigs using BLASTN<sup>43</sup> against nt and BLASTX<sup>43</sup> against nr (as  
429 of February 2017) and discarded all contigs with an e-value greater than 1E-4. We define viral contigs as contigs  
430 that hit a viral sequence, and we manually removed all reverse-transcriptase-like contigs due to their similarity to  
431 retrotransposon elements within the *Aedes aegypti* genome. We categorized viral contigs with less than 80% amino  
432 acid identity to their best hit as likely novel viral contigs. **Supplementary Table 4** lists the unique viral contigs we  
433 found, their best hit, and information scoring the hit.

#### 435 **Relationship between metadata and sequencing outcome**

436 To determine if available sample metadata are predictive of sequencing outcome, we tested the following variables:  
437 sample collection site, patient gender, patient age, sample type, and the number of days between symptom onset and  
438 sample collection (“collection interval”). To describe sequencing outcome of a sample  $S$ , we used the following  
439 response variable  $Y_S$ :

440  $\text{mean}(\{ I(R) * (\text{number of unambiguous bases in } R) \text{ for all amplicon sequencing replicates } R \text{ of } S \}),$   
441 where  $I(R)=1$  if median depth of coverage of  $R \geq 275$  and  $I(R)=0$  otherwise

442 This value is listed in **Supplementary Table 1** under “Dependent variable used in regression on metadata”. We  
443 excluded the saliva, cerebrospinal fluid, and whole blood sample types due to sample number ( $n=1$ ), and also  
444 excluded mosquito pool samples and rows with missing values. We excluded samples from one collection site  
445 (prefix “JAM\_2016\_WI-”) because most had missing values. We treated samples with type “Plasma EDTA” as  
446 having type “Plasma”. We treated the “collection interval” variable as categorical (0-1, 2-3, 4-6, and 7+ days).

447  
448 With a single model we underfit the zero counts, possibly because many zeros (samples without a replicate that  
449 passes ZIKV assembly) are truly ZIKV-negative. We thus view the data as coming from two processes: one

450 determining whether a sample is ZIKV-positive or ZIKV-negative, and another that determines, among the observed  
451 passing samples, how much of a ZIKV genome we are able to sequence. We modeled the first process, predicting  
452 whether a sample is passing, with logistic regression (in R using GLM<sup>47</sup> with binomial family and logit link); here,  
453 the observed passing samples are the samples  $S$  for which  $Y_S \geq 2500$ . For the second, we performed a beta  
454 regression, using only the observed passing samples, of  $Y_S$  divided by ZIKV genome length on the predictor  
455 variables. We implemented this in R using the betareg package<sup>48</sup> and transformed fractions from the closed unit  
456 interval to the open unit interval as the authors suggest.

457  
458 To test the significance of predictor variables, we used a likelihood ratio test. For variable  $X_i$  we compared a full  
459 model (with all predictors) against a model that uses all predictors except  $X_i$ . Results of these tests are shown in  
460 **Extended Data Fig. 1a and d**. We explore the effects of sample type and collection interval on obtaining a passing  
461 assembly in **Extended Data Fig. 1b and c**, respectively. Error bars are 95% confidence intervals derived from  
462 binomial distributions. We explore the effects of these same two variables on  $Y_S$  (in passing samples only) in  
463 **Extended Data Fig. 1e and f**.

#### 464 465 **Criteria for pooling across replicates**

466 We attempted to sequence one or more replicates of each sample and attempted to assemble a genome from each  
467 replicate. We discarded data from any replicates whose assembly showed high sequence similarity, in any part of the  
468 genome, to our assembly of the genome in a sample consisting of an African (Senegal) lineage (strain HD78788) of  
469 ZIKV. We used this sample as a positive control throughout this study, and considered its presence in the assembly  
470 of a clinical or mosquito pool sample to be evidence of contamination. Similarly, we discarded data from four  
471 replicates belonging to samples from the Dominican Republic because they yielded assemblies that were  
472 unexpectedly identical or highly similar to our assembly of the ZIKV isolate PE243 genome, another positive  
473 control used in this study. We also discarded data from replicates that showed evidence of contamination, at the  
474 RNA stage, by the baits used in hybrid capture; we detected these by looking for adapters that were added to these  
475 probes for amplification.

476  
477 For amplicon sequencing, we consider an assembly of a replicate to be “passing” if it contains at least 2500  
478 unambiguous base calls and has a median depth of coverage of at least 275x over its unambiguous bases (depth  
479 includes duplicate reads). For the unbiased and hybrid capture approaches, we consider an assembly of a replicate  
480 “passing” if it contains at least 4000 unambiguous base calls. For each approach, the unambiguous base threshold is  
481 based on an observed density of negative controls below the threshold (**Fig. 1a**). For amplicon sequencing  
482 assemblies, we added a coverage depth threshold because coverage depth was roughly binary across replicates, with  
483 negative controls falling in the lower class. Based on these thresholds, 0 of 99 negative controls used throughout our  
484 sequencing runs yield passing assemblies and 32 of 32 positive controls yield passing assemblies.

485  
486 We consider a sample to have a passing assembly if any of its replicates, by either method, yields an assembly that  
487 passes the above thresholds. For each sample with at least one passing assembly, we pooled read data across  
488 replicates for each sample, including replicates with assemblies that do not pass the assembly thresholds. When data  
489 was available from both amplicon sequencing and unbiased/hybrid capture approaches, we pooled amplicon  
490 sequencing data separately from data produced by the unbiased and hybrid capture approaches, the latter two of  
491 which were pooled together (henceforth, the “hybrid capture” pool). We then assembled a genome from each set of  
492 pooled data. When assemblies on pooled data were available from both approaches, we selected for downstream  
493 analysis the assembly from the hybrid capture approach if it had more than 10267 unambiguous base calls (95% of  
494 the reference genome used, GenBank accession KX197192.1); when this condition was not met, we selected the one  
495 that had more unambiguous base calls.

496  
497 The number of ZIKV genomes publicly available prior to this study is the result of an NCBI GenBank<sup>35</sup> search for  
498 ZIKV in February 2017. We filtered any sequences with length <4000 nt, excluded sequences that are being

499 published as part of this study or a companion paper<sup>10,11</sup>, excluded sequences from non-human hosts, and excluded  
500 sequences labeled as having been passaged. We counted fewer than 100 sequences, the precise number depending  
501 on details of the count.

502

### 503 **Visualization of coverage depth across genomes**

504 For amplicon sequencing data, we plotted coverage across the 110 samples that yielded a passing assembly by  
505 amplicon sequencing (**Fig. 1b**). With viral-ngs, we aligned depleted reads to the reference sequence KX197192.1  
506 using the novoalign aligner with options ‘-r Random -l 40 -g 40 -x 20 -t 100 -k’. Because of the nature of amplicon  
507 sequencing, duplicates were not identified or removed. We binarized depth at each nucleotide position, showing red  
508 if depth of coverage is at least 100x. Rows (samples) are hierarchically clustered to ease visualization.

509

510 For hybrid capture sequencing data, we plotted depth of coverage across the 37 samples that yielded a passing  
511 assembly (**Fig. 1c**). We aligned reads as described above for amplicon sequencing data, except we removed  
512 duplicates. For each sample, we calculated depth of coverage at each nucleotide position. We then scaled the values  
513 for each sample so that each would have a mean depth of 1.0. At each nucleotide position, we calculated the median  
514 depth across the samples, as well as the 20<sup>th</sup> and 80<sup>th</sup> percentiles. We plotted the mean of each of these metrics  
515 within a 200 nt sliding window.

516

### 517 **Multiple sequence alignments**

518 We aligned ZIKV consensus genomes using MAFFT v7.221<sup>49</sup> with the following parameters: ‘--maxiterate 1000 --  
519 ep 0.123 --localpair’.

520

521 In Supplementary Data, we provide sequences and alignments used in analyses.

522

### 523 **Analysis of within- and between-sample variants**

524 To measure overall per-base discordance between consensus genomes produced by amplicon sequencing and hybrid  
525 capture, we considered all sites where base calls were made in both the amplicon sequencing and hybrid capture  
526 consensus genomes of a sample, and we calculated the fraction in which the bases were not in agreement. To  
527 measure discordance at polymorphic sites, we took all of the consensus genomes generated in this study that we  
528 selected for downstream analysis and searched for positions with polymorphism (see **Criteria for pooling across**  
529 **replicates** for choosing among the amplicon sequencing and hybrid capture genome when both are available). We  
530 then looked at these positions in genomes that were available from both methods, and we calculated the fraction in  
531 which the alleles were not in agreement.

532

533 To measure discordance at minor alleles, we took all of the consensus genomes generated in this study that we  
534 selected for downstream analysis and searched for minor alleles. We then looked at all sites at which there was a  
535 minor allele and for which genomes from both methods were available, and we calculated the fraction in which the  
536 alleles were not in agreement. For these calculations, we tolerated partial ambiguity (e.g. ‘Y’ is concordant with  
537 ‘T’). If one genome had full ambiguity (‘N’) at a position and the other genome had an indel, we counted the site as  
538 discordant; otherwise, if one genome had full ambiguity, we did not count the site.

539

540 After assembling genomes, we determined within-sample allele frequencies for each sample by running V-Phaser  
541 2.0 via viral-ngs<sup>37</sup> on all pooled reads mapping to the sample assembly. When determining per-library allele counts  
542 at each variant position, we modified viral-ngs to require a minimum base (Phred) quality score of 30 for all bases,  
543 discard anomalous read pairs, and use per-base alignment quality (BAQ) in its calls to SAMtools<sup>50</sup> mpileup. This is  
544 particularly helpful for filtering spurious amplicon sequencing variants because all generated reads start and end at a  
545 limited number of positions (due to the pre-determined tiling of amplicons across the genome). Because amplicon  
546 sequencing libraries were sequenced using 250 bp paired-end reads, bases near the middle of the ~450 nt amplicons  
547 fall at the end of both paired reads, where quality scores drop and incorrect base calls are more likely. To determine

548 the overall frequency of each variant in a sample, we summed allele counts (calculated using SAMtools<sup>50</sup> mpileup  
549 via viral-ngs) across libraries.

550

551 When comparing variant frequencies between amplicon sequencing (7 technical replicates) and hybrid capture (7  
552 technical replicates) replicates of the PE243 positive control (**Fig. 1d**), we include only positions at which the mean  
553 (pooled) frequency across replicates within at least one method was  $\geq 1\%$ . When comparing allele frequencies  
554 between replicate libraries, we restricted the sample set to only samples with a passing assembly in both methods,  
555 and included only samples with two or more replicates. In contrast, when comparing alleles across methods we  
556 included samples that have a passing assembly by either method, with any number of replicates. For these  
557 comparisons, we only included positions with a minor variant; i.e. positions for which both libraries/methods had an  
558 allele at 100% were removed, even if the single allele differed between the two libraries/methods. Additionally, we  
559 considered any allele with frequency  $< 1\%$  as not found (0%).

560

561 When comparing allele frequencies across methods: let  $f_a$  and  $f_{hc}$  be frequencies in amplicon sequencing and hybrid  
562 capture, respectively. If both are non-zero, we only included an allele if the read depth at its position was  $\geq 1/\min(f_a,$   
563  $f_{hc})$  in both methods, and if depth at the position was at least 100 for hybrid capture and 275 for amplicon  
564 sequencing. If  $f_a=0$ , we required a read depth of  $\max(1/f_{hc}, 275)$  at the position in the amplicon sequencing method;  
565 similarly, if  $f_{hc}=0$  we required a read depth of  $\max(1/f_a, 100)$  at the position in the hybrid capture method. This was  
566 to eliminate lack of coverage as a reason for discrepancy between two methods. When comparing allele frequencies  
567 across sequencing replicates within a method, we imposed only a minimum read depth (275x for amplicon  
568 sequencing and 100x for hybrid capture), but required this depth in both libraries. In samples with more than two  
569 replicates, we only considered the two replicates with the highest depth at each plotted position.

570

571 We considered allele frequencies from hybrid capture sequencing “verified” if they passed the strand bias and  
572 frequency filters described in Gire et al. 2014<sup>25</sup>, with the exception that we imposed a minimum allele frequency of  
573 1% and allowed a variant identified in only one library if its frequency was  $\geq 5\%$ . In **Extended Data Table 3** and  
574 **Fig. 1f**, we considered variants “validated” if they were present at  $\geq 1\%$  frequency in both libraries or methods.

575 When comparing two libraries for a given method  $M$  (amplicon sequencing or hybrid capture): the proportion  
576 unvalidated is the fraction, among all variants in  $M$  at  $\geq 1\%$  frequency in at least one library, of the variants that are  
577 at  $\geq 1\%$  frequency in exactly one of the two libraries. Similarly, when comparing methods: the proportion  
578 unvalidated for a method  $M$  is the fraction, among all variants at  $\geq 1\%$  frequency in  $M$ , of the variants that are at  $\geq 1\%$   
579 frequency in  $M$  and  $< 1\%$  frequency in the other method.

580

581 We initially called SNPs on the aligned consensus genomes using Geneious version 9.1.7<sup>51</sup>. We converted all fully  
582 or partially ambiguous calls, which are treated by Geneious as variants, into missing data. We then removed all sites  
583 that were no longer polymorphic from the SNP set and re-calculated allele frequencies. A nonsynonymous SNP is  
584 shown on the tree (**Fig. 3b**) if it includes an allele that is nonsynonymous relative to the ancestral state (see  
585 **Molecular clock phylogenetics and ancestral state reconstruction** section below) and has a minor allele  
586 frequency of  $> 5\%$ ; all occurrences of nonsynonymous alleles are shown. (Two SNPs, at positions 2853 and 7229,  
587 had nominal derived allele frequencies over 95%; in both cases, the “ancestral” allele was seen only in a small clade  
588 within the tree, suggesting that the ancestral allele was incorrectly assigned.) We placed mutations at a node such  
589 that the node leads only to samples with the mutation or with no call at that site. Uncertainty in placement occurs  
590 when a sample lacks a base call for the corresponding SNP; in this case, we placed the SNP on the most recent  
591 branch for which we have available data. We also used this ancestral ZIKV state to count the frequency of each type  
592 of substitution over various regions of the ZIKV genome, per number of available bases in each region (**Fig. 3d** and  
593 **Supplementary Table 3**).

594

595 We quantified the effect of nonsynonymous SNPs using the original BLOSUM62 scoring matrix for amino acids<sup>52</sup>,  
596 in which positive scores indicate conservative amino acid changes and negative scores unlikely or extreme

597 substitutions. We assessed statistical significance for equality of proportions by  $\chi^2$  test (**Fig. 3c**, middle), and for  
598 difference of means by 2-sample *t*-test with Welch-Satterthwaite approximation of df (**Fig. 3c**, right). Error bars are  
599 95% confidence intervals derived from binomial distributions (**Fig. 3c**, left and middle; **Fig. 3d**) or Student's *t*-  
600 distributions (**Fig. 3c**, right).

601

#### 602 **Maximum likelihood estimation and root-to-tip regression**

603 We generated a maximum likelihood tree using a multiple sequence alignment that included genomes generated in  
604 this study, as well as a selection of other available sequences from the Americas, Southeast Asia, and the Pacific.  
605 The sequences are listed in Supplementary Information. We ran PhyML<sup>53</sup> with the GTR substitution model and 4  
606 gamma substitution rate categories; for the tree search operation, we used 'BEST' (best of NNI and SPR). In  
607 FigTree v1.4.2<sup>54</sup>, we rooted the tree on the oldest sequence used as input (GenBank accession EU545988.1).

608

609 We used TempEst v1.5<sup>55</sup>, which selects the best-fitting root with a residual mean squared function, to estimate root-  
610 to-tip distances. We performed regression in R with the lm function<sup>47</sup> of distances on dates. The relationship  
611 between root-to-tip divergence and sample dates (**Extended Data Fig. 2**) supports the use of a molecular clock  
612 analysis in this study.

613

614 In Supplementary Data, we provide the output of PhyML, as well as the dates and distances used for root-to-tip  
615 regression.

616

#### 617 **Molecular clock phylogenetics and ancestral state reconstruction**

618 For molecular clock phylogenetics, we made a multiple sequence alignment from the genomes generated in this  
619 study combined with a selection of other available sequences from the Americas. We did not use sequences from  
620 outside the outbreak in the Americas. Among ZIKV genomes published and publicly available on NCBI GenBank<sup>35</sup>,  
621 we selected 32 from the Americas that had at least 7000 unambiguous bases, were not labeled as having been  
622 passaged more than once, and had location metadata. We also used 32 genomes from Brazil published in a  
623 companion paper<sup>10</sup> that met the same criteria. The sequences are listed in Supplementary Information.

624

625 We used BEAST v1.8.4 to perform molecular clock analyses<sup>56</sup>. We used sampled tip dates to handle inexact dates<sup>57</sup>.  
626 Because of sparse data in non-coding regions, we used only the CDS as input. We used the SDR06 substitution  
627 model on the CDS, which uses HKY with gamma site heterogeneity and partitions codons into two partitions  
628 (positions (1+2) and 3)<sup>58</sup>. To perform model selection, we tested three coalescent tree priors: a constant-size  
629 population, an exponential growth population, and a Bayesian Skyline tree prior (10 groups, piecewise-constant  
630 model)<sup>59</sup>. For each tree prior, we tested two clock models: a strict clock and an uncorrelated relaxed clock with  
631 lognormal distribution (UCLN)<sup>60</sup>. In each case, we set the molecular clock rate to use a continuous time Markov  
632 chain rate reference prior<sup>61</sup>. For all six combinations of models, we performed path sampling (PS) and stepping-  
633 stone sampling (SS) to estimate marginal likelihood<sup>62,63</sup>. We sampled for 100 path steps with a chain length of 1  
634 million, with power posteriors determined from evenly spaced quantiles of a Beta(alpha=0.3; 1.0) distribution. The  
635 Skyline tree prior provided a better fit than the two other (baseline) tree priors (**Extended Data Table 2**), so we used  
636 this tree prior for all further analyses. Using a constant or exponential tree prior, a relaxed clock provides a better  
637 model fit, as shown by the log Bayes factor when comparing the two clock models. Using a Skyline tree prior, the  
638 log Bayes factor comparing a strict and relaxed clock is smaller than it is using the other tree priors, and it is similar  
639 to the variability between estimated log marginal likelihood from PS and SS methods. We chose to use a relaxed  
640 clock for further analyses, but we also report key findings using a strict clock.

641

642 For the tree and tMRCA estimates in **Fig. 2**, as well as the clock rate reported in main text, we ran BEAST with 400  
643 million MCMC steps using the SRD06 substitution model, Skyline tree prior, and relaxed clock model. We  
644 extracted clock rate and tMRCA estimates, and their distributions, with Tracer v1.6.0 and identified the maximum  
645 clade credibility (MCC) tree using TreeAnnotator v1.8.2. The reported credible intervals around estimates are 95%



646 highest posterior density (HPD) intervals. When reporting substitution rate from a relaxed clock model, we give the  
647 mean rate (mean of the rates of each branch weighted by the time length of the branch). Additionally, for the  
648 tMRCA estimates in **Fig. 2c** with a strict clock, we ran BEAST with the same specifications (also with 400M steps)  
649 except used a strict clock model. The resulting data are also used in the more comprehensive comparison shown in  
650 **Extended Data Fig. 3**.

651  
652 For the data with an outgroup in **Extended Data Fig. 3**, we ran BEAST the same as specified above (with strict and  
653 relaxed clock models), except with 100 million steps and with outgroup sequences in the input alignment. The  
654 outgroup sequences were the same as those used to make the maximum likelihood tree (see Supplementary  
655 Information). For the data excluding sample DOM\_2016\_MA-WGS16-020-SER in **Extended Data Fig. 3**, we ran  
656 BEAST the same as specified above (with strict and relaxed clocks), except we removed this sample from the input  
657 and ran 100 million steps.

658  
659 We used BEAST v1.8.4 to estimate transition and transversion rates with CDS and non-coding regions. The model  
660 was the same as above except that we used the Yang96 substitution model on the CDS, which uses GTR with  
661 gamma site heterogeneity and partitions codons into three partitions<sup>64</sup>; for the non-coding regions, we used a GTR  
662 substitution model with gamma site heterogeneity and no codon partitioning. There were four partitions in total: one  
663 for each codon position and another for the non-coding region (5' and 3' UTRs combined). We ran this for 200  
664 million steps. At each sampled step of the MCMC, we calculated substitution rates for each partition using the  
665 overall substitution rate, the relative substitution rate of the partition, the relative rates of substitutions in the  
666 partition, and base frequencies. In **Extended Data Fig. 4**, we plot the means of these rates over the steps; the error  
667 bars shown are 95% HPD intervals of the rates over the steps.

668  
669 We used BEAST v1.8.4 to reconstruct ancestral state at the root of the tree using CDS and non-coding regions. The  
670 model was the same as above except that, on the CDS, we used the HKY substitution model with gamma site  
671 heterogeneity and codons partitioned into three partitions (one per codon position). On the non-coding regions we  
672 used the same substitution model without codon partitioning. We ran this for 50 million steps and used  
673 TreeAnnotator v1.8.2 to find the state with the MCC tree. We selected the ancestral state corresponding to this state.

674  
675 In all BEAST runs, we discarded the first 10% of states from each run as burn-in.

676  
677 In Supplementary Data, we provide BEAST input (XML) and output files. We also provide the sequence of the  
678 reconstructed ancestral state.

679  
680 **Principal component analysis**  
681 We carried out principal component analysis using the R package FactoMineR<sup>65</sup>. We imputed missing data with the  
682 package missMDA<sup>66</sup> and we show the results in **Fig. 2d**.

683  
684 **Diagnostic assay assessment**  
685 We extracted primer and probe sequences from eight published RT-qPCR assays<sup>26-31</sup> and aligned to our ZIKV  
686 genomes using Geneious version 9.1.7<sup>51</sup>. We then tabulated matches and mismatches to the diagnostic sequence for  
687 all outbreak genomes, allowing multiple bases to match where the diagnostic primer and/or probe sequence  
688 contained nucleotide ambiguity codes (**Fig. 3e**).

689  
690 **Data availability**  
691 Sequence data that support findings of this study are deposited in NCBI GenBank<sup>35</sup> under BioProject accession  
692 PRJNA344504. Zika virus genomes have accession numbers KY014295-KY014327 and KY785409-KY785485.  
693 The dengue virus type 1 genome sequenced in this study has accession number KY829115. See **Supplementary**  
694 **Table 1** for a mapping of sample names to accession numbers.

## 695 Extended Data Figures

696

697 **Extended Data Figure 1 | Relationship between metadata and sequencing outcome.** Analysis of possible  
698 predictors of sequencing outcome: the site where a sample was collected, patient gender, patient age, sample type,  
699 and days between symptom onset and sample collection (“collection interval”). **(a)** Prediction of whether a sample  
700 passes assembly thresholds by sequencing. Rows show results of likelihood ratio tests on each predictor by omitting  
701 the variable from a full model that contains all predictors. Sample site and patient gender improve model fit, but  
702 sample type and collection interval do not. **(b)** Proportion of samples that pass assembly thresholds by sequencing,  
703 divided by sample type, across six sample sites. **(c)** Same as (b), except divided by collection interval. **(d)** Prediction  
704 of the genome fraction identified, using samples passing assembly thresholds. Rows show results of likelihood ratio  
705 tests, as in (a). Collection interval improves the model, but sample type does not. **(e)** Sequencing outcome for each  
706 sample, divided by sample type, across six sample sites. **(f)** Same as (e), except divided by collection interval.  
707 Samples collected 7+ days after symptom onset produced, on average, the fewest unambiguous bases, though these  
708 observations are based on a limited number of data points. While the sample site variable accounts for differences in  
709 cohort composition, the observed effects of gender and collection interval might be due to confounders in composition  
710 that span multiple cohorts. These results illustrate the effect of variables on sequencing outcome for the samples in  
711 this study; they are not indicative of ZIKV titer more generally. Other studies<sup>67,68</sup> have analyzed the impact of sample  
712 type and collection interval on ZIKV detection, sometimes with differing results.

713

714 **Extended Data Figure 2 | Maximum likelihood tree and root-to-tip regression.** **(a)** Tips are colored by sample  
715 source location. Labeled tips indicate those generated in this study; all other colored tips are other publicly available  
716 genomes from the outbreak in the Americas. Grey tips are samples from ZIKV cases in Southeast Asia and the  
717 Pacific. **(b)** Linear regression of root-to-tip divergence on dates. The substitution rate for the full tree, indicated by the  
718 slope of the black regression line, is similar to rates of Asian lineage ZIKV estimated by molecular clock analyses<sup>12</sup>.  
719 The substitution rate for sequences within the Americas outbreak only, indicated by the slope of the green regression  
720 line, is similar to rates estimated by BEAST ( $1.15 \times 10^{-3}$ ; 95% CI [ $9.78 \times 10^{-4}$ ,  $1.33 \times 10^{-3}$ ]) for this data set.

721

722 **Extended Data Figure 3 | Substitution rate and tMRCA distributions.** **(a)** Posterior density of the substitution rate.  
723 Shown with and without the use of sequences (outgroup) from outside the Americas. **(b-e)** Posterior density of the  
724 date of the most recent common ancestor (MRCA) of sequences in four regions corresponding to those in **Fig. 2c**.  
725 Shown with and without the use of outgroup sequences. The use of outgroup sequences has little effect on estimates  
726 of these dates. **(f)** Posterior density of the date of the MRCA of sequences in a clade consisting of samples from the  
727 Caribbean and continental US. Shown with and without the sequence of DOM\_2016\_MA-WGS16-020-SER, a  
728 sample from the Dominican Republic that has only 3037 unambiguous bases; this is the most ancestral sequence in  
729 the clade and its presence affects the tMRCA. In (a-f), all densities are shown as observed with a relaxed clock model  
730 and with a strict clock model.

731

732 **Extended Data Figure 4 | Substitution rates estimated with BEAST.** Substitution rates estimated in three codon  
733 positions and non-coding regions (5' and 3' UTRs). Transversions are shown in grey and transitions are colored by  
734 transition type. Plotted values show the mean of rates calculated at each sampled Markov chain Monte Carlo  
735 (MCMC) step of a BEAST run. These calculated rates provide additional evidence for the observed high C-to-T and  
736 T-to-C transition rates shown in **Fig. 3d**.

737

738 **Extended Data Figure 5 | cDNA concentration of amplicon primer pools predicts sequencing outcome.** cDNA  
739 concentration of amplicon pools (as measured by Agilent 2200 TapeStation) is highly predictive of amplicon  
740 sequencing outcome. On each axis, 1+primer pool concentration is plotted on a log scale. Each point is a technical  
741 replicate of a sample and colors denote observed sequencing outcome of the replicate. If a replicate is predicted to  
742 be passing when at least one primer pool concentration is  $\geq 0.8$  ng/ $\mu$ L, then sensitivity=98.71% and  
743 specificity=90.34%. An accurate predictor of sequencing success early in the sample processing workflow can save  
744 resources.

745

746 **Extended Data Figure 6 | Evaluating multiple rounds of Zika virus hybrid capture.** Genome assembly statistics  
747 of samples prior to hybrid capture (grey), and after one (blue) or two (red) rounds of hybrid capture. 9 individual

748 libraries (8 unique samples) were sequenced all three ways, had >1 million raw reads in each method, and generated  
749 at least one passing assembly. Raw reads from each method were downsampled to the same number of raw reads  
750 (8.5 million) before genomes were assembled. **(a)** Percent of the genome identified, as measured by number of  
751 unambiguous bases. **(b)** Median sequencing depth of ZIKV genomes, taken over the assembled regions.

752  
753 **Extended Data Table 1 | Viruses other than Zika uncovered by unbiased sequencing. (a)** Viral species other  
754 than Zika were found by unbiased sequencing of 38 samples. Column 3: number of reads in a sample belonging to a  
755 species as a raw count and a percent of total reads. Column 4: percent genome assembled based on the number of  
756 unambiguous bases called. We identified cell fusing agent virus (a flavivirus) and deformed wing virus-like genomes  
757 in mosquito pools, and dengue virus type 1, JC polyomavirus, and JC polyomavirus-like genomes in clinical samples.

758 All assemblies had  $\geq 95\%$  sequence identity to a reference sequence for the listed species, except cell fusing agent  
759 virus in USA\_2016\_FL-06-MOS (91%) and dengue virus type 1 in BLM\_2016\_MA-WGS16-006-SER (92%). The  
760 dengue virus type 1 genome showed  $\geq 95\%$  sequence identity to other available isolates of the virus. **(b)** Contigs  
761 assembled from unbiased sequencing data of 8 mosquito pools. Column 2: number of contigs assembled. Column 3:  
762 number of contigs classified by BLASTN/BLASTX<sup>43</sup>. Column 4: number of contigs hitting a viral species. Column 5:  
763 number of contigs hitting a viral species with <80% amino acid identity to the best hit. Each column is a subset of the  
764 previous column. Contigs in column 5 are considered to be likely novel. Last row lists counts, after removing duplicate  
765 contigs, for all mosquito pools combined. **Supplementary Table 4** lists the unique viral contigs and their best hit.

766  
767 **Extended Data Table 2 | Model selection for BEAST analyses. (a)** Marginal likelihoods calculated with path  
768 sampling (PS) and stepping-stone sampling (SS) for combinations of three coalescent tree priors (constant size  
769 population, exponential growth population, and Skyline) and two clock models (strict clock and uncorrelated relaxed  
770 clock with log-normal distribution). The Bayes factor is calculated against the baseline model, a constant size tree  
771 prior and strict clock. **(b)** Mean estimates and 95% credible intervals (CI) across evaluated models for the clock rate,  
772 date of tree root, and tMRCAs of the four regions shown in **Fig. 2c**. Under a Skyline tree prior, the use of strict and  
773 relaxed clock models yields similar estimates.

774  
775 **Extended Data Table 3 | Within-sample variant validation between and within sequencing methods. (a)** For  
776 each method (amplicon sequencing or hybrid capture), fraction of identified variants ( $\geq 1\%$ ) not identified at  $\geq 1\%$  by  
777 the other method (i.e. unvalidated). “Verified” hybrid capture variants are those passing strand bias and frequency  
778 filters, as described in Methods. **(b)** For each method, fraction of identified variants unvalidated in a second library. To  
779 pass the strand bias filter, a variant must meet filter criteria in both replicates.

780

## 781 **References**

- 782 1. *Zika situation report: Zika virus, Microcephaly and Guillain-Barré syndrome.* (World Health Organization,  
783 Feb 02 2017).
- 784 2. Reynolds, M. R. *et al.* Vital Signs: Update on Zika Virus-Associated Birth Defects and Evaluation of All U.S.  
785 Infants with Congenital Zika Virus Exposure - U.S. Zika Pregnancy Registry, 2016. *MMWR Morb. Mortal.*  
786 *Wkly. Rep.* **66**, 366–373 (2017).
- 787 3. de Vigilância em Saúde, S. *Protocolo de vigilância e resposta à ocorrência de microcefalia.* (Ministério da  
788 Saúde Brasília, 2016).
- 789 4. Schieffelin, J. S. *et al.* Clinical illness and outcomes in patients with Ebola in Sierra Leone. *N. Engl. J. Med.*  
790 **371**, 2092–2100 (2014).
- 791 5. Sardi, S. I. *et al.* Coinfections of Zika and Chikungunya Viruses in Bahia, Brazil, Identified by Metagenomic  
792 Next-Generation Sequencing. *J. Clin. Microbiol.* **54**, 2348–2353 (2016).
- 793 6. Martina, B. E. E., Koraka, P. & Osterhaus, A. D. M. E. Dengue virus pathogenesis: an integrated view. *Clin.*  
794 *Microbiol. Rev.* **22**, 564–581 (2009).
- 795 7. Fauci, A. S. & Morens, D. M. Zika Virus in the Americas — Yet Another Arbovirus Threat. *N. Engl. J. Med.*  
796 **374**, 601–604 (2016).
- 797 8. Quick, J. *et al.* Multiplex PCR method for MinION and Illumina sequencing of Zika and other virus genomes  
798 directly from clinical samples. *bioRxiv* 098913 (2017). doi:10.1101/098913
- 799 9. Matranga, C. B. *et al.* Enhanced methods for unbiased deep sequencing of Lassa and Ebola RNA viruses from  
800 clinical and biological samples. *Genome Biol.* **15**, 519 (2014).
- 801 10. Faria, N.R. *et al.* Epidemic establishment and cryptic transmission of Zika virus in Brazil and the Americas.  
802 Preprint at <https://doi.org/10.1101/105171> (2017).
- 803 11. Grubaugh, N.D. *et al.* Multiple introductions of Zika virus into the United States revealed through genomic  
804 epidemiology. Preprint at <https://doi.org/10.1101/104794> (2017).
- 805 12. Faria, N. R. *et al.* Zika virus in the Americas: Early epidemiological and genetic findings. *Science* **352**, 345–  
806 349 (2016).
- 807 13. Sall, A. A. *et al.* Yellow fever virus exhibits slower evolutionary dynamics than dengue virus. *J. Virol.* **84**,  
808 765–772 (2010).
- 809 14. First case of Zika virus reported in Puerto Rico. *Centers for Disease Control and Prevention* (2015).
- 810 15. Pan American Health Organization. *Zika: Epidemiological Report Honduras.* (World Health Organization,  
811 2016).
- 812 16. Pan American Health Organization. *Epidemiological Update: Zika virus infection.* (World Health Organization,  
813 Oct 16 2015).
- 814 17. Pan American Health Organization. *Zika: Epidemiological Report Dominican Republic.* (World Health  
815 Organization, 2016).
- 816 18. Nunes, M. R. T. *et al.* Emergence and potential for spread of Chikungunya virus in Brazil. *BMC Med.* **13**, 102  
817 (2015).
- 818 19. Tsetsarkin, K. A., Vanlandingham, D. L., McGee, C. E. & Higgs, S. A single mutation in chikungunya virus  
819 affects vector specificity and epidemic potential. *PLoS Pathog.* **3**, e201 (2007).

- 820 20. Piantadosi, A. *et al.* HIV-1 evolution in gag and env is highly correlated but exhibits different relationships  
821 with viral load and the immune response. *AIDS* **23**, 579–587 (2009).
- 822 21. Villabona-Arenas, C. J. *et al.* Dengue Virus Type 3 Adaptive Changes during Epidemics in São Jose de Rio  
823 Preto, Brazil, 2006–2007. *PLoS One* **8**, e63496 (2013).
- 824 22. Brinton, M. A. & Basu, M. Functions of the 3' and 5' genome RNA regions of members of the genus  
825 Flavivirus. *Virus Res.* **206**, 108–119 (2015).
- 826 23. Duchêne, S., Ho, S. Y. W. & Holmes, E. C. Declining transition/transversion ratios through time reveal  
827 limitations to the accuracy of nucleotide substitution models. *BMC Evol. Biol.* **15**, 36 (2015).
- 828 24. Corman, V. M. *et al.* Clinical comparison, standardization and optimization of Zika virus molecular detection.  
829 *Bull. World Health Organ.* (2016).
- 830 25. Gire, S. K. *et al.* Genomic surveillance elucidates Ebola virus origin and transmission during the 2014  
831 outbreak. *Science* **345**, 1369–1372 (2014).
- 832 26. Pyke, A. T. *et al.* Imported zika virus infection from the cook islands into australia, 2014. *PLoS Curr.* **6**,  
833 (2014).
- 834 27. Lanciotti, R. S. *et al.* Genetic and serologic properties of Zika virus associated with an epidemic, Yap State,  
835 Micronesia, 2007. *Emerg. Infect. Dis.* **14**, 1232–1239 (2008).
- 836 28. Faye, O. *et al.* One-step RT-PCR for detection of Zika virus. *J. Clin. Virol.* **43**, 96–101 (2008).
- 837 29. Faye, O. *et al.* Quantitative real-time PCR detection of Zika virus and evaluation with field-caught mosquitoes.  
838 *Virol. J.* **10**, 311 (2013).
- 839 30. Balm, M. N. D. *et al.* A diagnostic polymerase chain reaction assay for Zika virus. *J. Med. Virol.* **84**, 1501–  
840 1505 (2012).
- 841 31. Tappe, D. *et al.* First case of laboratory-confirmed Zika virus infection imported into Europe, November 2013.  
842 *Euro Surveill.* **19**, (2014).

843

## 844 **References cited in Methods and Extended Data**

- 845 32. *Zika Virus Response Updates from FDA.* (U.S. Food and Drug Administration, 2017).
- 846 33. Morlan, J. D., Qu, K. & Sinicropi, D. V. Selective depletion of rRNA enables whole transcriptome profiling of  
847 archival fixed tissue. *PLoS One* **7**, e42882 (2012).
- 848 34. Worobey, M. *et al.* 1970s and 'Patient 0' HIV-1 genomes illuminate early HIV/AIDS history in North  
849 America. *Nature* **539**, 98–101 (2016).
- 850 35. Clark, K., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J. & Sayers, E. W. GenBank. *Nucleic Acids Res.* **44**, D67–  
851 72 (2016).
- 852 36. Park, D. J. *et al.* Ebola Virus Epidemiology, Transmission, and Evolution during Seven Months in Sierra  
853 Leone. *Cell* **161**, 1516–1526 (2015).
- 854 37. Tomkins-Tinch, C. *et al.* *broadinstitute/viral-ngs: vl.13.3.* (2016). doi:10.5281/zenodo.200428
- 855 38. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using exact alignments.  
856 *Genome Biol.* **15**, R46 (2014).
- 857 39. O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and

- 858 functional annotation. *Nucleic Acids Res.* **44**, D733–45 (2016).
- 859 40. Aurrecochea, C. *et al.* PlasmoDB: a functional genomic database for malaria parasites. *Nucleic Acids Res.* **37**,  
860 D539–43 (2009).
- 861 41. Yarza, P. *et al.* The All-Species Living Tree project: a 16S rRNA-based phylogenetic tree of all sequenced type  
862 strains. *Syst. Appl. Microbiol.* **31**, 241–250 (2008).
- 863 42. Brister, J. R., Ako-Adjei, D., Bao, Y. & Blinkova, O. NCBI viral genomes resource. *Nucleic Acids Res.* **43**,  
864 D571–7 (2015).
- 865 43. NCBI Resource Coordinators. Database resources of the National Center for Biotechnology Information.  
866 *Nucleic Acids Res.* **44**, D7–19 (2016).
- 867 44. Grabherr, M. G. *et al.* Full-length transcriptome assembly from RNA-Seq data without a reference genome.  
868 *Nat. Biotechnol.* **29**, 644–652 (2011).
- 869 45. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation sequencing  
870 data. *Bioinformatics* **28**, 3150–3152 (2012).
- 871 46. Hyatt, D., LoCascio, P. F., Hauser, L. J. & Uberbacher, E. C. Gene and translation initiation site prediction in  
872 metagenomic sequences. *Bioinformatics* **28**, 2223–2230 (2012).
- 873 47. R Core Team. R: A Language and Environment for Statistical Computing. *R Foundation for Statistical*  
874 *Computing* (2016). Available at: <https://www.R-project.org/>.
- 875 48. Cribari-Neto, F. & Zeileis, A. Beta Regression in R. *J. Stat. Softw.* **34**, 1–24 (2010).
- 876 49. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements in  
877 performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
- 878 50. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
- 879 51. Kears, M. *et al.* Geneious Basic: an integrated and extendable desktop software platform for the organization  
880 and analysis of sequence data. *Bioinformatics* **28**, 1647–1649 (2012).
- 881 52. Henikoff, S. & Henikoff, J. G. Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. U.*  
882 *S. A.* **89**, 10915–10919 (1992).
- 883 53. Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the  
884 performance of PhyML 3.0. *Syst. Biol.* **59**, 307–321 (2010).
- 885 54. Rambaut, A. FigTree. Version 1.4.2. *Edinburgh, UK: Inst. Evol. Biol., Univ. Edinburgh.* (2014). Available at:  
886 <http://tree.bio.ed.ac.uk/software/figtree/>.
- 887 55. Rambaut, A., Lam, T. T., Max Carvalho, L. & Pybus, O. G. Exploring the temporal structure of heterochronous  
888 sequences using TempEst (formerly Path-O-Gen). *Virus Evol.* **2**, vew007 (2016).
- 889 56. Drummond, A. J., Suchard, M. A., Xie, D. & Rambaut, A. Bayesian phylogenetics with BEAUti and the  
890 BEAST 1.7. *Mol. Biol. Evol.* **29**, 1969–1973 (2012).
- 891 57. Shapiro, B. *et al.* A Bayesian phylogenetic method to estimate unknown sequence ages. *Mol. Biol. Evol.* **28**,  
892 879–887 (2011).
- 893 58. Shapiro, B., Rambaut, A. & Drummond, A. J. Choosing appropriate substitution models for the phylogenetic  
894 analysis of protein-coding sequences. *Mol. Biol. Evol.* **23**, 7–9 (2006).
- 895 59. Drummond, A. J., Rambaut, A., Shapiro, B. & Pybus, O. G. Bayesian coalescent inference of past population

- 896 dynamics from molecular sequences. *Mol. Biol. Evol.* **22**, 1185–1192 (2005).
- 897 60. Drummond, A. J., Ho, S. Y. W., Phillips, M. J. & Rambaut, A. Relaxed phylogenetics and dating with  
898 confidence. *PLoS Biol.* **4**, e88 (2006).
- 899 61. Ferreira, M. A. R. & Suchard, M. A. Bayesian analysis of elapsed times in continuous-time Markov chains.  
900 *Can. J. Stat.* **36**, 355–368 (2008).
- 901 62. Baele, G. *et al.* Improving the accuracy of demographic and molecular clock model comparison while  
902 accommodating phylogenetic uncertainty. *Mol. Biol. Evol.* **29**, 2157–2167 (2012).
- 903 63. Baele, G., Li, W. L. S., Drummond, A. J., Suchard, M. A. & Lemey, P. Accurate model selection of relaxed  
904 molecular clocks in bayesian phylogenetics. *Mol. Biol. Evol.* **30**, 239–243 (2013).
- 905 64. Yang, Z. Maximum-Likelihood Models for Combined Analyses of Multiple Sequence Data. *J. Mol. Evol.* **42**,  
906 587–596 (1996).
- 907 65. Lê, S., Josse, J. & Husson, F. FactoMineR: an R package for multivariate analysis. *J. Stat. Softw.* (2008).
- 908 66. Josse, J. & Husson, F. missMDA: A Package for Handling Missing Values in Multivariate Data Analysis. *J.*  
909 *Stat. Softw.* **70**, 1–31 (2016).
- 910 67. Ann-Claire Gourinat, Olivia O’Connor, Elodie Calvez, Cyrille Goarant & Myrielle Dupont-Rouzeyrol.  
911 Detection of Zika Virus in Urine. *Emerging Infectious Disease journal* **21**, 84 (2015).
- 912 68. Paz-Bailey, G. *et al.* Persistence of Zika Virus in Body Fluids — Preliminary Report. *N. Engl. J. Med.* Epub  
913 ahead of print at <https://doi.org/10.1056/NEJMoa1613108> (2017).
- 914







